

Developmentally-plausible Working Memory Shapes a Critical Period for Language Acquisition

Anonymous ACL submission

Abstract

Large language models possess general linguistic abilities but acquire language less efficiently than humans. This study proposes a method for integrating the developmental characteristics of working memory during the critical period, a stage when human language acquisition is particularly efficient, into the training process of language models. The proposed method introduces a mechanism that initially constrains *working memory* during the early stages of training and gradually relaxes this constraint in an exponential manner as learning progresses. Targeted syntactic evaluation shows that the proposed method outperforms conventional methods without memory constraints or with static memory constraints. These findings not only provide new directions for designing data-efficient language models but also offer indirect evidence supporting the role of the developmental characteristics of working memory as the underlying mechanism of the critical period in language acquisition.

1 Introduction

Large language models (LLMs) exhibit general linguistic abilities comparable to those of humans; however, their efficiency in language acquisition remains far inferior. It has been noted that LLMs require data quantities that are three to four orders of magnitude larger than those needed for humans to achieve comparable performance across many evaluation metrics (Warstadt et al., 2023). This disparity in data efficiency reflects the current reliance of LLMs on scaling and suggests not only a significant potential for improving learning efficiency but also the possibility of drawing *insights* from human language processing and acquisition.

An important theoretical framework for understanding the efficiency of human language acquisition is the **Critical Period Hypothesis (CPH)** (Lenneberg, 1967). The CPH posits that there is a specific period during which language

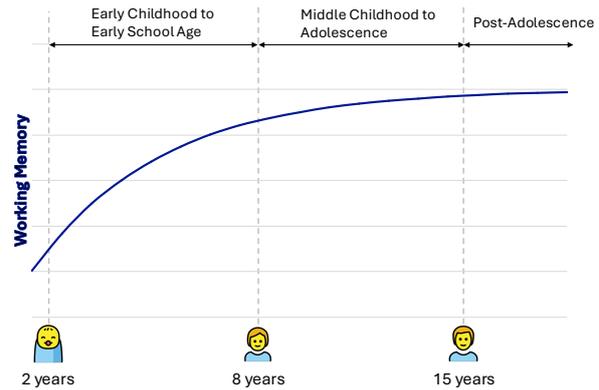


Figure 1: Developmental trajectory of human working memory

can be acquired efficiently, and that this ability diminishes thereafter. Various studies, including cases of limited first language (L_1) exposure during childhood and age-related effects on second language (L_2) acquisition, support the existence of a critical period (CP) (Fromkin et al., 1974; Curtiss, 1977; Johnson and Newport, 1989). However, the reasons why children acquire language more efficiently than adults remain partially unresolved. One compelling explanation for the CP in L_1 acquisition is the **Less-is-More Hypothesis** (Newport, 1990), which argues that children’s cognitive limitations (e.g., working memory capacity and attentional scope) are advantageous for language learning. According to this hypothesis, children’s limited processing capacities enable them to efficiently extract fundamental patterns and structures (e.g., grammatical rules) from linguistic input, whereas adults, with their greater cognitive capacities, are more likely to be distracted by complex information, thereby hindering rule acquisition.

Inspired by the “Less-is-More” hypothesis, we use language models (LMs) to study the CP for language acquisition, focusing on L_1 acquisition and investigating whether integrating human cognitive developmental characteristics, particularly the

developmental properties of *working memory* (Figure 1), into LMs can facilitate efficient language acquisition. Specifically, we propose a method for incorporating the exponential increase in working memory capacity that corresponds to the CP into LMs and analyze its impact on learning efficiency. Using a GPT-2 model (Radford et al., 2019) trained on a Child-Directed Speech (CDS) dataset (Huebner and Willits, 2021), we conduct evaluation experiments with Zorro (Huebner et al., 2021), a targeted syntactic evaluation benchmark specialized for CDS. The results demonstrate that a cognitively plausible model, which initially restricts working memory and gradually relaxes this constraint exponentially as training progresses, outperforms models without memory constraints or with static memory constraints. These findings provide new insights into designing data-efficient LMs, contributing to the field of **natural language processing**, while also offering indirect evidence supporting the role of the developmental characteristics of working memory as the underlying mechanism of the CPH in human language acquisition, contributing to the field of **cognitive science**.

