

Mission Accomplished?

Recovering Information from ‘Impossible’ Languages with LLMs

Anonymous ACL submission

Abstract

Understanding whether large language models operate under constraints comparable to human linguistic cognition remains a central question in AI and cognitive science. While there has been some research into whether LLMs can learn linguistically possible and impossible languages, it is less clear whether they can systematically recover linguistic structure and meaning from systematically degraded input. In this work, we investigate whether LLMs can translate and recover impossible languages back into possible forms and whether there is any difference between different types of impossible languages for recovery. By fine-tuning GPT-2 on several perturbation types, we find that models can reconstruct grammatically well-formed output, with performance systematically modulated by the nature of the perturbation. Models trained on longer sentences benefit from richer training contexts, although longer sequences also increase the difficulty of resolving non-local dependencies. Overall, our findings indicate that LLMs display a preference for local over distant dependencies, yet can still overcome structural violations that render input unintelligible, revealing a partial alignment between neural architectural constraints and human linguistic biases.

https://anonymous.4open.science/r/impossible_translation-9C4D

1 Introduction

Human cognitive constraints determine which languages humans can learn, understand, and use (Chomsky, 1957, 1986; Jackendoff, 2002; Pinker, 1994). One goal of linguistics is to distinguish natural languages from other systems (de Saussure, 1983; Hockett, 1958; Lyons, 1968). Hence, linguists have proposed various descriptive frameworks, with some theories of universal grammar focusing on word order patterns and relationships among sentence elements, drawing on cross-linguistic typology (Greenberg, 1963; Hawkins,

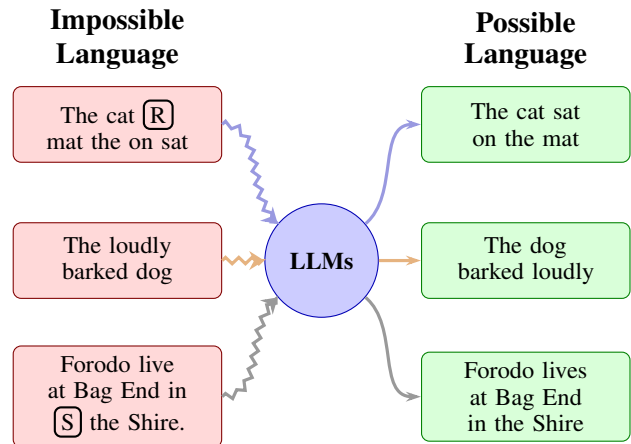


Figure 1: Translation task overview. Three separate LLMs are fine-tuned on specific perturbation types (LOCALSHUFFLE, PARTIALREVERSE, WORDHOP). Each model translates its corresponding impossible language input (left) into possible language output (right), recovering disrupted natural structure.

1983; Dryer, 1992), formal semantics (Barwise and Cooper, 1981; Keenan and Stavi, 1986), syntactic theory (Vennemann, 1974; Primus, 2001; Hawkins, 1994), and generative grammar (Chomsky, 1981, 1995).

A central criticism of LLMs is that they fail to capture constraints defining human linguistic competence. Chomsky and colleagues argue that LLMs cannot distinguish between possible and impossible languages, as they can accommodate systems violating universal grammatical principles, thereby undermining their status as models of the human language faculty (Chomsky, 2023; Chomsky et al., 2023). Moro et al. (2023) argue that impossible languages exist only for humans, as neurobiological architecture constrains learnable grammars, whereas LLMs lack such constraints. Extending this view, Bolhuis et al. (2024) contend that LLMs do not model human language representation or acquisition, since their behaviour arises from statistical optimisation rather than biological mechanisms.

064 While some argue that LLMs process impossible
065 and possible languages similarly (Ziv et al., 2025),
066 Kallini et al. (2024) demonstrate systematic learn-
067 ing differences, showing a clear bias toward natu-
068 ral language structures. This suggests that while
069 LLMs can technically learn impossible languages,
070 the principles defining human linguistic feasibility
071 still influence their performance.

072 A fundamental principle underlying human lan-
073 guage is *information locality*, which posits that
074 semantically and syntactically related elements ap-
075 pear close together in linear order (Gibson et al.,
076 2000; Futrell et al., 2020). The *Dependency Lo-
077 cality Theory* (DLT) formalizes this, predicting
078 that processing difficulty increases as distance be-
079 tween dependent elements grows, reflecting cog-
080 nitive constraints on working memory during in-
081 cremental comprehension (Gibson, 1998). Cross-
082 linguistic research supports this through depen-
083 dency length minimisation (DLM), showing natu-
084 ral languages systematically organise word order to
085 keep grammatically related elements close together.
086 (Hawkins, 1994; Liu, 2008; Futrell et al., 2020;
087 Temperley, 2018). When information locality is
088 violated through scattered dependent elements with
089 intervening unrelated material, structures become
090 cognitively inaccessible to humans despite intact
091 lexical content, rendering text incomprehensible
092 because information cannot be efficiently accessed
093 (Levy, 2008; Hahn et al., 2020).

094 This study examines whether LLMs can trans-
095 late linguistically impossible languages into possi-
096 ble forms by recovering information rendered inac-
097 cessible through structural violations. We applied
098 perturbation functions to the BabyLM dataset that
099 systematically altered sentence structures, intro-
100 ducing linguistic inconsistencies absent in natural
101 languages (Warstadt et al., 2023; Sinha et al., 2021).
102 We then fine-tuned three separate GPT-2 models
103 on these transformed datasets to reconstruct valid
104 forms (Radford et al., 2019; Ebrahimi et al., 2020).
105 As Figure 1 illustrates, each model translates im-
106 possible language inputs into possible language
107 outputs, recovering disrupted natural structure (one
108 model for each impossible language).

109 Some perturbations employed in this study sys-
110 tematically violate information locality. By scram-
111 bling word order or inserting distance-based mark-
112 ers, these transformations increase dependency
113 lengths and disrupt local relationships, enabling in-
114 cremental comprehension (Gibson et al., 2000; Liu,
115 2008). Since intelligible natural language requires

116 semantically related words to appear close together,
117 translating an impossible language requires mod-
118 els to identify scattered semantic dependencies and
119 restore characteristic proximity. This raises a criti-
120 cal question: can LLMs, not subject to the same
121 working memory limitations as humans, recover in-
122 formation made inaccessible by locality violations?

