SOMETIMES I AM A TREE: DATA DRIVES UNSTABLE HIERARCHICAL GENERALIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Neural networks often favor shortcut heuristics based on surface-level patterns. Language models (LMs), for example, behave like n-gram models early in training. However, to correctly apply grammatical rules, LMs must instead rely on hierarchical syntactic representations rather than on surface-level heuristics derived from n-grams. In this work, we use cases studies of English grammar to explore how latent structures in training data drives models toward improved out-of-distribution (OOD) generalization. We then investigate how data composition can lead to inconsistent behavior across random seeds. Our results show that models stabilize in their OOD behavior only when they commit to either a surface-level linear rule or a hierarchical rule. The hierarchical rule, furthermore, is induced by grammatically complex sequences with deep embedding structures, whereas the linear rule is induced by simpler sequences. When the data contains a mix of simple and complex examples, potential rules compete; each independent training run either stabilizes by committing to a single rule or remains unstable in its OOD behavior. We also identify an exception to the relationship between stability and generalization: Models which memorize patterns from homogeneous training data can overfit stably, with different rules for memorized and unmemorized patterns. While existing works have attributed similar generalization behavior to training objective and model architecture, our findings emphasize the critical role of training data in shaping generalization patterns and how competition between data subsets contributes to inconsistent generalization outcomes.

031 032

033

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

1 INTRODUCTION

034 Neural networks often learn shortcut heuristics which reflect simple, surface-level patterns in data. In the case of language models (LMs) trained on next-token prediction objectives, this simplicity bias can lead models to behave like n-gram models, relying heavily on local dependencies without 037 fully capturing the deeper, more complex structures of language (Choshen et al., 2022; Geirhos et al., 2020; Saphra & Lopez, 2018). However, LMs are also capable of breakthroughs in generalization, shifting from these simple heuristics to more sophisticated behaviors (Choshen et al., 2022; Chen et al., 2023; McCoy et al., 2020). Such transitions suggest that under certain training conditions, LMs 040 can eventually overcome spurious shortcuts and use linguistic structures to generalize beyond surface-041 level patterns. Previous works often attribute this ability to model architecture and training objectives 042 (Ahuja et al., 2024; McCoy et al., 2020). In this work, we investigate how data characteristics 043 influence the generalization rules learned, especially when multiple solutions fit the training data 044 equally well. We also examine the instabilities associated with generalization behaviors. 045

To understand when and why a model favors learning latent structures over surface-level heuristics, we use case studies in learning English grammar rules. Grammatically correct sentences in English must follow a set of rules that operate on a sequence's latent tree-like structure (Chomsky, 2015; Crain & Nakayama, 1987). When trained on next-token prediction, an LM may approximate these rules from surface-level statistics, acting as an n-gram model by applying a *linear rule*. However, such a model struggles to generalize to unseen grammatical patterns. Figure 1 (*bottom right*) shows that a LM can use a linear bigram model to capture the relationship between a subject noun and its verb by *inflecting* the verb with the same plurality as the subject. This LM would fail to generalize when a *distractor* noun, e.g., from a prepositional phrase, appears between subject and verb. In contrast,



Figure 1: Data plays a critical role in generalization behaviors and training stability. *Left:* Along 071 the data diversity x-axis, Low data diversity (as measured by variation in syntactic structure) leads 072 the model to memorize unreliable sample-specific patterns, whereas high data diversity promotes 073 commitment to a general rule. Along the data complexity y-axis, high data complexity (i.e., with 074 more center embeddings) induces the hierarchical rule, while simpler data (right-branching sentences) 075 induces the surface-level linear rule. Mixing these data types results in unstable OOD training 076 behaviors. Upper Right: A model that captures hierarchical structure can generalize grammatical 077 rules OOD by correctly identifying the subject as the noun closest to the root on the syntax tree graph. *Lower Right:* A model that uses the linear rule will treat the most recent noun as the target verb's subject and thereby fail to generalize to unseen sentence compositions. 079

081

083

Figure 1 (*upper right*) shows a model that instead uses a latent tree structure can learn the correct syntactic rule (i.e., the *hierarchical rule*), enabling it to generalize to novel sentence compositions.

Building on previous work (McCoy et al., 2018; 2020; Ahuja et al., 2024), we use two tasks question formation and tense inflection—to investigate whether the model learns the hierarchical rule or defaults to the surface-level linear rule. We train models on ambiguous data, which is compatible with both rules, and evaluate them on out-of-distribution (OOD) data which is compatible only with the hierarchical rule. We first find that a preference for OOD hierarchical generalization is induced by training on samples with center embeddings, where the subject is modified by an relative clause. This result mirrors a celebrated claim from linguistics (Wexler, 1980) that center embeddings are responsible for human syntax acquisition.

Models trained on the same data exhibit inconsistent OOD behaviors across random seeds. By examining training dynamics, we identify a connection between training stability and rule commitment: Only runs that commit to a rule can exhibit stable OOD performance during training. Connecting back to data, we show that training dynamics can be categorized into three distinct regimes that depend on data complexity and diversity, illustrated in Figure 1 *left*. Data diversity determines whether a model learns a general rule, while data complexity determines which rule is preferred. Models trained on a mix of hierarchical-inducing samples and linear-inducing samples are most unstable during training and exhibit the largest inconsistency across random seeds.

100

103

104

105

106

107

092

Taken together, our findings demonstrate that data composition plays a critical role in shaping model's
 OOD generalization behaviors. Our contributions are as follows:

- We show that sentences with complex grammatical structure—specifically center embeddings drive LMs to favor hierarchical syntactic representations over surface-level n-gram heuristics, enabling correct OOD generalization of grammatical rules (see Section 4).
- We demonstrate that models stabilize in OOD performance only when they commit to either a surface-level heuristic or a hierarchical rule (see Section 5).

• We show that when the training data mixes complex and simple grammatical structures, the resulting rules are inconsistent across random seeds and many models fail to stabilize OOD behavior by the end of training (see Section 5).

• We identify an exception to the relationship between stability and rule learning: Models trained on insufficiently diverse data stabilize in a memorization regime without learning either rule, highlighting another way that data can drive generalization failures (see Section 6).

2 RELATED WORK

We include the most relevant work in this section. For an extended discussion of related work, please refer to Appendix A.

108

109

110

111

112

113

114 115

116 117

2.1 SYNTAX AND HIERARCHICAL GENERALIZATION

122 McCoy et al. (2018) first used the question formation task to study hierarchical generalization in 123 neural networks, showing that RNNs trained with a seq-to-seq objective exhibit limited hierarchical generalization. However, adding attention mechanisms improved performance on the generalization 124 set. Later, McCoy et al. (2020) found that tree-structured architectures consistently induce hierarchical 125 generalization. Petty & Frank (2021) and Mueller et al. (2022) further investigated inductive biases 126 and concluded that, like RNNs, transformers tend to generalize linearly. This view was challenged by 127 Murty et al. (2023), who attributed the failure of prior attempts to insufficient training, demonstrating 128 that decoder-only transformers can generalize hierarchically, but only after in-distribution perfor-129 mance has plateaued. Expanding on this, Ahuja et al. (2024) showed that hierarchical generalization 130 is achieved only with models trained on a language modeling objective. Previous work primarily 131 attributed the source of inductive bias toward hierarchical rule to model architecture, whereas our 132 study highlights the impact of data, and we further provide a precise measure of data complexity and 133 data diversity. In addition, the inconsistency in generalization behaviors are observed in (McCoy 134 et al., 2018; McCoy et al.). While they have pointed out that models trained on different random seeds can manifest very different generalization behaviors, they did not further studies the distributions of 135 model behaviors. 136

137

Similar to our work, Papadimitriou & Jurafsky (2023) and Papadimitriou & Jurafsky (2020) also
 studied how training data could introduce an inductive bias to affect language acquisition. specifically
 identified that by pretraining models on data with a recursive structure, finetuning them on natural
 language yields superior performances. This finding is closely related to our conclusions around
 center embeddings since the center embedding structure in language is recursive in nature.