2 Related Work

2.1 Critical Period for Language Acquisition

The CPH posits that language acquisition is most efficient within a specific developmental window, after which it declines. CP effects are observed in both L₁ and L₂ acquisition, suggesting a shared underlying mechanism.

Critical Period for L₁ Acquisition Research in neurolinguistics and cognitive science suggests that there is a biologically determined CP for acquiring an L₁, beyond which full native-like proficiency is unattainable if exposure to language is delayed. Studies on late L₁ learners, such as deaf individuals who acquire sign language after early childhood, indicate severe deficits in grammatical proficiency compared to those exposed to language from birth (Mayberry and Fischer, 1989; Newport, 1990). These findings suggest that neural plasticity, essential for L₁ acquisition, diminishes with age, limiting the ability to develop full linguistic competence. From a theoretical perspective, the existence of the CP for L₁ acquisition is often attributed to biological constraints. Nativist theories propose that L₁ acquisition relies on an innate language faculty that operates most effectively during the CP (Penfield, 1965; Chomsky, 1965; Pinker, 1994). On

the other hand, empiricist perspectives argue that the decline in L₁ learning ability may result from environmental factors, such as a reduced need for language learning mechanisms once fundamental linguistic structures have been internalized (Elman et al., 1996; Seidenberg and Zevin, 2006). Despite extensive research, the precise boundary and mechanisms of the CP for L₁ remain a subject of debate.

Critical Period for L₂ Acquisition CP effects are also observed in L₂ acquisition, where late learners struggle with pronunciation, morphology, and syntax (Johnson and Newport, 1989; Hartshorne et al., 2018). While biological constraints play a role, entrenchment—where prior exposure to L₁ limits flexibility in learning new linguistic structures—is also a factor (Ellis and Lambon Ralph, 2000; Seidenberg and Zevin, 2006). Although the CP for L₂ acquisition is an important topic, this study focuses on the CP for L₁ acquisition, since our goal is to design data-efficient LMs by exploring the mechanisms of CP in L₁ acquisition.

2.2 The Role of Language Models in Acquisition Theories

In recent years, computational models have played a crucial role in elucidating the mechanisms of language acquisition. These models enable controlled investigations of learning mechanisms and environments, which are difficult to achieve with human participants, and they are used to test theoretical claims such as the “poverty of the stimulus” (Clark and Lappin, 2011). For instance, McCoy et al. (2020), Wilcox et al. (2024), and Warstadt et al. (2023) have employed LMs to directly test hypotheses about language acquisition, demonstrating that such models can provide proof-of-concept evidence for *learnability*. These studies have attracted attention as efforts to deepen theoretical discussions on language acquisition through computational modeling, including research on the CP.

Constantinescu et al. (2025) investigated CP phenomena in L₂ acquisition and L₁ attrition,¹ assuming a shared underlying mechanism for CP effects across L₁ and L₂. They simulated L₂ exposure at varying ages to examine how LMs differ from human learners, finding that LMs do not naturally exhibit CP effects. To artificially induce such effects, they employed Elastic Weight Consolidation (Kirk-

¹The phenomenon in which earlier cessation of L₁ exposure increases the likelihood of L₁ forgetting.

patrick et al., 2017), a regularization method for mitigating catastrophic forgetting, thereby mimicking a maturational decline in plasticity. Their findings suggest that CP effects are not an inevitable outcome of statistical learning but may instead involve innate mechanisms.