123 Investigating whether neural language models
124 can translate these structurally degraded inputs
125 back into forms respecting information locality pro-
126 vides insight into LLM architectural biases and
127 their relationship to cognitive constraints defining
128 possible human languages (Kallini et al., 2024; Xu
129 et al., 2025; Tran et al., 2018; Khandelwal et al.,
130 2018). Beyond theoretical implications for under-
131 standing neural architectures and human linguistic
132 constraints, this information recovery capability
133 has practical applications in processing corrupted
134 text, restoring degraded documents, and handling
135 noisy linguistic input. This study provides insights
136 into limits and capabilities of LLMs in language
137 learning, demonstrating their potential to interpret
138 and transform linguistically impossible inputs into
139 coherent outputs, contributing to understanding ma-
140 chine learning applications in linguistics, cognitive
141 science, and artificial intelligence¹.

142 2 Background and Previous Work

143 Chomsky’s Universal Grammar (UG) framework
144 distinguishes possible from impossible human lan-
145 guages, proposing that all human languages share
146 constrained structural principles from innate cog-
147 nitive mechanisms (Chomsky, 1965). Crucially, for-
148 mal inadequacy (whether grammatical formalisms
149 describe natural language structure) differs from
150 cognitive impossibility (whether humans can learn
151 a linguistic system) (Moro, 2016; Musso et al.,
152 2003). Impossible languages are defined by un-
153 learnable grammars which humans cannot acquire
154 despite exposure due to cognitive constraints, rather
155 than failure to fit formal grammatical models. Em-
156 pirical evidence suggests languages relying solely
157 on position-based rules without hierarchical struc-
158 ture are impossible, conflicting with core syntactic
159 processing mechanisms and fundamentally concern-
160 ing learnability rather than formal expressiveness
161 (Musso et al., 2003; Smith and Tsimpli, 1993).

162 Nefdt (2024) argues identifying possible gram-
163 mar boundaries requires grounding in theoretical

¹All data and source code for this paper are available on
GitHub

Language	Example 1	Example 2
ORIGINAL TEXT	It is nice in there	we 'd need to look at it again , would n't we
LOCALSHUFFLE	there It in is nice	we 'd need to it look again at , would n't we
PARTIALREVERSE	It is R there in nice	we 'd need R we n't would , again it at look to
WORDHOP	It be nice in there S	we 'd need to look at P it P again , would n't we P

Table 1: Example data for each perturbation function. For ease of comparison, we use a similar visualisation to Kallini et al. (2024).

linguistic competence, not solely distributional regularities. Smith and Tsimpli (1993)’s work on the polyglot savant Christopher showed successful acquisition of an artificial language consistent with universal principles but failure to learn one violating them, suggesting certain grammars remain unlearnable even for exceptionally linguistically talented individuals. Neurolinguistic studies show brain regions for syntactic processing selectively respond to natural language-compatible structures, supporting that impossible languages lack neurobiological realisation (Moro, 2016; Musso et al., 2003). However, empirical evidence remains remarkably limited. Smith and Tsimpli (1993) focus on a single individual and most theoretical claims lack systematic empirical testing. This gap highlights computational approaches’ potential value, as LLMs offer systematic ways to test which linguistic systems are truly unlearnable, providing empirical grounding for theoretical claims about possible human language boundaries.

Against this theoretical and cognitive backdrop, artificial neural networks examined as linguistic competence models reveal important discrepancies. Hahn (2020) shows standard self-attention architectures face intrinsic limitations modelling unbounded hierarchical dependencies, challenging their adequacy as full natural language syntax models. Mitchell and Bowers (2020) demonstrate that recurrent neural networks learn highly unnatural patterns humans do not acquire, highlighting misalignment between human and neural inductive biases. While approximating complex string distributions, they are not automatically constrained by possible human languages and may generalise divergently from human cognition.

Studies of transformer capabilities show mixed results. Abdou et al. (2022) show that models trained on shuffled text retain non-trivial word order information, while Huang et al. (2023) demonstrate that lexically invariant transformer vari-

ants match standard models under long contexts. Ebrahimi et al. (2020) show that self-attention networks learn certain hierarchical languages with appropriate cues, though Delétang et al. (2023) argue standard architectures face principled challenges representing complex hierarchical dependencies without explicit structured memory.

Recent research comparing human learners and LLMs shows that both benefit from compositional, structured input (Galke et al., 2024), with transformers partially encoding hierarchical phrase structure though not fully matching generative grammar expectations (Alleman et al., 2021). Kallini et al. (2024) demonstrate that LLMs exhibit systematic learnability differences between possible and impossible languages, displaying biases toward natural patterns while struggling with highly unnatural systems. Xu et al. (2025) find that GPT-2 shows learning slowdowns aligned with typological universals when acquiring implausible counterfactual languages, though it can ultimately learn them. While LLMs can learn highly unnatural patterns, their learning is shaped by architectural and data-driven biases that only partially overlap with the cognitive constraints defining human language, motivating examination of whether LLMs can systematically recover possible languages from impossible ones by reversing perturbations violating natural linguistic structure.

3 Method

The main goal of this research is to train an LLM to translate impossible language into possible language. Formally, we define a perturbation function $f : A \rightarrow A'$ where A represents standard English text and A' is the perturbed version. Our objective is to train a model \mathcal{M} such that $\mathcal{M}(A') := A$, effectively reversing the perturbation and recovering the original linguistic information.

A critical constraint is that \mathcal{M} must not have

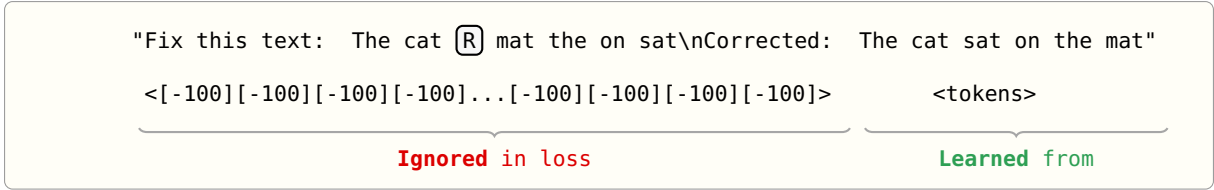


Figure 2: The mechanism of masking labels and calculating loss

244 been pre-trained on A . If \mathcal{M} had prior exposure to
 245 standard English, it could leverage that knowledge
 246 rather than genuinely learning to recover linguistic
 247 structure from perturbed input, confounding our
 248 evaluation of information recovery. Therefore, \mathcal{M}
 249 must be pre-trained on an impossible language and
 250 then fine-tuned for the translation task. To satisfy
 251 this requirement, we adopt the models pre-trained
 252 by Kallini et al. (2024) as the starting point for
 253 our fine-tuning process, since these models were
 254 trained from scratch on impossible languages, gen-
 255 erated by applying different perturbation functions
 256 to standard English sentences.