143 144

2.2 TRAINING DYNAMICS AND GROKKING

145 Grokking refers to the phenomenon where a neural network, after achieving seemingly poor per-146 formance for a long period, suddenly generalize on unseen data. Power et al. (2022) first observed 147 this behavior in simple arithmetic tasks. Since then, the exact mechanism of grokking has been widely studied. Different from existing grokking work, we studied a different types of grokking: 148 "structural grokking" (Murty et al., 2023). In classic grokking, the model transitions from memoriza-149 tion to generalization, allowing it to achieve non-trivial performance on unseen data. In structural 150 grokking, a model transitions from the simple linear rule to the hierarchical rule, leading to non-trivial 151 performance on OOD data. However, the findings of this study is potentially related to those in 152 classic grokking. Zhu et al. (2024) studies the role of data and finds that grokking only occurs when 153 training data is sufficiently large. Berlot-Attwell et al. (2023) broadly studies how data diversity and 154 complexity leads to different generalization behaviors. Liu et al. (2022) shows that grokking can be 155 induced with different weight norms, associating generalization with a specific goldilocks zone weight 156 norm value. Huang et al. (2024) and Varma et al. (2023) have shown that during training, different 157 circuits compete, and models trained on different random seeds can lead to distinct generalization 158 behaviors depending on which circuits dominate. While these works primarily attribute generalization 159 differences to circuit formation, our findings highlight that this competition between circuits can also destabilize learning dynamics. Importantly, we characterize the unstable regime in both data diversity 160 and data complexity and we address connections between training stability and consistency under 161 random variation.

162 Table 1: Examples from two grammar case studies. *Top*: In the question formation task, the model 163 moves the main auxiliary verb to the front to form a question. *Bottom*: In the tense inflection task, the 164 model inflects the main verb from past to present tense, while respecting subject-verb agreement.

Dataset	Task Type	Examples
Quastian Earmatian	Quest	Input: My unicorn does move the dogs that do wait.
Question Formation	(Ambiguous)	Output: Does my unicorn move the dogs that do wait?
		Input: My unicorn who doesn't sing does move.
	Quest	Linear Output: Doesn't my unicorn who sing does move?
	(Unambiguous)	Hierarchical Output: Does my unicorn who doesn't sing move?
	Present	Input: My zebra behind the peacock smiled.
	(Ambiguous)	Output: My zebra behind the peacock smiles.
Tense Inflection	Present	Input: My zebra behind the peacocks smiled.
		Linear output: My zebra behind the peacocks smile.
	(Unambiguous)	Hierarchical output: My zebra behind the peacocks smiles.

17 176 177

178

187

188

2.3 RANDOM VARIATION

Although choices like hyperparameter settings, architecture, and optimizer all shape model outcomes, 179 training remains inherently stochastic. Models are sensitive to random initialization and the order of 180 training examples. Several studies (Zhou et al., 2020; D'Amour et al., 2022; Naik et al., 2018) have 181 reported significant performance differences across model checkpoints and Zhou et al. (2020) noted 182 that instability extends throughout the training curve. Dodge et al. (2020) found that both weight 183 initialization and data order contribute equally to out-of-sample performance variation. Unlike prior work, which focuses on the experimental implications of random variations, we investigate the source 185 of these training inconsistencies and link them to characteristics of the training data.

3 EXPERIMENTAL SETUP

The question formation task and the tense inflection task are first proposed by Frank & Mathis (2007) 189 and Linzen et al. (2016) as canonical tasks to assess a model's language modeling ability. In this 190 study, we use the synthetic dataset constructed by McCoy et al. (2018) (for question formation) and 191 McCoy et al. (2020) (for tense inflection). 192

193 3.1 QUESTION FORMATION TASK

194 In the question formation (QF) task, a declarative sentence is transformed into a question (see 195 Table 1) by moving the main auxiliary verb (such as "does" in "does move") to the front. Our training 196 data permits two strategies for choosing which verb to move: (1) move first: a linear rule that moves 197 the first auxiliary, or (2) move main: a hierarchical rule—the correct rule in English grammar—based on the sentence's syntactic structure. This syntactic structure links each word into a tree-like structure 199 in which edges specify syntactic dependencies (e.g., subject, preposition, object), as shown in Figure 2. The model leverages this tree representation to determine which auxiliary to move. 200

201

In Table 1, the first example is considered **ambiguous** because both the hierarchical and linear rules 202 produce the correct outcome. In contrast, the second example is **unambiguous** because only the 203 hierarchical rule produces the correct outcome. The training data contains only ambiguous samples, 204 while the OOD generalization set includes only unambiguous samples. If a model uses a hierarchical 205 representation of syntax, it should achieve 100% accuracy on both the in-distribution (ambiguous 206 questions) and OOD generalization (unambiguous questions) sets. Conversely, if a model rise on 207 linear rules, it will score 0% on the OOD generalization set, but still score 100% accuracy on the 208 in-distribution set. We therefore use the model's accuracy on the OOD generalization set as a metric 209 for hierarchical generalization.

210 3.2 TENSE INFLECTION TASK 211

212 In the **tense inflection** (**TI**) task, we provide the model with a sentence in the past tense, and the 213 model transforms it into the present tense. Since past-tense verbs in English do not differentiate between singular and plural forms, the model must identify the subject to determine whether the 214 present-tense verb should be inflected as singular or plural. The TI task tests whether the model 215 follows the hierarchical or linear rule for subject-verb agreement. The linear rule inflects the verb



Figure 2: Sentence Examples. *Left:* Right-branching sentence examples. The linear progression of the main phrase is not interrupted by the relative clause. *Right:* Center-embedded sentence examples. When the relative clause modifies the subject, it interrupts the linear progression of the main clause.

based on the most recent noun, while the hierarchical rule correctly inflects the verb according to
the subject. Like in the QF task, the training data contains ambiguous samples (example in Table 1),
where the subject noun (i.e., "*zebra*") and the most recent noun (i.e., "*peacock*") always share the
same plurality and therefore either rule produces the correct answer. The OOD generalization set
includes unambiguous examples, where the subject and the most recent noun differ in plurality and
therefore only the hierarchical rule produces the correct answer. Similar to the QF task, we use the
model's main-verb prediction accuracy on the OOD set as a metric for hierarchical generalization.

236 3.3 MODEL, DATA AND TRAINING

237 We use a decoder-only Transformer architecture with 12M parameters: 6 layers of 8 heads with 238 a 512-dimensional embedding for QF. For TI, we use the same transformer architecture but with 239 4 layers. All models are trained from scratch on a causal language modeling objective for 300K steps. We use the Adam optimizer (Kingma & Ba, 2014), a learning rate of 1e-4, and a linear decay 240 schedule. All the hyperparamter settings are directly adopted from existing works (Ahuja et al., 2024; 241 Murty et al., 2023). We run all experiments on the same 50 random seeds. We use the original 242 training, validation and OOD generalization data proposed by McCoy et al. (2018) and McCoy et al. 243 (2020). To create variations on the training data, we mimic the data generation process used for the 244 original QF and TI task. Specifically, the original TI and QF data are generated with Context-Free 245 Grammar (CFG) using a simplified set of grammatical rules, and we reuse the same CFG rules to 246 create variations of the training data. We use a word-level tokenizer with a vocabulary of size 72. 247

248

224

225

226

227 228

4 DATA COMPLEXITY DETERMINES RULE PREFERENCE

We begin by analyzing center-embedded sentences in Section 4.1. We then show that center-embedded sentences drives hierarchical generalization in both QF task (Section 4.2) and TI task (Section 4.3).