While this study shares the broader objective of enhancing the cognitive plausibility of LMs as models of human language acquisition, it differs from Constantinescu et al. (2025) in both *focus* and *methodology*. Rather than modeling CP effects through dataset manipulation or post-CP plasticity constraints, this study explicitly addresses the **developmental processes unfolding during the CP itself**. Specifically, we integrate a mechanism to simulate the progressive growth of working memory capacity throughout the CP, a factor considered crucial for L₁ acquisition but previously unmodeled in LM-based research. By incorporating developmental constraints, this study aims to provide a more fine-grained computational model of early L₁ acquisition and its cognitive underpinnings, advancing the developmental plausibility of LMs.

3 Language Model with Developmentally-plausible Working Memory

3.1 Modeling Developmental Trajectory of Human Working Memory

Human working memory undergoes substantial developmental changes, progressing through three distinct stages: early childhood to early school age (2–7 years), middle childhood to early adolescence (8–14 years), and post-adolescence (15 years and older). During early childhood, both information retention capacity and processing ability improve rapidly, reflecting a significant expansion of cognitive resources (Cowan et al., 1999; Gathercole et al., 2004). This rapid growth begins to decelerate during middle childhood and early adolescence as the brain approaches maturation (Luna et al., 2004; Gathercole et al., 2004). By post-adolescence, working memory capacity plateaus, reaching adult-level performance (Sowell et al., 2002; Luna et al., 2004).

Based on these observations, we characterized the growth trajectory of working memory, as illustrated in Figure 1, using an exponential model of the form $y = b - a^x$ ($0 < a < 1$). In this model, b represents the asymptotic upper limit of working memory capacity, corresponding to adult-level per-

formance, while a determines the rate of growth. Specifically, smaller values of a result in steeper early growth, reflecting the rapid cognitive development observed during early childhood, whereas larger values of a indicate a slower rate of change.

This modeling approach is justified for several reasons. First, the horizontal asymptote inherent in the exponential function accurately represents the biological ceiling of adult working memory capacity. Second, the rapid initial increase observed during early childhood is consistent with the steep growth predicted by this exponential form. Finally, alternative models, such as logarithmic or linear growth, fail to account for both the early rapid development and the eventual plateau: logarithmic models imply unbounded growth, while linear models oversimplify the deceleration phase. Thus, the exponential model $y = b - a^x$ offers a concise and biologically plausible representation of the developmental trajectory of human working memory, aligning well with observed patterns and theoretical considerations.

3.2 Integrating Human Working Memory into Language Models

In this study, Attention with Linear Biases (ALiBi) (Press et al., 2022) is employed to model the constraints of human working memory. ALiBi is a method for Transformer (Vaswani et al., 2017) models that does not use positional embeddings but instead applies a distance-dependent linear penalty to attention scores. Specifically, the attention score for an input sequence of length L is calculated as follows:

$$\text{Attention Score} = \text{softmax} \left(q_i K^\top + m \cdot B \right),$$

$$B = \left[-(i-1) \quad -(i-2) \quad \dots \quad 0 \right]. \quad (1)$$

Here, $q_i \in \mathbb{R}^{1 \times d}$, $K \in \mathbb{R}^{L \times d}$, $m \in \mathbb{R}_{[0,1]}$, and $B \in \mathbb{R}^{1 \times L}$ represent the query, the key, a scalar slope specific to each attention head, and a bias matrix encoding the relative distances between queries and keys, respectively, where B_i is defined as the negative absolute difference between the query position i and each key position. The values of m are set geometrically for each head. For example, in an 8-head model, the values of m are assigned as follows: $m = 1, \frac{1}{2}, \frac{1}{4}, \dots, \frac{1}{128}$. The slope m takes values in the range $[0, 1]$, ensuring a consistent interpretation of its influence on attention

scores. By penalizing attention scores for query-key pairs with greater distances, ALiBi introduces a *recency bias* to the model. Originally, ALiBi was proposed to enhance the extrapolation capability of Transformer models. More recently, Clark et al. (2025) has shown that incorporating it into attention score computation during training allows for the estimation of surprisal patterns resembling human reading times. This suggests its potential for modeling human-like memory decay and cognitive limitations.