257 3.1 Data Generation

258 The BabyLM (Warstadt et al., 2023) dataset was
 259 used for these experiments, specifically selecting
 260 the *bnc_spoken* and *Gutenberg* subsets to anal-
 261 yse the impact of contrasting sentence lengths on
 262 translation performance. Applying perturbations
 263 to these subsets yields paired examples (perturbed,
 264 original) for training.

265 Following Kallini et al. (2024), we apply three
 266 perturbation functions that systematically violate
 267 natural characteristics of language.

268 **LOCALSHUFFLE:** This perturbation randomly
 269 reorders words within a local window, disrupt-
 270 ing the sequential arrangement while maintaining
 271 words within bounded distances.

272 **PARTIALREVERSE:** A random starting point
 273 is selected within the sentence, and an [R] token is
 274 placed in this position, and subsequent tokens are
 275 reversed in order. This creates a partition where
 276 the initial segment remains unchanged, the final
 277 segment reversed, and the [R] marks the boundary.

278 **WORDHOP:** This perturbation violates the
 279 principle that no human language requires counting
 280 words for grammatical operations (Musso et al.,
 281 2003) by adding markers ([S] for singular, [P]
 282 for plural) at fixed distances after verbs based on
 283 subject-verb agreement.

284 Table 1 presents examples of each function.

285 We created three datasets, each containing paired

286 texts (impossible, possible), by applying one per-
 287 turbation strategy to each. Preprocessing included
 288 tokenization, lemmatization, POS tagging, and con-
 289 stituency parsing, enabling grammatical structure-
 290 dependent perturbations such as WORDHOP’s
 291 verb phrase identification and marker placement.
 292 To investigate training scale effects, we created
 293 datasets of 10K and 100K sample pairs from
 294 *bnc_spoken*.

295 3.2 Model Fine-tuning

296 The models were fine-tuned by combining the
 297 causal language modelling (CLM) paradigm of
 298 GPT-2 (Radford et al., 2019) with a masked la-
 299 bel strategy. Each training sample consists of a
 300 pair: an impossible text and its corresponding pos-
 301 sible version. These pairs are formatted into the
 302 following instruction-response structure:

303 Fix this text: <impossible_text>
 304 Corrected: <possible_text><|endoftext|>

305 To ensure the model focuses on producing cor-
 306 rect output rather than memorising instructions,
 307 masked labelling is applied. The instruction por-
 308 tion (from sequence start to the word Corrected)
 309 receives label value $[-100]$, while the response
 310 portion (correct text and <|endoftext|> token)
 311 retains its tokens as labels. Only these tokens con-
 312 tribute to the cross-entropy loss, enabling the model
 313 to learn accurate corrections conditioned on the in-
 314 struction. Figure 2 displays the label masking and
 315 loss calculation.

316 3.3 Experimental Design

317 We conduct three experiments testing whether
 318 LLMs exhibit locality constraints when translating
 319 impossible languages. We hypothesise that transla-
 320 tion performance degrades as perturbations inten-
 321 sify locality violations. Our perturbations systemat-
 322 ically violate this by scattering related words across
 323 sentences, altering the normal subject-verb agree-
 324 ment rule. If LLMs exhibit locality constraints,

translation performance should degrade as violations intensify: longer sentences increase difficulty by expanding distances between related elements and enlarging reordering search spaces, while perturbation types should vary by locality disruption severity.

In all experiments, models were trained on NVIDIA H100 96GB GPUs. We train with a batch size of 512. Training the 100K dataset requires approximately 1 GPU hour per model.

Experiment 1: Effect of Training Sample Size.

This experiment investigates how training data size influences translation performance. Models were fine-tuned on 10K and 100K subsets of *bnc_spoken*, with learning progress evaluated at multiple checkpoints using Exact Match (EM), measuring perfect reconstruction (Rajpurkar et al., 2016), and BLEU score quantifying n-gram overlap (Papineni et al., 2002). These metrics suit the task because impossible and possible texts share vocabulary but differ in token order, making both perfect reconstruction (EM) and partial accuracy (BLEU) informative. The hypothesis is that larger datasets will improve performance across all perturbations, but with different learning rates: WORDHOP should be learned fastest due to preserved word order, while LOCALSHUFFLE and PARTIALREVERSE will require more examples to recover disrupted locality.

Experiment 2: Effect of Sentence Length.

Building on the locality principle introduced above, this experiment directly tests how sentence length affects the difficulty of recovering grammatical structure. Results were compared between the *bnc_spoken* dataset (shorter sentences, avg. 12 tokens) and the *Gutenberg* dataset (longer sentences, avg. 40 tokens). The hypothesis is that longer sentences provide more training signal, potentially improving overall performance, but simultaneously increase difficulty for perturbations disrupting information locality. Specifically, LOCALSHUFFLE and PARTIALREVERSE are expected to show complex effects: benefits from additional training content may be offset by challenges from related words scattered over greater distances, which exponentially expand the space of possible reorderings and integration costs. WORDHOP is expected to benefit more straightforwardly from longer sentences since it preserves underlying word order and only adds markers at fixed distances, maintaining most local dependencies intact.

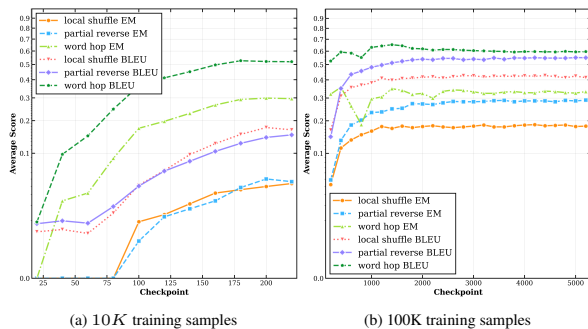


Figure 3: Training progress on *bnc_spoken* with (a) 10K and (b) 100K samples. WORDHOP learns fastest, reaching $\sim 60\%$ BLEU and $\sim 30\%$ EM on 100K data, while LOCALSHUFFLE and PARTIALREVERSE converge more slowly at lower performance levels. For larger plots, see Figure 6.

Experiment 3: Quality of Generated Text.

This experiment evaluates the extent to which generated texts are similar to human-written text using perplexity, computed under pretrained GPT-2 to measure how well generated text aligns with natural language distributions (Jelinek et al., 2005). A lower perplexity value indicates that a text is more consistent with natural language distributions (Radford et al., 2019). We hypothesise that prediction perplexity will converge toward that of the original text as training progresses, with rates varying by perturbation type. Specifically, LOCALSHUFFLE and PARTIALREVERSE predictions are expected to show steeper perplexity reductions as these perturbations more severely disrupt natural word order, while WORDHOP predictions maintain slightly higher perplexity even after convergence due to potential artifacts from marker removal that do not fully restore natural linguistic flow.