2 4.1 CENTER EMBEDDING

253 Center embedding occurs when a clause—often acting as a modifier—is placed within another clause 254 or phrase. Figure 2 (*left*) illustrates two examples of center-embedded sentences, where the embedded clause disrupts syntactic dependencies, such as the subject-verb-object relationship. Moreover, center 256 embeddings exhibit a recursive structure: inside the relative clause, one can find the same structure as the entire sentence. In contrast, sentences without center embeddings are exclusively right branching. 257 Right-branching structures may also include modifying clauses, but these can only be appended at 258 the end of the main clause, maintaining its linear flow (see Figure 2, *right*). Center embedding has 259 been central to linguistic studies on the types of data required to learn grammatical rules. According 260 to Chomsky's generative grammar framework (Chomsky, 2015), center-embedded clauses give rise 261 to hierarchical, tree-like syntactic structures. Additionally, Wexler (1980) posits that all English 262 syntactic rules can be learned from "degree 2" sentences, which contain exactly one embedded clause. 263

While center embeddings are crucial for human language acquisition, in this study we investigate
whether the same type of data can lead a LM to acquire the hierarchical grammar rule. To correctly
predict the distribution of next tokens, LMs must track dependencies between sentence components.
In right-branching sentences, LMs can rely on linear proximity to identify dependencies; for example,
as shown in Figure 2, a simple bigram model suffices to capture the subject-verb relationship. In
contrast, center embeddings introduce relative clauses of various lengths, making linear n-gram
models inefficient for capturing subject-verb dependencies. Furthermore, the recursive nature means



Figure 3: **Components of training data drive different generalization behaviors.** *Left:* Centerembedded sentences, which in the QF training data only appear in declaration copying examples, induce hierarchical generalization. *Right:* Models are trained on different data mixes and evaluated on two OOD sets: unambiguous right-branching sentences (*green*) and unambiguous center-embedded sentences (*red*). For center-embedded sentences, the hierarchical rule is preferred regardless of data mixes. For right-branching sentences, the model's preference for the hierarchical rule is exclusively driven by having a large mix of center-embedded sentences in the training data.

that the model needs to keep track of multiple subject-verb dependencies: one for the main clause
 and a separate one for the embedded relative clause. In these cases, modeling those subject-verb
 relationships with a tree structure is more compact and efficient.

4.2 QUESTION FORMATION

293 As specified in Section 3.1, the training data for QF must be ambiguous between the linear rule (i.e., moving the first auxiliary) and the hierarchical rule (i.e., moving the main auxiliary). Center-295 embedded sentences do not meet this ambiguity requirement and, therefore, cannot appear in question 296 formation training samples. To ensure the model is exposed to diverse sentence types, McCoy 297 et al. (2018) introduces a secondary task to the QF training dataset: declaration copying. Like 298 question formation, the declaration-copying sample starts with a declarative sentence, but instead of 299 transforming it, the model simply repeats it. Since the ambiguity requirement does not apply to the declaration-copying task, center-embedded sentences are included in this secondary task. Concrete 300 examples of both tasks can be found in Appendix **B**. 301

We train models on three subsets of the original training data, varying the composition of the 303 declaration-copying examples. In *Ouest Only*, we remove all declaration copying examples. In *Center* 304 embed, we only keep center-embedded examples. In Right branch, we only keep right-branching 305 examples. For all three subsets, the question formation samples remain unchanged. Each setup 306 reaches 100% in-distribution validation accuracy; however, the OOD generalization performance, 307 shown in Figure 3 (left), differs significantly across the subsets. When the declaration-copying task is 308 removed, none of the 50 runs achieve an OOD generalization accuracy above 75%, indicating that 309 declaration copying examples are essential for inducing the hierarchical rule. When trained solely on center-embedded sentences in the declaration-copying task, models exhibit a strong preference for 310 the hierarchical rule. In contrast, training only on right-branching sentences leads to poor hierarchical 311 generalization. This evidence suggests that center-embedded sentences direct a model towards the 312 hierarchical rule. 313

4.3 TENSE INFLECTION

We now analyze hierarchical generalization in the tense inflection task, demonstrating the generality of our findings across grammatical rules. Linzen et al. (2016) first proposed the idea to use a verb inflections to assess the model's grammatical capabilities. McCoy et al. (2020) then adopted the question formation data to the tense inflection task by creating the TI dataset using a set of CFG rules and vocabularies similar to the ones used for QF.

320

281

282

283

284

285

286

287 288

302

In the TI training data, both right-branching and center-embedded sentences are made ambiguous
 by ensuring the distractor noun shares the same plurality as the main subject. For right-branching
 sentences, since there isn't a relative clause modifying the subject, a preposition phrase modifying
 the subject provides the distractor noun. In contrast, for center-embedded sentences, since there is a

324 relative clause modifying the subject, either the subject or the object of the modifying clause can act 325 as the distractor noun. We list examples below: 326

- 1. **Right Branching**: The noun in the prepositional phrase (e.g., "to the cabinet") acts as the distractor in the TI task. Example A (ID): The keys to the cabinet are on the table.
- 328 330

334

327

- Example B (OOD): The keys to the **cabinets** are on the table. 2. Center Embedding: Either the subject or the object inside the relative clause acts as the 331 332 distractor in the TI task. 333
 - Example C (ID): The keys that unlock the cabinets are on the table.

Example D (OOD): The keys that unlock the cabinet are on the table.

335 We create variations of the TI training data by adjusting the ratio of right-branching to center-336 embedded sentences while keeping the total training size constant. The model is trained on nine 337 different data mixes, and its generalization behavior is tested on two OOD sets: one containing 338 unambiguous right-branching sentences (e.g., Example B) and the other containing unambiguous 339 center-embedded sentences (e.g., Example D). 340

341 Generalization accuracies are shown in Figure 3 (right). When the training data is dominated by 342 ambiguous right-branching sentences, the model fails to learn the hierarchical rule, as indicated 343 by low OOD generalization accuracy on the left side of the green line. However, increasing the proportion of center-embedded sentences biases the model toward applying the hierarchical rule, 344 even on right-branching sentences. This shift in behavior is reflected by improved generalization 345 accuracy on the right side of the green line. The red line in Figure 3 (right) represents the model's 346 generalization accuracy on unambiguous center-embedded sentences. Regardless of the data mix, the 347 model consistently treats center embeddings as hierarchical and applies the hierarchical rule to OOD 348 data. In contrast, the model only applies the hierarchical rule to right-branching sentences after being 349 exposed to a sufficient quantity of center-embedded sentences during training. These observations 350 suggest that center embeddings drive the model's overall preference for tree structures. 351

352 In Appendix C, we further partition center-embedded sentences based on the syntactic role of the 353 main subject within the modifying clause. We show that while both subtypes induce the hierarchical rule in the QF task, one subtype provides a stronger hierarchical bias in the TI task. 354

355 356

5 TRAINING STABILIZES IF A MODEL COMMITS TO A RULE

357 Why do some runs fail to generalize hierarchically even when trained on hierarchical-inducing data? 358 In this section, we will show that these failures are consequences of training instability; models only 359 stabilize OOD if they commit to a general rule. 360

361

5.1 INSTABILITY DURING TRAINING

362 When training models on both QF and TI, some random seeds lead to highly unstable OOD behavior, 363 with generalization accuracy often undergoing large swings during training. Furthermore, the unstable 364 behavior is not consistent across different seeds. In Appendix F, we show examples of different OOD behaviors during training. We also show that both the instability and inconsistency in OOD behavior 366 are significant only after the ID performance has converged. We measure instability across training 367 time using total variation (TV). Specifically, we checkpoint the model every 2K steps and measure 368 the generalization accuracy at each checkpoint, denoting as Acc_i . The total variation is defined as:

Total Variation (TV) =
$$\frac{1}{|\text{ckpts}|} \sum_{i \in \text{ckpts}} |\text{Acc}_i - \text{Acc}_{i-1}|$$
, where $\text{ckpts} = \{2\text{K}, 4\text{K}, 6\text{K}, \dots\}$

371 372 373

369 370

5.2 TRAINING STABILITY TIES TO RULE COMMITMENT

374 We now demonstrate the connection between stable OOD behavior and rule commitment. We con-375 struct QF training datasets such that they contain different proportions of hierarchical-inducing (i.e., 376 center-embedded) and linear-inducing (i.e., right-branching) declarations, while keeping questions 377 constant. Further details on the dataset can be found in Appendix D.



Figure 5: **Total variation across training v.s. final generalization accuracy for QF task.** OOD behavior stabilizes during training if a model commits to a simple rule. By mixing data that induces the linear and hierarchical rules, we can create conditions that allow models to stabilize in either rule. "Linear" denotes the proportion of linear-inducing declarations in the data. Grey line indicates the smoothed average curve across all runs and all five datasets.