However, since the slope m in ALiBi is fixed for each attention head, the approach does not inherently reflect the developmental increase in working memory capacity (i.e., reduced decay) over time (Figure 1). Therefore, this study proposes a method, DYNAMICLIMIT-EXP, which replicates the developmental characteristics of working memory during the CP, specifically its exponential growth. This is achieved by exponentially decreasing the slope m in ALiBi as training epochs progress. In this method, the slope m in the ALiBi mechanism is updated at each epoch t as follows:

$$m_t = m_0 \cdot r^t, \quad (2)$$

where m_0 represents the initial slope, $r \in (0, 1)$ is the decay rate, and t denotes the current epoch. In this study, the model’s working memory capacity w_t is formulated as follows:

$$w_t := 1 - m_t. \quad (3)$$

This definition establishes a direct relationship between the dynamically decaying slope m_t and the model’s working memory capacity w_t . As m_t decreases exponentially over time, w_t , representing working memory, grows correspondingly, allowing the model to retain broader contextual information as training progresses. By mimicking this developmentally plausible growth of working memory, the model prioritizes attention to short-range dependencies during the early stages of training, gradually shifting its focus to long-range dependencies as training progresses.

Furthermore, a key distinction between ALiBi and DYNAMICLIMIT-EXP lies in how the slope m is assigned across attention heads. While ALiBi applies a fixed per-head bias, enforcing a predetermined recency bias throughout training, DYNAMICLIMIT-EXP instead shares the slope m across all heads. This ensures that the model maintains a globally coherent bias that evolves dynam-

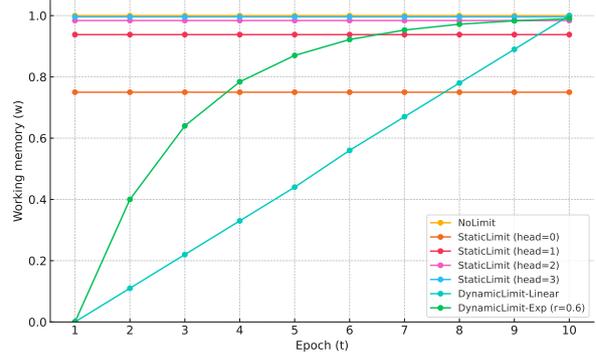


Figure 2: Trajectory of working memory capacity for each model (num. of epochs = 10)

ically over the course of training. In other words, ALiBi imposes a head-specific static recency bias, whereas DYNAMICLIMIT-EXP introduces a dynamically changing proximity bias that governs the entire learning schedule. This shift enables the model to more accurately simulate the adaptive nature of human working memory development, potentially capturing the CP of cognitive maturation.

4 Experiments

This study explores whether LMs trained from scratch can achieve more efficient L_1 acquisition by incorporating the developmental characteristics of human working memory. Specifically, we aim to determine whether this approach can replicate the increased efficiency of L_1 acquisition observed during the CP in L_1 acquisition, focusing on the developmental advantages before the end of this period.

4.1 Configurations

Models We used the transformers (Wolf et al., 2020) implementation of the GPT-2 (Radford et al., 2019) as the base LM. While some studies utilize RoBERTa (Liu et al., 2019) as a base model (Huebner et al., 2021; Warstadt et al., 2023), we selected GPT-2 for two primary reasons: (1) its unidirectional (left-to-right) predictions more effectively capture human working memory constraints, and (2) GPT-based architectures dominate modern LLMs (OpenAI, 2023; Touvron et al., 2023b).

Dataset We used AO-CHILDES (Huebner and Willits, 2021)² as the training dataset, which is derived from the CHILDES dataset (Macwhinney,