4 Experimental results

Experiment 1: Figure 3 shows that training size critically affects performance, with 100K samples substantially outperforming 10K across all perturbations. However, improvement rates vary: WORDHOP’s high performance (0.6 EM) reflects the tractability of marker removal when word order is preserved, while LOCALSHUFFLE’s low ceiling (0.18 EM) even with 100K samples indicates that additional data cannot fully overcome the difficulty of recovering severely disrupted locality, suggesting fundamental architectural limitations in reconstructing scrambled dependencies. Qualitative examples illustrating the learning progression at different checkpoints are provided in Appendix

Dataset		Size	Average Length
<i>bnc_spoken</i>	train	10K	12.54
<i>bnc_spoken</i>	train	100K	12.77
<i>bnc_spoken</i>	test	1K	11.72
<i>Gutenberg</i>	train	10K	40.56
<i>Gutenberg</i>	train	100K	40.61
<i>Gutenberg</i>	test	1K	46.78

Table 2: Average sentence lengths for the training and test datasets at 10K and 100K sample sizes for training and 1K for testing.

(Tables 3, 4, and 5).

Experiment 2: Table 2 shows average sentence lengths for both datasets. Comparing final performance at 100K samples (Figures 4 and 3) reveals substantial sentence length effects that vary by perturbation type. On *Gutenberg* (avg. 40 tokens), WORDHOP achieves approximately 78% BLEU compared to 60% on *bnc_spoken* (avg. 12 tokens), representing an 18 percentage point improvement. PARTIALREVERSE shows even larger gains, improving from 55% BLEU on short sentences to 82% on long sentences (27 percentage points). In contrast, LOCALSHUFFLE demonstrates minimal improvement, rising from 18% EM to only 25% EM (7 percentage points), remaining challenging regardless of sentence length. Analysis of performance improvement by input text length class is shown in Appendix Figure 9.

This differential response indicates that WORDHOP’s preserved word order and PARTIALREVERSE’s partial structure enable better exploitation of additional training content in longer sentences, whereas LOCALSHUFFLE’s complete locality disruption presents fundamental recovery difficulties that additional context cannot overcome. Figure 4 further confirms that while model performance improves on longer texts, learning progress is notably slower than on shorter sentences.

For perturbations that directly disrupt intra-sentential dependencies, sentence length becomes a crucial factor: as sentences grow longer, dependency distances increase and the perturbation affects a larger portion of the sequence, thereby enlarging the effective search space of possible reorderings. This is consistent with prior work showing that neural language models rely more heavily on nearby context and struggle with long-range dependencies and hierarchical structure, especially when architectural capacity does not scale with in-

put length (Khandelwal et al., 2018; Hahn, 2020; Tran et al., 2018). It also aligns with findings that models trained on shuffled or perturbed word orders retain only partial and noisy order information, with performance degrading as structures become more globally disrupted (Sinha et al., 2021; Abdou et al., 2022).

Experiment 3: Although our model is fine-tuned from GPT-2, we use the pre-trained (base) GPT-2 model for evaluation to measure how well the generated text aligns with general language patterns learned from large-scale pretraining data: lower perplexity indicates that the text is more consistent with natural language distributions.

Figure 5 compares impossible language inputs with model predictions against baseline unperturbed text (perplexity = 129.34). Results reveal striking differences across perturbation types. LOCALSHUFFLE input exhibits extremely high perplexity (3902.34), which model predictions reduce significantly to 212.40, representing a 94.6% decrease. This substantial reduction demonstrates successful reconstruction approaching natural language distributions, though predictions remain above baseline, indicating incomplete restoration. PARTIALREVERSE shows similar patterns with less extreme values. Input perplexity (1500.38) reduces to 146.56, a 90.2% reduction. Predictions come closer to baseline than LOCALSHUFFLE (146.56 vs. 212.40), suggesting more successful structural recovery, aligning with higher exact match scores (0.29 vs. 0.18). WORDHOP presents a contrasting pattern. Input perplexity (168.11) is already near baseline, as this perturbation preserves underlying word order. However, predictions show increased perplexity (243.81), a 45.0% increase. While the model successfully removes inserted markers, the reconstruction introduces subtle artefacts deviating from standard distributions, likely from imperfect verb agreement handling. Concrete examples of text quality at different training stages are provided in Appendix Tables 3, 4, and 5.

5 Discussion

Recent research demonstrates that LLMs exhibit disparities in their capacity to acquire impossible languages, attributed to underlying inductive biases (Kallini et al., 2024; Xu et al., 2025). Building on this foundation, we adopt a reverse approach, examining whether large language models can recover accessible linguistic structure from systemat-

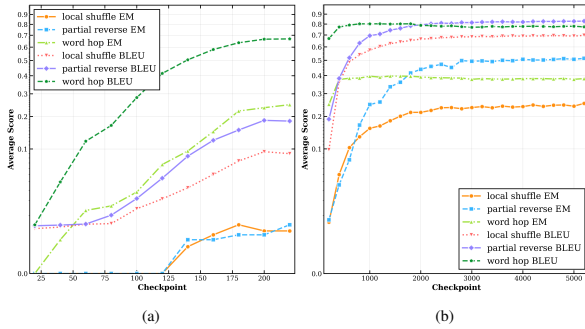


Figure 4: Training progress on *Gutenberg* with (a) 10K and (b) 100K samples. Longer sentences yield higher performance than *bnc_spoken*. For larger plots, see Figure 7.

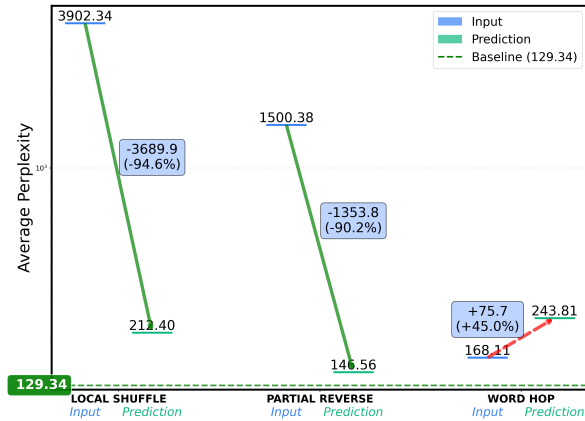


Figure 5: Perplexity of impossible inputs (blue) versus model predictions (green) relative to baseline unperturbed text (129.34, dashed line), evaluated using pretrained GPT-2. For larger plots, see Figure 8.

ically degraded input. Specifically, we investigate whether models are influenced by universal principles constraining human linguistic feasibility and whether their ability to recover information is affected by the locality of the perturbation. We argue that this serves as a test of the Information Locality hypothesis, which holds that natural languages tend to group related elements closer together. Our experiments with GPT-2, pre-trained on impossible languages by Kallini et al. (2024) and fine-tuned on three perturbation types (LOCALSHUFFLE, PARTIALREVERSE, WORDHOP), reveal important findings about neural language models’ capabilities and limitations in recovering linguistically possible forms from impossible ones.