394 Figure 4 shows the relationship between data homogeneity 395 and training stability. When the training data is domi-396 nated by either linear-inducing (linear=99%) or hierarchy-397 inducing (linear=0%) examples, more random seeds lead to stable OOD curves. When the training data is a hetero-398 geneous mix instead, potential rules compete, leading to 399 unstable training. Figure 5 shows the relationship between 400 training stability and generalization performance. 401



Across the five data mixes, the final generalization accuracy for all the stable models is either 100% or 0%, indicating that the stable models have all committed to a general rule. While models can stabilize in either rule, data composition determines how likely a run is stable and for stable runes which rule is favored. Interestingly, when



the data heterogeneous (e.g., linear=10% case), the final generalization accuracy for stable runs is
bimodally distributed, clustering around 100% or 0%. This bimodality suggests that training stability
is always associated with a commitment to either the linear or the hierarchical rule, even when the
data mix does not favor either rule.

In summary, with heterogeneous training data, competition between rules leads to more unstable training runs. Even with heterogeneous data mixes, however, some runs can still stabilize if they commit to one of the competing rules. In Appendix E.2, we replicate this analysis for the TI task, showing similar results.

416 417 418

387

388

389

390

391

392 393

6 DATA DIVERSITY LEADS TO GENERALIZATION

We have linked training stability to rule commitment. But why can't networks stabilize without committing to a rule? In this section, we will explore the non-monotonic relationship between data diversity, training instability, and rule commitment.

423 6.1 MEASURING DATA DIVERSITY

424 In order to measure the diversity of our training data, we must compute the syntactic similarity 425 between different example sentences. We describe a sentence pair's similarity by the tree-edit 426 distance (TED) of their latent tree representations (Chomsky, 2015). When two sentences share the 427 same syntax tree, transforming one into the other requires only leaf-node (i.e., vocabulary) changes. 428 For example, "My unicorn entertains her tyrannosaurus," and, "Your zebra eats some apples," have 429 different vocabulary but identical syntax trees. We define data diversity as the number of unique syntactic trees in the training data. This way of using syntax TED to measure diversity data has 430 been used in both natural language (Huang et al., 2023; Gao & He, 2024; Ramírez et al., 2022) and 431 code (Song et al., 2024). We will show that when the model is exposed to a fewer unique syntax



Figure 6: **Inverse U-shaped relationship between training stability and data diversity.** At low data diversity, training is stable but the model memorize individual syntactic patterns rather than committing to a rule. With moderate data diversity, training becomes unstable. As diversity increases further, the model commits to a rule and training stabilizes again.

trees during training, it memorizes those patterns without extrapolating any rules to unseen sentencestructures. Consequently, the model fails to commit to a general rule.

450 451 6.2 INVERSE U-SHAPED SCALING

Commitment to hierarchical rule We first control data diversity on datasets that induce hierarchical generalization in QF. We construct variations of the QF training data, each with 50K question samples and 50K hierarchical-inducing declarations, while varying the diversity of the declaration examples. We train 20 random seeds for each training set variation. To measure intra-run instability, we use total variation, and to assess hierarchical rule commitment, we report the proportion of runs achieving generalization accuracy >95%. Figure 6 (left) shows the distribution of total variation across 20 seeds and the corresponding hierarchical generalization ratios.

459 We observe an inverse U-shaped relationship between data diversity and training instability, revealing 460 three distinct regimes. In the low-diversity regime, training is stable but the model fails to commit to a rule. In Appendix G, we further investigate this failure to rule commitment. We show that in trained 461 on low diversity data, model can memorize specific syntax patterns and apply the hierarchical rule to 462 to those structures but it cannot extrapolate the rule to unseen syntax structures. In the mid-diversity 463 regime, training becomes unstable due to variation across batches. Overall, with insufficient diversity, 464 relatively few runs can learn the hierarchical rule. Finally, in the high-diversity regime, training 465 stabilizes again as the model commits to the hierarchical rule, indicated by a high hierarchical 466 generalization ratio. 467

Commitment to linear rule In Figure 5 right-most panel, the model has a strong preference to 468 apply linear rule OOD when the training data contains 99% linear-inducing data (i.e., right branching 469 sentences). However, Figure 3 (red violin) shows that when the training data contains *exclusively* 470 linear-inducing sentences, models suddenly fail to apply the linear rule OOD either. We can use data 471 diversity to explain the failure to rule commitment: right-branching sentences lack syntactic variation, 472 as the main auxiliary always follows the subject noun. This lack of syntax diversity prevents rule 473 extrapolation. By introducing as little as 1% of center-embedded sentences, we introduce the diversity 474 necessary to consistent apply a rule OOD and the skewed ratio between the hierarchical and linear 475 inducing sentences determine that the linear rule is preferred over the hierarchical rule.

476

443

444

445

446 447

To confirm that data diversity is also key to learning the linear rule, we create variations of QF
training data with 50K questions and 50K declarations, including 99% right-branching and 1%
center-embedded sentences. We control the diversity of *center-embedded* sentences as before and
use the proportion of runs achieving generalization accuracy below 5% to quantify the likelihood of
committing to the linear rule. As shown in Figure 6 (*right*), we observe a similar U-shaped scaling
behavior, confirming that models only commit to a rule when trained on diverse data.

482 483 484

7 DISCUSSION AND CONCLUSIONS

By exploring the role of data structure in determining OOD generalization rules, we have also revealed which settings allow us to predict model behavior. We show that complex grammatical structures

guide models toward hierarchical rules, while mixed data compositions lead to unstable dynamics
 and inconsistent rule commitment. These findings emphasize the importance of understanding how
 data diversity shapes both stability and generalization in neural networks. Our findings have a number
 of implications across machine learning and even formal linguistics.

490 **Clusters of generalization behavior across seeds** While errors are often treated as Gaussian noise 491 in the theoretical literature, our findings suggest that errors may only be distributed unimodally for a 492 given compositional solution. Our work joins the growing literature that suggests random variation not 493 only has an effect, but can create clusters of OOD behaviors. Previously, clustered distributions have 494 been documented in text classification heuristics (Juneja et al., 2022) and training dynamics (Hu et al., 495 2023). In our case, we note that generalization accuracy is only clearly multimodally distributed when 496 specifically considering stable training runs. We suggest that research on compositional variation in 497 training consider training stability in the future.

498 **Implications for formal linguistics** Our findings have potential implications for linguistics debates 499 about the poverty of the stimulus (McCoy et al., 2018; Berwick et al., 2011). Linguists have 500 extensively studied the question of what data is necessary and sufficient to learn grammatical rules. 501 In particular, Wexler (1980) argue that all English syntactic rules are learnable given "degree 2" data: 502 sentences with only one embedded clause nested within another clause. Our mixed scoping results 503 show that without a stronger architectural inductive bias-the very subject of the poverty of the 504 stimulus debate—degree 1 data alone cannot induce a preference for hierarchical structure. However, our work also supports the position of Lightfoot (1989) that lower degree data is adequate for a child 505 to learn a specific rule, as the LM generalizes ID degree 1 QF rule examples to OOD degree 2 by 506 using the hierarchical inductive bias induced by declaration examples. 507

Grokking, instability, and latent structure Murty et al. (2023), exploring the same data setting
 we do, call the transition from linear generalization to hierarchical generalization rules during training
 structural grokking. Classic grokking (Power et al., 2022), however, is different: Rather than a
 transition between generalization rules, it describes a transition from memorization to generalization.

512 Our findings clarify both scenarios. We link structural grokking to the instability formed by com-513 petition between linear- and hierarchical-inducing training subsets. Without competing subsets, the 514 model immediately learns either the linear or the hierarchical rule without the gradual transition of 515 structural grokking. This instability could represent the same phenomenon of circuit competition 516 described by Ahuja et al. (2024). We find a similar pattern of instability in our study of data diversity, 517 with implications for classic grokking. In this case, the competition is not between two rules, but 518 instead between memorized heuristics—sufficient for modeling syntactically homogeneous training 519 data—and simple OOD rules—required to efficiently model diverse training data. Yet again, while a 520 strict memorization regime is relatively stable, the regime between memorization and generalization is unstable, leading to potential grokking. 521

Our findings suggest that memorization is just another rule that the model can adopt when it is
the simplest way of capturing the training distribution. Such a framework unifies the grokking
literature with other phenomena such as emergence (Schaeffer et al., 2023) and benign interpolation
(Theunissen et al., 2020).