²<https://github.com/UIUCLearningLanguageLab/AOCHILDES>

Model	OVERALL	D-N AGR	S-V AGR	ANA. AGR	ARG. STR	BINDING	CASE	ELLIPSIS	FILLER. GAP	IRREGULAR	ISLAND	LOCAL. ATR	QUANTIFIERS	NPI
NOLIMIT	56.5	49.8	49.7	49.9	44.8	61.8	70.8	73.3	72.1	51.7	61.7	47.1	47.9	53.9
STATICLIMIT	56.8	50.2	49.9	49.8	44.4	60.5	70.3	71.4	74.7	52.2	62.9	45.3	52.3	54.4
DYNAMICLIMIT-LINEAR	61.6	51.0	49.6	49.5	64.3	60.3	88.6	47.6	90.8	53.0	57.0	47.9	56.8	84.3
DYNAMICLIMIT-EXP	62.2	50.8	50.0	49.6	67.7	58.7	95.2	43.1	93.6	52.2	53.6	51.3	57.6	85.0

Table 1: Accuracy (%) of models trained on AO-CHILDES dataset. OVERALL represents the macro average of the scores across all grammar items.

Model	OVERALL	D-N AGR	S-V AGR	ANA. AGR	ARG. STR	BINDING	CASE	ELLIPSIS	FILLER. GAP	IRREGULAR	ISLAND	LOCAL. ATR	QUANTIFIERS	NPI
NOLIMIT	54.7	50.3	50.0	47.2	68.4	62.6	73.4	60.8	42.9	53.4	51.1	42.7	41.2	42.6
STATICLIMIT	54.7	50.4	50.0	47.1	73.7	61.2	87.4	57.3	56.1	52.3	53.0	40.8	42.0	38.9
DYNAMICLIMIT-LINEAR	58.6	50.0	50.5	48.4	71.9	58.8	96.9	38.7	82.7	51.6	57.9	59.6	41.5	53.4
DYNAMICLIMIT-EXP	59.1	49.8	50.4	46.0	71.5	59.3	97.7	37.4	86.5	51.1	58.0	60.5	42.2	53.9

Table 2: Accuracy (%) of models trained on Wikipedia dataset. OVERALL represents the macro average of the scores across all grammatical items.

2000) and records CDS from conversations between children and adults. AO-CHILDES contains 5 million words of speech directed at English-speaking children aged 1–6 years and controls for external factors such as age group, speaker variation, and situational context. As a preprocessing step, following Haga et al. (2024), all sentences were converted to lowercase, and sentences shorter than three words were excluded. Since the AO-CHILDES dataset contains only about 5 million words, training a standard GPT-2 model would likely result in overfitting. To mitigate this, we followed existing studies on small language models (SLMs) trained with CDS datasets (Huebner et al., 2021; Haga et al., 2024) and constructed an SLM with 4 layers, 4 attention heads, and 256 embedding dimensions for the base model. Details of the training configuration for the base model are provided in Appendix A.

Furthermore, to determine whether the CP effect stems from exposure to specific linguistic stimuli, such as CDS, or from the model’s cognitive developmental properties independent of input, we conducted a complementary experiment using Wikipedia (written language, adult-oriented) as training data. Following Huebner et al. (2021), 500,000 sentences were randomly sampled from the English Wikipedia corpus. We used the latest version of Wikipedia, as of January 2025,³ and preprocessed it using WikiExtractor.⁴

³<https://dumps.wikimedia.org/enwiki/latest/enwiki-latest-pages-articles.xml.bz2>

⁴<https://github.com/attardi/wikiextractor>

Evaluation We evaluate the grammatical abilities of these models using a developmentally inspired targeted syntactic evaluation benchmark, Zorro (Huebner et al., 2021). Zorro is designed for assessing the syntactic and grammatical knowledge of LMs in child-directed language and consists of 13 mid-level categories and 23 subcategories. Each subcategory contains 2,000 sentence pairs, with one grammatically acceptable and one unacceptable sentence per pair. Below is an example of a minimal pair from the “Subject-verb agreement (S-V AGR)” category.⁵

- (1) a. The **lie** on the foot is flat.
- b. *The **lies** on the foot is flat.

By inputting both the acceptable and unacceptable sentence into the model and calculating the proportion of pairs where the model assigns a higher probability to the acceptable sentence, we obtain the grammaticality judgment score (Accuracy). In this study, we report scores for each mid-level category (henceforth, *grammatical items*) as well as their macro-average.