Information Locality Our results demonstrate that model performance directly reflects the degree to which perturbations violate information locality (Gibson et al., 2000). LOCALSHUFFLE and PARTIALREVERSE, which disperse grammati-

cally related elements across sentences and create non-local dependencies, proved most difficult (0.18 and 0.29 EM respectively), requiring models to reconstruct hierarchical structures from scrambled input. Conversely, WORDHOP, which preserves underlying word order and maintains most local dependencies, achieved substantially higher performance (>0.6 EM) despite introducing distance-based markers. Crucially, perplexity analysis reveals that models genuinely recover underlying structure rather than memorizing mappings. The dramatic perplexity reductions for LOCALSHUFFLE (94.6%) and PARTIALREVERSE (90.2%) toward the baseline indicate successful reconstruction of locality-respecting configurations, where scattered elements are reordered into interpretable structures restoring local dependencies (Hahn et al., 2020; Hahn, 2020). This convergence demonstrates that translation effectively recovers information rendered inaccessible by locality violations, though WORDHOP’s increased perplexity suggests artifacts from imperfect marker removal.

Performance across perturbation functions

Our results demonstrate GPT-2 can perform information recovery, with performance varying systematically by perturbation type. WORDHOP proved easiest, achieving exact match scores above 0.6 on 100K samples, as it preserves underlying word order despite introducing distance-based markers. In contrast, LOCALSHUFFLE and PARTIALREVERSE plateaued around 0.18 and 0.29 EM respectively, requiring reconstruction of hierarchical syntactic structures from scrambled input where local dependencies become non-local. This difficulty aligns with theoretical arguments about transformer architectural limitations in capturing unbounded hierarchical dependencies (Tran et al., 2018; Delétang et al., 2023; Hahn, 2020; Khandelwal et al., 2018). Our work examines the reverse task: whether models pre-trained on impossible languages can recover possible language structure through translation. The pattern of difficulty we observe in translation—WORDHOP easiest, LOCALSHUFFLE hardest—mirrors the pattern Kallini et al. (2024) found for learning these languages, suggesting the same architectural biases that make certain impossible languages harder to acquire also make them harder to translate. Perplexity analysis confirms models genuinely recover underlying structure rather than memorizing mappings. For LOCALSHUFFLE and PARTIALREVERSE, prediction

518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568

perplexity decreased toward original unperturbed text, with fine-tuned models producing sequences increasingly consistent with pre-trained GPT-2 distributions (Radford et al., 2019; Holtzman et al., 2020), demonstrating effective recovery of information rendered inaccessible by locality violations.

Importance of sentence length Our second experiment revealed strong sentence length effects relating to information locality and data availability. Models trained on *Gutenberg* outperform *bnc_spoken* as longer sentences provide more training content. However, longer sequences exponentially expand possible reorderings and distances between originally adjacent elements (Gibson, 1998), increasing integration costs for reconstructing relationships as processing costs grow with dependency length (Futrell et al., 2020; Gibson et al., 2000). The model must recover information over longer distances while identifying which scattered elements belong together, similar to the challenge of processing deeply nested structures (Levy, 2008). This reveals important implications for LLM inductive biases. While longer sequences provide more training signal and thereby enable better performance, they simultaneously present greater challenges for recovering non-local dependencies. Though LLMs lack human cognitive-neurobiological constraints, they exhibit architectural constraints where learning efficiency balances data availability against dependency distance, creating biases partially aligned with natural language properties and preferring local over non-local dependencies (Futrell et al., 2020; Khandelwal et al., 2018).

6 Conclusion

This study demonstrates that LLMs can translate linguistically impossible inputs into possible forms by recovering information obscured through information locality violations. Performance varies systematically (WORDHOP easiest, LOCALSHUFFLE hardest), with longer sentences providing more training content while increasing recovery difficulty. While LLMs are sensitive to locality constraints, they overcome intractable violations.

Our findings reveal incomplete alignment between LLMs and human linguistic competence. Models show sensitivity to information locality yet operate with distinct inductive biases, indicating that putative constraints on human learnability do not necessarily impose bounds on learnability by

neural networks. The same architectural properties constrain both acquisition and translation, revealing biases partially aligned with cognitive constraints.

This raises intriguing questions about child language acquisition: Could similar approaches help explain how children extract linguistic structure from limited and imperfect input (Chomsky, 1986; Pinker, 1994)? More speculatively, these methods might help analyse extraterrestrial communication systems violating human universals, with models distinguishing incomprehensible from unfamiliar aspects. Systematic investigation of impossible languages through computational models reveals boundaries of linguistic possibility beyond human cognition.

This study’s approach could be applied to some real-world problems, such as recovering disrupted information locality, enabling error correction, text normalisation, and robustness to noisy input. Future work should investigate scaling to larger architectures, generalisation to novel perturbations, and cross-linguistic extension, contributing to understanding the computational basis of language and boundaries defining linguistic possibility.

7 Limitations and Future Directions

It is important to acknowledge several limitations of this study. First, our experiments focus on a relatively small model (GPT-2) and a limited set of perturbation functions. Larger, more recent architectures with enhanced capabilities for modeling long-range dependencies might exhibit different learning patterns or overcome some of the limitations we observed. Second, we evaluate performance primarily through exact match, BLEU scores, and perplexity, which capture surface-level accuracy but may not fully reflect whether the model has genuinely learned the underlying principles of information locality restoration or merely exploited distributional shortcuts.