526 527

527 ETHICS STATEMENT 528

This research does not present any direct ethical concerns. The work involves empirical studies of
machine learning models and their behavior in language tasks. No human subjects, sensitive data, or
high-stakes applications were involved in this research. Therefore, no specific ethical considerations
were necessary for this work.

533 534 REPRODUCIBILITY STATEMENT

All relevant details regarding the experimental setup including model architecture, hyperparameters, and data preprocessing, are included in the main text (Section 3.3) and appendices (Section D). Additionally, the code and scripts used to run the experiments are provided in the supplementary material and will be made publicly available upon acceptance.

539

540 REFERENCES

578

579

580

587

588

- Kabir Ahuja, Vidhisha Balachandran, Madhur Panwar, Tianxing He, Noah A Smith, Navin Goyal, and
 Yulia Tsvetkov. Learning syntax without planting trees: Understanding when and why transformers
 generalize hierarchically. *arXiv* [*cs.CL*], April 2024.
- Boaz Barak, Benjamin L Edelman, Surbhi Goel, Sham Kakade, Eran Malach, and Cyril Zhang.
 Hidden progress in deep learning: SGD learns parities near the computational limit. *arXiv [cs.LG]*, July 2022.
- Ian Berlot-Attwell, Kumar Krishna Agrawal, A Michael Carrell, Yash Sharma, and Naomi Saphra.
 Attribute diversity determines the systematicity gap in VQA. *arXiv [cs.LG]*, November 2023.
- Robert C Berwick, Paul Pietroski, Beracah Yankama, and Noam Chomsky. Poverty of the stimulus revisited. *Cogn. Sci.*, 35(7):1207–1242, September 2011.
- Angelica Chen, Ravid Shwartz-Ziv, Kyunghyun Cho, Matthew L Leavitt, and Naomi Saphra. Sudden drops in the loss: Syntax acquisition, phase transitions, and simplicity bias in MLMs. *arXiv* [*cs.CL*], September 2023.
- Noam Chomsky. *Aspects of the theory of syntax*. The MIT Press. MIT Press, London, England, 50 edition, 2015.
- Leshem Choshen, Guy Hacohen, Daphna Weinshall, and Omri Abend. The grammar-learning trajectories of neural language models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8281–8297, Stroudsburg, PA, USA, January 2022. Association for Computational Linguistics.
- Stephen Crain and Mineharu Nakayama. Structure dependence in grammar formation. *Language* (*Baltim.*), 63(3):522, September 1987.

566 Alexander D'Amour, Katherine Heller, Dan Moldovan, Ben Adlam, Babak Alipanahi, Alex Beutel, 567 Christina Chen, Jonathan Deaton, Jacob Eisenstein, Matthew D Hoffman, Farhad Hormozdiari, 568 Neil Houlsby, Shaobo Hou, Ghassen Jerfel, Alan Karthikesalingam, Mario Lucic, Yian Ma, Cory 569 McLean, Diana Mincu, Akinori Mitani, Andrea Montanari, Zachary Nado, Vivek Natarajan, 570 Christopher Nielson, Thomas F Osborne, Rajiv Raman, Kim Ramasamy, Rory Sayres, Jessica 571 Schrouff, Martin Seneviratne, Shannon Sequeira, Harini Suresh, Victor Veitch, Max Vladymyrov, Xuezhi Wang, Kellie Webster, Steve Yadlowsky, Taedong Yun, Xiaohua Zhai, and D Sculley. 572 Underspecification presents challenges for credibility in modern machine learning. Journal of 573 Machine Learning Research, 23(226):1–61, 2022. 574

- Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali Farhadi, Hannaneh Hajishirzi, and Noah Smith.
 Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping. *arXiv* [cs.CL], February 2020.
 - Robert Frank and Donald Mathis. Transformational networks. *Models of Human Language Acquisition*, 22, 2007.
- Nan Gao and Qingshun He. A dependency distance approach to the syntactic complexity variation in
 the connected speech of alzheimer's disease. *Humanit. Soc. Sci. Commun.*, 11(1), August 2024.
- Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias
 Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nat. Mach. Intell.*, 2 (11):665–673, November 2020.
 - Katherine L Hermann and Andrew K Lampinen. What shapes feature representations? exploring datasets, architectures, and training. *arXiv* [cs.LG], June 2020.
- John Hewitt and Christopher D Manning. A structural probe for finding syntax in word representations. In *Proceedings of the 2019 Conference of the North*, pp. 4129–4138, Stroudsburg, PA, USA, 2019. Association for Computational Linguistics.
- 593 Michael Y Hu, Angelica Chen, Naomi Saphra, and Kyunghyun Cho. Latent state models of training dynamics. *arXiv* [*cs.LG*], August 2023.

594 595 596	Kuan-Hao Huang, Varun Iyer, I-Hung Hsu, Anoop Kumar, Kai-Wei Chang, and Aram Galstyan. ParaAMR: A large-scale syntactically diverse paraphrase dataset by AMR back-translation. In
597	1: Long Papers), Stroudsburg, PA, USA, 2023. Association for Computational Linguistics.
598	Yufei Huang Shengding Hu Xu Han Zhiyuan Liu and Maosong Sun Unified view of grokking
299	double descent and emergent abilities: A perspective from circuits competition. arXiv [cs.LG].
601	February 2024.
601	
602	Jeevesh Juneja, Rachit Bansal, Kyunghyun Cho, João Sedoc, and Naomi Saphra. Linear connectivity
603	reveals generalization strategies. arXiv [cs.LG], May 2022.
604	
606	December 2014.
607 608 609	David Lightfoot. The child's trigger experience: Degree-0 learnability. <i>Behavioral and brain sciences</i> , 12(2):321–334, 1989.
610 611	Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg. Assessing the ability of LSTMs to learn syntax-sensitive dependencies. <i>Trans. Assoc. Comput. Linguist.</i> , 4:521–535, December 2016.
612	Ziming Lin Eric I Michaud and May Termory, Omnigraly, Chaldring bayond algorithmic date
614	arXiv [cs.LG], October 2022.
615	Brian MacWhinney. The childes project: Tools for analyzing talk, volume II: The database. Psychol-
616	ogy Press, London, England, 3 edition, January 2014.
618 619 620	R Thomas McCoy, Robert Frank, and Tal Linzen. Revisiting the poverty of the stimulus: hierarchical generalization without a hierarchical bias in recurrent neural networks. <i>arXiv</i> [cs.CL], February 2018.
621 622 623	R Thomas McCoy, Ellie Pavlick, and Tal Linzen. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. <i>arXiv [cs.CL]</i> , February 2019.
624 625 626 627	R Thomas McCoy, Robert Frank, and Tal Linzen. Does syntax need to grow on trees? sources of hierarchical inductive bias in sequence-to-sequence networks. <i>Trans. Assoc. Comput. Linguist.</i> , 8: 125–140, December 2020.
628 629 630	Tom McCoy, Robert Frank, and Tal Linzen. Does syntax need to grow on trees? sources of hierarchical inductive bias in sequence-to-sequence networks. https://rtmccoy.com/rnn_hierarchical_biases.html. Accessed: 2024-11-20.
631 632 633	William Merrill, Nikolaos Tsilivis, and Aman Shukla. A tale of two circuits: Grokking as competition of sparse and dense subnetworks. <i>arXiv [cs.LG]</i> , March 2023.
634 635 636	Aaron Mueller and Tal Linzen. How to plant trees in language models: Data and architectural effects on the emergence of syntactic inductive biases. <i>arXiv [cs.CL]</i> , May 2023.
637 638 639 640	 Aaron Mueller, Robert Frank, Tal Linzen, Luheng Wang, and Sebastian Schuster. Coloring the blank slate: Pre-training imparts a hierarchical inductive bias to sequence-to-sequence models. In <i>Findings of the Association for Computational Linguistics: ACL 2022</i>, Stroudsburg, PA, USA, 2022. Association for Computational Linguistics.
641 642 643 644 645 646	Aaron Mueller, Albert Webson, Jackson Petty, and Tal Linzen. In-context learning generalizes, but not always robustly: The case of syntax. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pp. 4761–4779, Stroudsburg, PA, USA, 2024. Association for Computational Linguistics.
647	Shikhar Murty, Pratyusha Sharma, Jacob Andreas, and Christopher D Manning. Characterizing intrinsic compositionality in transformers with tree projections. <i>arXiv</i> [cs.CL], November 2022.