Third, our impossible languages, while designed to violate linguistic universals, represent only a small subset of the space of possible perturbations. Future work could explore a broader range of impossible language constructions, including those that violate other aspects of linguistic structure beyond word order, and evaluate whether consistent patterns emerge in model behavior. Additionally, investigating whether models can generalize their translation abilities to novel perturbation types or combinations would provide insight into whether

669	they learn abstract principles of information local-	Noam Chomsky, Ian Roberts, and Jeffrey Watumull.	718
670	ity recovery or task-specific mappings.	2023. The false promise of chatgpt. <i>The New York</i>	719
671	The notion of "impossible language" itself re-	<i>Times</i> . Guest essay, March 8, 2023.	720
672	quires careful interpretation. While linguistic the-	Ferdinand de Saussure. 1983. <i>Course in General Lin-</i>	721
673	ory identifies certain systems as impossible for hu-	<i>guistics</i> . Duckworth, London. R. Harris (Trans.).	722
674	man acquisition based on violations of informa-	Original work published 1916.	723
675	tion locality and other neurobiological and cogni-	Grégoire Delétang, Anian Ruoss, Jordi Grau-Moya, Tim	724
676	tive constraints, these systems remain well-defined	Genewein, Kevin Li Wenliang, Elliot Catt, Chris	725
677	computational problems that neural networks can,	Cundy, Marcus Hutter, Shane Legg, Joel Veness,	726
678	in principle, solve. Our results highlight an im-	and Pedro A. Ortega. 2023. Neural networks and	727
679	portant distinction between cognitive impossibility,	the chomsky hierarchy. In <i>Proceedings of the In-</i>	728
680	which arises from violations of information local-	<i>ternational Conference on Learning Representations</i>	729
681	ity, and computational possibility. Recognizing this	(<i>ICLR 2023</i>).	730
682	gap is essential for understanding both the potential	Matthew S. Dryer. 1992. The greenbergian word order	731
683	and the limits of LLMs as models of human lan-	correlations. <i>Language</i> , 68(1):81–138.	732
684	guage and as practical language processing tools.	Javid Ebrahimi, Dhruv Gelda, and Wei Zhang. 2020.	733
685	References	How can self-attention networks recognize dyck-n	734
686	Mostafa Abdou, Vinit Ravishankar, Artur Kulmizev, and	languages? In <i>Findings of the Association for Com-</i>	735
687	Anders Søgaard. 2022. Word order does matter (and	<i>putational Linguistics: EMNLP 2020</i> , pages 4307–	736
688	shuffled language models know it). In <i>Proceedings</i>	4313.	737
689	<i>of the 60th Annual Meeting of the Association for</i>	Richard Futrell, Roger Levy, and Edward Gibson. 2020.	738
690	<i>Computational Linguistics (ACL 2022)</i> , pages 6904–	Dependency locality as an explanatory principle for	739
691	6919.	word order. <i>Language</i> , 96(2):372–413.	740
692	Matteo Alleman, Jonathan Mamou, Miguel A. Del Rio,	Lukas Galke, Yoav Ram, and Limor Raviv. 2024. Deep	741
693	Hanlin Tang, Yoon Kim, and SueYeon Chung. 2021.	neural networks and humans both benefit from com-	742
694	Syntactic perturbations reveal representational corre-	positional language structure. <i>Nature Communica-</i>	743
695	lates of hierarchical phrase structure in pretrained lan-	<i>tions</i> , 15:10816.	744
696	guage models. In <i>Proceedings of the 6th Workshop</i>	Edward Gibson. 1998. Linguistic complexity: locality	745
697	<i>on Representation Learning for NLP (RePLANLP-</i>	of syntactic dependencies. <i>Cognition</i> , 68(1):1–76.	746
698	<i>2021)</i> , pages 263–276.	Edward Gibson et al. 2000. The dependency locality	747
699	Jon Barwise and Robin Cooper. 1981. Generalized	theory: A distance-based theory of linguistic com-	748
700	quantifiers and natural language. <i>Linguistics and</i>	plexity. <i>Image, language, brain</i> , 2000:95–126.	749
701	<i>Philosophy</i> , 4(2):159–219.	Joseph H. Greenberg. 1963. Some universals of gram-	750
702	Johan J Bolhuis, Stephen Crain, Sandiway Fong, and	mar with particular reference to the order of mean-	751
703	Andrea Moro. 2024. Three reasons why ai doesn't	ingful elements. In Joseph H. Greenberg, editor,	752
704	model human language. <i>Nature</i> , 627(8004):489.	<i>Universals of Language</i> , pages 73–113. MIT Press,	753
705	Noam Chomsky. 1957. <i>Syntactic Structures</i> . Mouton,	Cambridge, MA.	754
706	The Hague.	Michael Hahn. 2020. Theoretical limitations of self-	755
707	Noam Chomsky. 1965. <i>Aspects of the Theory of Syntax</i> .	attention in neural sequence models. <i>Transactions of</i>	756
708	MIT Press, Cambridge, MA.	<i>the Association for Computational Linguistics</i> , 8:156–	757
709	Noam Chomsky. 1981. <i>Lectures on Government and</i>	171.	758
710	<i>Binding</i> . Foris, Dordrecht.	Michael Hahn, Dan Jurafsky, and Richard Futrell. 2020.	759
711	Noam Chomsky. 1986. <i>Knowledge of Language: Its</i>	Universals of word order reflect optimization of gram-	760
712	<i>Nature, Origin, and Use</i> . Praeger, New York.	mars for efficient communication. <i>Proceedings of the</i>	761
713	Noam Chomsky. 1995. <i>The Minimalist Program</i> . MIT	<i>National Academy of Sciences</i> , 117(5):2347–2353.	762
714	Press, Cambridge, MA.	John A. Hawkins. 1983. <i>Word Order Universals</i> . Aca-	763
715	Noam Chomsky. 2023. Noam chomsky on language,	ademic Press, New York.	764
716	left libertarianism, and progress . Podcast episode.	John A. Hawkins. 1994. <i>A Performance Theory of Or-</i>	765
717	Conversations with Tyler, Episode 182.	<i>der and Constituency</i> . Cambridge University Press,	766
		Cambridge.	767
		Charles F. Hockett. 1958. <i>A Course in Modern Linguis-</i>	768
		<i>tics</i> . Macmillan, New York.	769

Table 3: Learning process for translation LOCALSHUFFLE at different learning checkpoints in a successful example. The **red** colour determines the difference between the current text and with original text, and the **green** colour determines similarity.

Input Text	Type
Reggie was able to find out for them the place at which Fleming was putting up in Boston. Having ascertained from the clerk that he was still staying there, the next thing was to get hold of Louis Anderson.	original text
able Reggie was to find out the for them place at which putting Fleming was Boston in up ascert. Havingained the from that clerk he still was staying the there, next was thing hold to get Louis of Anderson .	LOCALSHUFFLE
“Steve had to find out what they were going to put him in the hotel, the next thing was that was that he was going to put on to Tom. Having he was staying in the Boston car, which was the next thing	Checkpoint-200
Reggie was able to find out for the them place at which Fleming was putting up in Boston. Having ascertained from the clerk that he was still staying there, the next thing was to get on Louis Anderson.	Checkpoint-600
Reggie was able to find out for them the place at which Fleming was putting up in Boston. Having ascertained from the clerk that he was still staying there, the next thing was to get to Louis Anderson of.	Checkpoint-1000
Reggie was able to find out for them the place at which Fleming was putting up in Boston. Having ascertained from the clerk that he was still staying there, the next thing was to get hold of Louis Anderson.	Checkpoint-1400

A Larger Plots

Because the images embedded in the main text are relatively small, this appendix provides enlarged versions to improve readability and highlight finer details.

B Examples

During training, we evaluated model performance at multiple checkpoints to monitor learning progress. We tracked generated text quality across checkpoints for each impossible language type. This section provides examples illustrating the improvement trajectory throughout training.

Additionally, we investigate how performance varies by input text length for the LOCALSHUFFLE perturbation. Since LOCALSHUFFLE uses a window size of 3, we classify input texts by the number of shuffling windows they contain. Specifically, we define $\text{class}(\text{text}_i) = \lceil \text{len}(\text{text}_i) / 3 \rceil$, where $\text{len}(\text{text}_i)$ is the number of tokens in the text. This groups texts by how many 3-token windows were shuffled during perturbation. Figure 9 shows model improvement across these classes, revealing how performance scales with the number of shuffled segments.

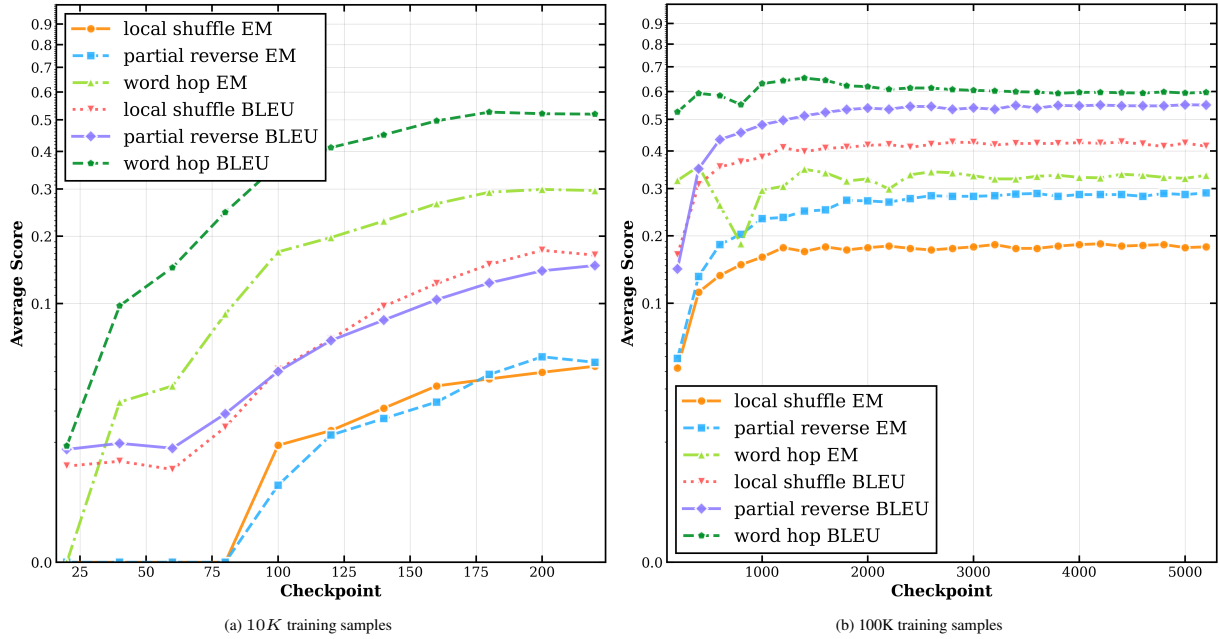


Figure 6: Training progress on *bnc_spoken* with (a) 10K and (b) 100K samples. WORDHOP learns fastest, reaching ~60% BLEU and ~30% EM on 100K data, while LOCALSHUFFLE and PARTIALREVERSE converge more slowly at lower performance levels.

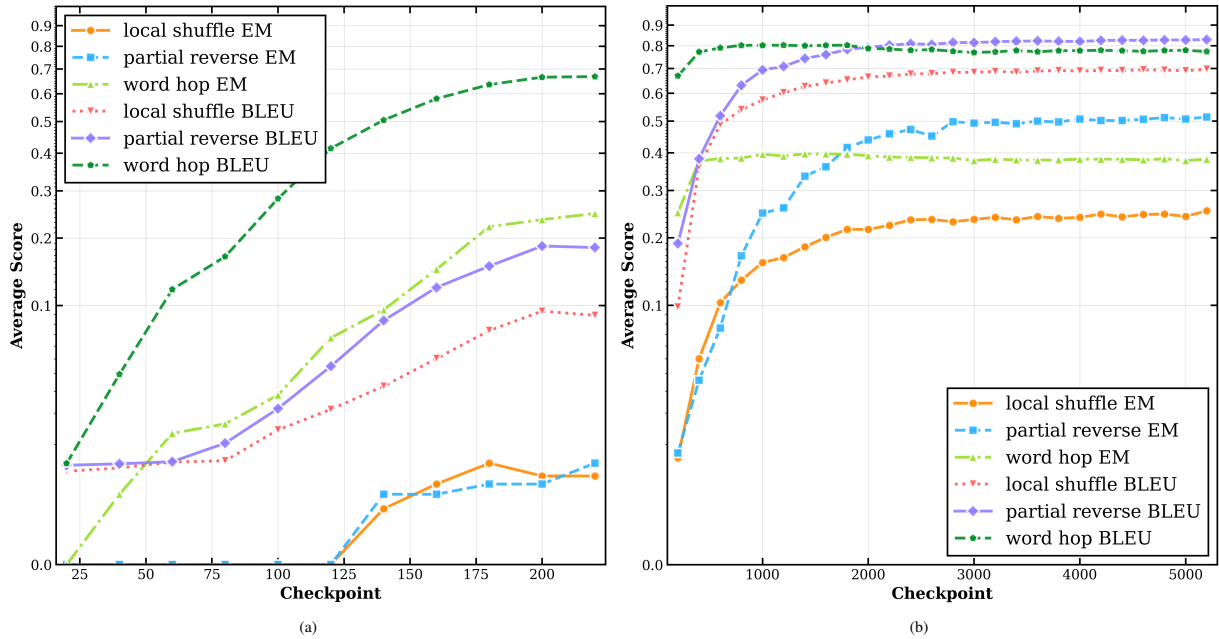


Figure 7: Training progress on *Gutenberg* with (a) 10K and (b) 100K samples. Longer sentences yield higher performance than *bnc_spoken*