648 649	Shikhar Murty, Pratyusha Sharma, Jacob Andreas, and Christopher Manning. Grokking of hierarchical structure in vanilla transformers. In <i>Proceedings of the 61st Annual Meeting of the Association for</i>
650 651	<i>Computational Linguistics (Volume 2: Short Papers)</i> , Stroudsburg, PA, USA, 2023. Association for Computational Linguistics
652	
653	Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. Stress
654	test evaluation for natural language inference. In <i>Proceedings of the 27th International Conference</i> on Computational Linguistics, pp. 2340–2353, 2018.
655	
656 657	Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures for grokking via mechanistic interpretability. <i>arXiv</i> [<i>cs.LG</i>], January 2023.
658	Cethering Oleson Nelson Elhago Neel Nanda Nicholes Jesonh Neve DesSerme Tem Hanishan
659	Ben Mann Amanda Askell Yuntao Bai Anna Chen Tom Conerly Dawn Drain Deen Ganguli
660	Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane
661	Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish,
662 663	and Chris Olah. In-context learning and induction heads. arXiv [cs.LG], September 2022.
664	Isabel Papadimitriou and Dan Jurafsky. Learning music helps you read: Using transfer to study
665	linguistic structure in language models. In <i>Proceedings of the 2020 Conference on Empirical</i>
666	Methods in Natural Language Processing (EMNLP), Stroudsburg, PA, USA, 2020. Association for
667	Computational Linguistics.
669	
000	Isabel Papadimitriou and Dan Juratsky. Injecting structural hints: Using language models to study
009	inductive blases in language learning. arXiv [cs.CL], April 2023.
674	Jackson Petty and Robert Frank. Transformers generalize linearly. arXiv [cs.CL], September 2021.
670	
673	Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: General- ization beyond overfitting on small algorithmic datasets. <i>arXiv</i> [cs.LG], January 2022.
674	Less Derde Marco Der Alle Dere De der Derdellet erd Edit Gerdi Gredienerie
675	Jorge Ramirez, Marcos Baez, Auday Berro, Boualem Benatalian, and Fabio Casati. Crowdsourcing
676	Information Systems Engineering: 34th International Conference, CAISE 2022, Lawan, Belaium
677 678	June 6–10, 2022, Proceedings, pp. 253–269, Berlin, Heidelberg, 2022. Springer-Verlag.
679	Naomi Saphra and Adam Lopez. Understanding learning dynamics of language models with SVCCA.
621	<i>uraiv [cs.cL]</i> , November 2018.
682	Naomi Saphra and Adam Lopez. Understanding learning dynamics of language models with. In
683	Proceedings of the 2019 Conference of the North, pp. 3257–3267, Stroudsburg, PA, USA, January
684	2019. Association for Computational Linguistics.
685	Rylan Schaeffer, Brando Miranda, and Sanmi Koveio. Are emergent abilities of large language
686	models a mirage? arXiv [cs.AI], April 2023.
687	Thibault Sellam, Steve Yadlowsky, Ian Tenney, Jason Wei, Naomi Saphra. Alexander D'Amour. Tal
688	Linzen, Jasmijn Bastings, Iulia Raluca Turc, Jacob Eisenstein, Dipanjan Das, and Ellie Pavlick.
689	The MultiBERTs: BERT reproductions for robustness analysis. In International Conference on
690	Learning Representations, October 2021.
691	Verifier Calinian Paris The Target Philippine and and the Philippine
692	Yewei Song, Cedric Lotnritz, Daniel lang, legawende F Bissyande, and Jacques Klein. Revisiting
693	code similarity evaluation with abstract syntax tree edit distance. arxiv [cs.CL], April 2024.
694	Marthinus Wilhelmus Theunissen, Marelie Davel, and Etienne Barnard. Benign interpolation of noise
695	in deep learning. S. Afr. Comput. J., 32(2), December 2020.
696	Vila (V. D. I. Cl. 1. 7. L. R. G. L. K. K. J. D. K. L. D. K. L. L. L.
697	vikrani varina, Konin Snan, Zachary Kenton, Janos Kramar, and Kamana Kumar. Explaining
698	grokking unough cheun enciency. <i>arkiv [cs.LG]</i> , September 2025.
699	Kenneth Wexler. Formal principles of language acquisition, 1980.
700	Vaizhang Zhang and Dannig Shasha. Simple fast algorithms for the addition distance between trees
701	and related problems. <i>SIAM J. Comput.</i> , 18(6):1245–1262, December 1989.

Xuekai Zhu, Yao Fu, Bowen Zhou, and Zhouhan Lin. Critical data size of language models from a grokking perspective. arXiv [cs.CL], January 2024. Yes
Advance Endy, Taky Te, Dowent Shou, and Endomation. Circled data size of hanguage models from a grokking perspective. arXiv [cs.CL], January 2024. 709 711 712 713 714 715 716 717 718 719 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 730 731 732 733 734 735 736 737 738
200 200 200 200 200 200 200 200 200 200 200 200 201 200 201 200 201 200 201 200 201 200 201 200 201 200 202 200 203 200 204 200 205 200 206 200 207 200 208 200 209 200 200 200 201 200 202 200 203 200 204 200 205 200 206 200 207 200 208 200 209 200 200 200 201 200 202
708 709 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 730 734 735 736 737 738 739 730 731 732 733 734 735 736 737 738 739 734 735 736 737 7
709 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 745 746 747
710 712 713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748
111 112 113 114 115 116 117 118 119 120 121 122 123 124 125 126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144 145 146 147 148
112 113 114 115 116 117 118 119 120 121 122 123 124 125 126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144 145 145 146 147 148
713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 745 746 747
114 115 116 117 118 119 120 121 122 123 124 125 126 127 128 129 120 121 122 123 124 125 126 127 128 129 120 121 122 123 124 125 126 127 128 129 129 121 122 123 124 125 126 127 128 129 129 120 121 122 123 124 124 1
113 116 117 118 119 120 121 122 123 124 125 126 127 128 129 130 131 132 133 134 135 136 137 138 139 141 142 143 144 145 146 147 148
110 111 112 121 122 123 124 125 126 127 128 129 130 131 132 133 134 135 136 137 138 139 134 135 136 137 138 139 130 131 132 133 134 135 136 137 138 139 130 131 132 133 134 135 136 137 138 139 130 131 132 133 134 1
711 712 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 745 746 747 748
710 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748
122 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748
722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748
723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748
724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747
725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748
726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748
727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748
728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748
729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748
730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748
731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747
732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748
733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748
734 735 736 737 738 739 740 741 742 743 744 745 746 747 748
735 736 737 738 739 740 741 742 743 744 745 746 747 748
736 737 738 739 740 741 742 743 744 745 746 747 748
737 738 739 740 741 742 743 744 745 746 747 748
738 739 740 741 742 743 744 745 746 747 748
739 740 741 742 743 744 745 746 747 748
740 741 742 743 744 745 746 747 748
741 742 743 744 745 746 747 748
742 743 744 745 746 747
743 744 745 746 747 748
744 745 746 747 748
745 746 747 748
747 748
747
7/0
750
751
752
753
754

756 A RELATED WORK EXTENDED

758 A.1 SYNTAX AND HIERARCHICAL GENERALIZATION

While works mentioned in Section 2.1 focused on models trained from scratch, another line of
research examined the inductive bias of pretrained models. Mueller et al. (2024); Mueller & Linzen
(2023) pretrained transformers on text corpora such as Wikipedia and CHILDES (MacWhinney,
2014) before fine-tuning them on the question formation task. They found that exposure to large
amounts of natural language data enables transformers to generalize hierarchically.

764

Instead of using the question formation task as a probe, Hewitt & Manning (2019); Murty et al. (2022) directly interpreted model's internal representation to understand whether transformers constrain their computations to to follow tree-structure patterns. Hewitt & Manning (2019) demonstrated that the syntax tress are embedded in model's representation space. Similarly, Murty et al. (2022) projects transformers into a tree-structured network, and showed that transformers become more tree-like over the course of training on language data.