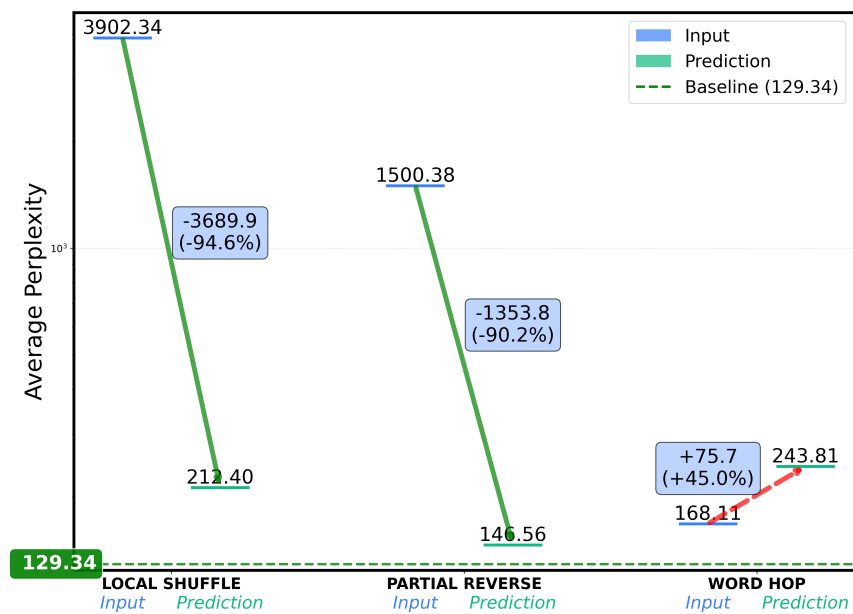


Figure 8: Perplexity of impossible inputs (blue) versus model predictions (green) relative to baseline unperturbed text (129.34, dashed line), evaluated using pretrained GPT-2.

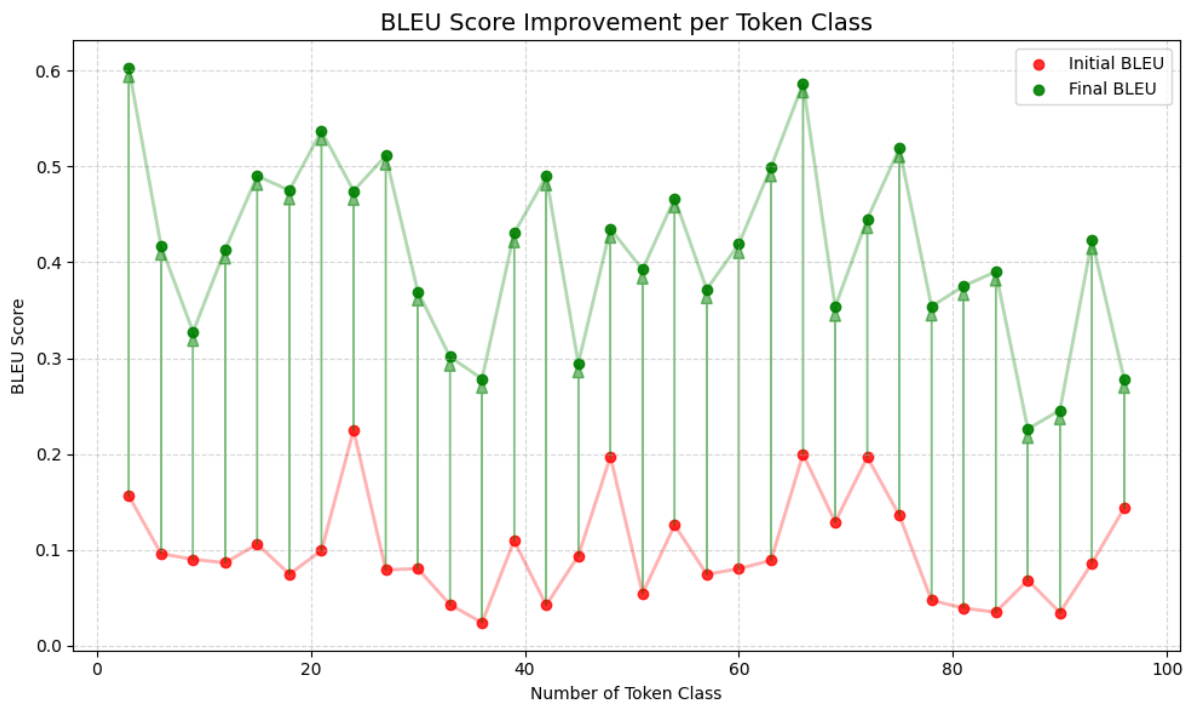


Figure 9: Improvement of the BLEU score for each class. The red line corresponds to the BLEU score before translating and for shuffled text, and the green line is for texts after translation by the final model.

Table 4: Learning process for translation PARTIALREVERSE at different learning checkpoints in a successful example. The **red** colour determines the difference between the current text and with original text, and the **green** colour determines similarity.

Input Text	Type
Bunny and Sue, so impatient they could hardly keep still, waited. They heard the front door open. They heard their father talking. Then came a funny, squeaking, whining sound.	original text
Bunny and Sue , so impatient (R) . sound whining , aking sque , funny a came Then . talking father their heard They . open door front the heard They . waited , still keep hardly could they	PARTIALREVERSE
Bessie and Sue, so impatient and impatiently. The children were not allowed to keep them shut up. They were not allowed to keep them shut. They were not allowed to keep them shut, and they were not allowed to	Checkpoint-200
Bucky and Sue, so impatient they could hardly keep still, waited, waited, waited, waited. They heard their father talking. Then Then came a funny, squeaky voice.	Checkpoint-600
Bumper and Sue, so impatient they could hardly keep still, waited. They heard the front door open. They heard their father talking. Then came a funny, squeaking, whining sound.	Checkpoint-1000
Bunny and Sue, so impatient they could hardly keep still, waited. They heard the front door open. They heard their father talking. Then came a funny, squeaking, whining sound.	Checkpoint-1400

Table 5: Learning process for translation PARTIALREVERSE at different learning checkpoints in a successful example. The **red** colour determines the difference between the current text and with original text, and the **green** colour determines similarity.

Input Text	Type
"That's what I say!" came from Mart. "This is no night for a little boy to be out. It's snowing harder than ever."	original text
"That be what I say!" came (S) from Mart. "This (P) be no night for a (S) little boy to be out. It be snowing harder than ever (S) ."	WORDHOP
"That's what I say!" came (S) from Mart. "This** what I say!" came (S) from Mart. "This## no night for a (S) little boy to be out. It's snowing	Checkpoint-200
"That's what I say!" came (S) 'from Mart. "This is no night for a (S) little boy to be out. It's snowing harder than ever"	Checkpoint-600
"That's what I say!" came from Mart. "This is no night for a little boy to be out. It's snowing harder than ever."	Checkpoint-1000