A.2 RANDOM VARIATION

772 Specific training choices, such as hyperparameters, are crucial to model outcomes. However, even 773 when controlling for these factors, training machine learning models remains inherently stochas-774 tic—models can be sensitive to random initialization and the order of training examples. Zhou et al. (2020); D'Amour et al. (2022); Naik et al. (2018) reported significant performance differences across 775 model checkpoints on various analysis and stress test sets. Zhou et al. (2020) further found that 776 instability extends throughout the training curve, not just in final outcomes. To investigate the source 777 of this inconsistency, Dodge et al. (2020) compared the effects of weight initialization and data order, 778 concluding that both factors contribute equally to variations in out-of-sample performance. 779

Similarly, Sellam et al. (2021) found that repeating the pre-training process on BERT models can result
 in significantly different performances on downstream tasks. To promote more robust experimental
 testing, they introduced a set of 25 BERT-BASE checkpoints to ensure that experimental conclusions
 are not influenced by artifacts, such as specific instances of the model. In this work, we also
 observe training inconsistencies across runs on OOD data, both during training and at convergence.
 Unlike prior studies that focus on implications of random variations on experimental design, we
 study the source of training inconsistencies and link these inconsistencies to simplicity bias and the
 characteristics of the training data.

788 A.3 SIMPLICITY BIAS

Models often favor simpler functions early in training, a phenomenon known as simplicity bias (Hermann & Lampinen, 2020), which is also common in LMs. Choshen et al. (2022) found that early LMs behave like n-gram models, and Saphra & Lopez (2019) observed that early LMs learn simplified versions of the language modeling task. McCoy et al. (2019) showed that even fully trained models can rely on simple heuristics, like lexical overlap, to perform well on Natural Language Inference (NLI) tasks. Chen et al. (2023) further explored the connection between training dynamics and simplicity bias, showing that simpler functions learned early on can continue to influence fully trained models, and mitigating this bias can have long-term effects on training outcomes.

797

Phase transitions have been identified as markers of shifts from simplistic heuristics to more complex model behavior, often triggered by the amount of training data or model size. In language models, Olsson et al. (2022) showed that the emergence of induction heads in autoregressive models is linked to handling longer context sizes and in-context learning. Similar phase transitions have been studied in non-language domains, such as algorithmic tasks (Power et al., 2022; Merrill et al., 2023) and arithmetic tasks (Nanda et al., 2023; Barak et al., 2022).

00

In the context of hierarchical generalization, Ahuja et al. (2024) used a Bayesian approach to analyze the simplicity of hierarchical versus linear rules in modeling English syntax. They argued that transformers favor the hierarchical rule because it is simpler than the linear rule. However, their model fails to explain (1) why learning the hierarchical rule is delayed (i.e., after learning the linear rule) and (2) why hierarchical generalization is inconsistent across runs. In this work, we offer a different perspective, showing that a model's simplicity bias towards either rule is driven by the characteristics of the training data.



Figure 7: **Components of the original QF and TI training data.** *Left:* QF training data contains samples of two tasks types: question formation and declaration copying. We further break down samples in the declaration copying task by branching type. We also breakdown center-embedded sentences based on whether the main subject serves the subject or object in the embedded clause. *Right:* TI training data also contains samples of two task types: tense inflection and past tense copying. Similar to QF, we further breakdown tense inflection samples by branching types, and center-embedded sentences (in the tense inflection samples) by subject/object type.

824 825 826

827

830

831 832

833

834

835

836

837

838

839

840

843

844

845

846 847

819

820

821

822

823

B TRAINING DATA SAMPLES

828 B.1 QUESTION FORMATION 829

When we mention "declarations," we are referring to the declaration copying task, and "questions" refer to the question formation task. Here are two examples randomly taken from the training data:

- Declaration Example: our zebra doesn't applaud the unicorn . decl our zebra doesn't applaud the unicorn .
- Question Example: some unicorns do move . quest do some unicorns move ?

Both tasks begin with an input declarative sentence, followed by a task indicator token (decl or quest), and end with the output. During training, the entire sequence is used in the causal language modeling objective. The in-distribution validation set and the OOD generalization set only contain question formation samples.

841 B.2 TENSE INFLECTION

- Past Example: our peacocks above our walruses amused your zebras . PAST our peacocks above our walruses amused your zebras .
- Present Example: your unicorns that our xylophones comforted swam . PRESENT your unicorns that our xylophones comfort swim .

The tense inflection task is indicate by the PRESENT token, and in Section 4.3, we only used tense inflection samples during training. In Appendix E.1, we further explore the use of a secondary copying task to achieve OOD generalization. Similar to the question formation training data, the secondary task only requires repeating the given sentence, which is always in the past tense, and the copying task is marked by the PRESENT token.

853 854

855

C FURTHER PARTITIONS ON CENTER-EMBEDDED SENTENCES

C.1 Two subtypes of center-embedded sentences

In Section 4, we showed that center-embedded sentences drive hierarchical generalization in both the
QF and TI tasks. Here, we further partition center-embedded sentences based on the syntactic role of
the *main subject* (i.e., the subject of the main clause) within the modifying clause. Specifically, we
classify them into two types:

- Subject-type: The main subject serves as the subject within the clause. Example: *The keys that unlock the cabinet are on the table.*
- 2. Object-type: The main subject serves as the object within the clause. Example: *The keys that the bear uses are on the table.*

This partition is motivated by their distinct subject-verb dependency patterns. In subject-type sentences, both the main verb (from the main clause) and the embedded verb (from the relative clause) depend on the main subject. In contrast, object-type sentences exhibit a nested subject-verb structure. Our goal is to investigate whether differences in subject-verb dependency patterns influence the model's preference for the hierarchical rule.



Figure 8: Both subtypes of center-embedded sentences induces hierarchical generalization in **OF.** We train models on datasets containing different ratios of object-type v.s. subject-type centerembedded sentences. We then evaluate on models on two OOD generalization set, one containing unambiguous object-type center-embedded sentences and the other unambiguous subject-type center-embedded sentences.

C.2 QF TASK

We first investigate whether the two subtypes of center-embedded sentences differentially influence the model's preference for the hierarchical rule in the QF task. For all training data variants, we fix 50K question formation samples and 50K declaration copying samples, with the latter containing only center-embedded sentences but varying the ratio between the two subtypes. To analyze generalization behavior on a more granular level, we partition the generalization set (composed solely of center-embedded sentences) into the two subtypes as well. Models are trained on 30 random seeds, and results are shown in Figure 8. Regardless of the data mix, the model consistently favors the hierarchical rule across both partitions of the generalization set. This suggests that, for question formation, both subtypes of center-embedded sentences equally contribute to the model's ability to identify the main auxiliary.



Figure 9: Object-type center-embedded sentences gives a stronger bias towards hierarchical generalization in TI. We train models on datasets containing different ratios of object-type v.s. subject-type center-embedded sentences. We then evaluate on models on three OOD generalization set, one containing unambiguous object-type center-embedded sentences, one unambiguous subjecttype center-embedded sentences, and one unambiguous right-branching sentences.

C.3 TI TASK

We repeat a similar experiment for the TI task, fixing the total number of tense inflection samples to 100K. As shown in Section 4.3, models exhibit the strongest hierarchical generalization when trained on primarily center-embedded sentences. Therefore, in the following data variants, 99% of the samples are center-embedded sentences, with the remaining 1% being right-branching sentences. Within the center-embedded samples, we vary the ratio between the two subtypes. To evaluate

918 generalization, we split the generalization set into three groups: the two subtypes of center-embedded 919 sentences and right-branching sentences. Models trained on 30 random seeds show that, across 920 all three generalization sets, accuracy is positively correlated with the proportion of object-type 921 center-embedded sentences (Figure 9). However, even when models are trained predominantly on 922 subject-type center-embedded sentences (teal violins in Figure 9), they still show a strong tendency toward hierarchical generalization. Thus, while both subtypes drive hierarchical generalization in 923 TI, object-type center-embedded sentences have a stronger effect. Notably, the original TI training 924 data includes a higher proportion of right-branching sentences (shown in 7) and a higher ratio of 925 subject-type center-embedded sentences-both of which are suboptimal for inducing hierarchical 926 generalization. 927

928 929

D VARYING DATA RATIOS FOR QUESTION FORMATION

930 **Data composition details** We construct variations of the training data using the following procedure. Each new 931 dataset contains 50K questions (reused from the original 932 data) and 50K declarations, where we control the ratio 933 between center-embedded and right-branching sentences. 934 These datasets are used for the experiments in Section 5.2. 935 To generate additional declarations, we keep the distri-936 bution of the unique syntax structures in original dataset. 937 Specifically, for each sentence in the original data, we 938 extract the syntax tree using the CGF rules and resample 939 words from the vocabulary to create new sentence samples.





Figure 10: Hierarchical generalization in QF is sensitive to compositions of declaration-copying samples.

945 Figure 10. First, note that there is a sharp performance drop between the blue bar and the right-most 946 green bar. This sharp transition indicates that mixing in as little as 1% of right-branching declarations 947 significantly reduces the model's likelihood of generalizing hierarchically. Interestingly, when the dataset is predominantly right-branching declarations, models consistently achieve 0% generalization 948 accuracy, indicating a strong preference for the linear rule across all training runs. However, note 949 that there is another sharp transition between the red bard and the left-most green bar. This transition 950 indicates that as soon as we remove the 1% of center-embedded sentences, the model fails to learn 951 either the linear rule or the hierarchical rule. As a result, the generalization accuracy is close to 952 random guess ($\sim 25\%$). This transition is closely studied in Section 6.1, where we examine how data 953 diversity leads to rule commitment.

954 955 956

E ADDITIONAL RESULTS ON TENSE INFLECTION

957 E.1 A SECONDARY TASK IS NOT NECESSARY

In the original of TI training data (McCoy et al., 2020), 958 a secondary task is also included to mimic the question 959 formation training data. In this secondary task, instead of 960 transforming a sentence from the past tense to the present 961 tense, the model simply needs to repeat it. For concrete 962 examples, see Appendix B. Figure 7 (right) shows a break-963 down of the two tasks in the original TI training data. In 964 experiments conducted in Section 3.2, we have eliminated 965 the used of this secondary task because center-embedded 966 sentences can be included in the tense inflection train-967 ing samples *without* violating the ambiguity requirement. 968 Here, we use the training data originally proposed by Mc-Coy et al. (2020) to confirm that the use of secondary task 969 is indeed not necessary. Specifically, we remove all the 970 past-tense-copying samples from the original training data 971 and train models on the tense-inflection task only. We



Figure 11: **Past-copy task is not neces**sary to induce hierarchical generalization in TI.

evaluate the model's generalization performance on two OOD set containing unambiguous right-973 branching and unambiguous center-embedded sentences, shown in Figure 11. We can see that the 974 model's OOD performances are the same with or without the secondary task.



Figure 12: Total Variation v.s. final generalization accuracy for TI task. Similar to Figure 5, we observe the same horseshoe shaped behavior between training stability and final generalization accuracy on right-branching sentences for the TI task.

We repeat the same total variation analysis in Section 5 for the tense inflection task. We use the data mixes from Section 4.3. Specifically, we include only tense inflection samples and vary the ratio between linear-inducing (i.e., right-branching) and hierarchical-inducing (i.e., center-embedded) 992 sentences. In Section 4.3, we have already concluded that the hierarchical rule is *always* preferred for center-embedded sentences regardless of data mixes. For this reason, we are interested in examining the rule preference and training stability for unambiguous right-branching sentences. In Figure 995 12 we visualize the relationship between total variation and the final generalization accuracy on unambiguous right-branching sentences. The qualitative behavior is similar to what we have observed 997 in QF (Section 5.2). 998

F TRAINING INSTABILITY



Figure 13: Each training run either stabilizes in a simple OOD generalization rule or oscillates 1012 in its OOD accuracy. The OOD generalization behaviors can be either stable or unstable when 1013 trained on different seeds. We use total variation to quantify the instability within one training run. 1014

In Figure 14, we visualize the training dynamics for 30 independent runs when trained on the original 1016 QF data. Each run differs in both model initialization and data order. Notice that the training 1017 dynamics for runs exhibit grokking behaviors: OOD generalization is delayed when compared 1018 to training loss convergence and validation performance convergence. These runs share a similar 1019 progression in training loss, validation accuracy, and generalization accuracy up until moment when 1020 the training loss converges. Interestingly, after convergence on training loss, all runs reach 0% on 1021 the generalization set, indicating that the model strictly prefers linear rules on OOD data. After that, models start to achieve non-trivial performance in generalization accuracy. However, for many runs the generalization accuracy does not increase monotonically. Instead, we observe massive 1023 swings in generalization accuracy during this training period as well as large inconsistency across 1024 different seeds. Overall, training is *always* stable for ID data while the performance for OOD data is 1025 inconsistent across seeds. We visualize runs with different of total variation values in Figure 13.



986

987

988 989 990

991

993

994

996

999

1000

1001

1002

1003

1004

1007

1008

1009

1010

1011

1015

19



Figure 14: **Training Dynamics on original question formation data.** Training loss and indistribution validation accuracy is stable during training and consistent across random seeds. In contrast, the model's performance on OOD data is both unstable during training and inconsistent across seeds. The instability and inconsistency is most prominent during grokking (i.e., when training loss has converged).

G DATA DIVERSITY AND MEMORIZATION PATTERNS



Figure 15: **OOD generalization v.s. syntax similarity to training data.** At low data diversity, model memorizes syntax patterns and applies the hierarchical rule only syntax structures similar to ones in the training data. With higher data diversity, model extrapolates rules and can apply the hierarchical rule even to unseen syntax structures that are dissimilar to training data.

1062

1044 1045

1046 1047

1048

1049

1050

1051

1052 1053

1054

1055

1056

1057

We investigate model behavior when trained on data with limited diversity. By analyzing a model's generalization accuracy across different syntactic types, we aim to distinguish patterns indicative of either memorization or generalization.

Measuring data similarity Building on the diversity measure from Section 6.1, we now use Tree-1067 Edit Distance (TED) as a measure of sentence similarity. As before, we first construct syntax trees 1068 using CFG rules, then calculate TED using the Zhang-Shasha Tree-Edit Distance algorithm (Zhang & 1069 Shasha, 1989). We define TED=0 for sentences that share the same syntax structure but differ only in vocabulary. This similarity measure allows us to quantify, for each sample in the OOD generalization 1070 set, the closest matching sentence type in the training data. In the memorization regime, where the 1071 model encounters only a few syntax types, we suspect it cannot extrapolate rules to syntactically 1072 distinct OOD sentences. In contrast, with a more diverse syntax exposure, rule extrapolation may enable the model to apply rules even to OOD sentence types. 1074

Experiment To verify our intuition about memorization and generalization, we train models on
 two variations of the QF data. In the first variation, the declaration-copying task has data diversity set
 to 1, meaning only one syntax type appears, and we specifically choose one with center embedding.
 In the second variation, the declaration-copying task has diversity set to 5, with all 5 types containing
 center embeddings. For both datasets, the question-formation task remains unchanged, consisting
 solely of right-branching sentences. For the diversity=1 dataset, we calculate TED for each unique

syntax type in the OOD set against the single syntax type in the declaration-copying task. For the diversity=5 dataset, we compute TED between each OOD sample and the five syntax types in the declaration-copying task, taking the minimum. This TED score provides a measure of similarity between the OOD samples and those encountered during training. Our goal is to determine, based on training with these datasets, which OOD syntax types the model applies the hierarchical rule to.

Result In Figure 15, we visualize the final generalization accuracy for each OOD syntax type against its TED relative to the training data. When trained on low-diversity data (Figure 15, *left*), generalization accuracy is negatively correlated with TED. For syntax types seen in the declarationcopying task (TED=0) and those similar to it, the model applies the hierarchical rule. However, for syntax types with high TED, the model's behavior is random (25%), indicating failure to follow any rule. As data diversity increases slightly (Figure 15, *right*), generalization accuracy no longer correlates with TED, suggesting that once the model begins to extrapolate the hierarchical rule, it can apply this rule to a wider range of OOD syntax types.