4.2 Baselines

We prepared the following three baseline models to precisely analyze the learning effects of different working memory limitation strategies:

- **NOLIMIT**: A model with no memory constraints. Working memory remains constant

⁵See Appendix B for the full list of grammatical categories.

Model	OVERALL	D-N AGR	S-V AGR	ANA. AGR	ARG. STR	BINDING	CASE	ELLIPSIS	FILLER. GAP	IRREGULAR	ISLAND	LOCAL. ATR	QUANTIFIERS	NPI
AO-CHILDES														
DYNAMICLIMIT-EXP (↑)	62.2	50.8	50.0	49.6	67.7	58.7	95.2	43.1	93.6	52.2	53.6	51.3	57.6	85.0
DYNAMICLIMIT-EXP (↓)	56.5	49.9	49.7	50.1	44.7	61.9	70.6	73.3	72.0	51.8	61.9	47.0	48.1	54.1
Δ (↑, ↓)	5.7	0.9	0.3	-0.5	23.0	-3.2	24.6	-30.1	21.6	0.4	-8.3	4.4	9.5	30.8
Wikipedia														
DYNAMICLIMIT-EXP (↑)	59.1	49.8	50.4	46.0	71.5	59.3	97.7	37.4	86.5	51.1	58.0	60.5	42.2	53.9
DYNAMICLIMIT-EXP (↓)	52.9	50.4	50.1	47.4	68.7	62.3	74.4	60.2	44.2	53.2	51.7	42.7	40.6	42.2
Δ (↑, ↓)	6.1	-0.6	0.3	-1.4	2.9	-3.0	23.3	-22.8	42.3	-2.2	6.3	17.8	1.7	11.7

Table 3: Performance difference when changing the direction of the cognitive constraints in DYNAMICLIMIT-EXP

from the early stages of training, simulating the mature working memory observed post-adolescence. This configuration is equivalent to a vanilla GPT-2 (Radford et al., 2019).

- **STATICLIMIT**: A model applying standard ALiBi (Press et al., 2022) during attention score calculation, where memory constraints remain fixed throughout training.
- **DYNAMICLIMIT-LINEAR**: A model in which the ALiBi slope m decreases linearly over the course of training.

To ensure a fair comparison between the linear and exponential growth curves of working memory, we controlled the initial and final values of working memory capacity w_t in DYNAMICLIMIT-LINEAR and DYNAMICLIMIT-EXP to be as similar as possible. Specifically, we set the number of training epochs to 10 and configured both models with an initial slope of $m = 1.0$ and a final slope of $m = 0.0$. Figure 2 illustrates the trajectory of working memory capacity for each model. All models were trained using three different seeds, and we report the average results across these runs.

4.3 Results

Developmentally-plausible working memory shapes the CP for L_1 acquisition Table 1 presents the accuracy of each model trained on the AO-CHILDES. Compared to NOLIMIT and STATICLIMIT, which do not account for developmental changes in working memory, DYNAMICLIMIT-LINEAR and DYNAMICLIMIT-EXP, which simulate its gradual growth, achieve significantly higher overall performance. Among them, DYNAMICLIMIT-EXP attains the highest overall accuracy, supporting the effectiveness of a cognitively plausible mechanism. The comparable performance of STATICLIMIT to NOLIMIT suggests

that the gradual introduction of working memory constraints throughout training is crucial, rather than their static application. These results indicate that DYNAMICLIMIT-EXP effectively replicates the CP effect observed in human L_1 acquisition.

The CP depends on the child’s learning algorithm, not the input stimulus Table 2 presents the accuracy of models trained on Wikipedia, showing trends similar to those observed in Table 1, where the models were trained on AO-CHILDES. Specifically, DYNAMICLIMIT-LINEAR and DYNAMICLIMIT-EXP outperform NOLIMIT and STATICLIMIT in overall accuracy, with DYNAMICLIMIT-EXP achieving the highest performance, further supporting the efficacy of incorporating developmental working memory constraints. These findings suggest that the CP effect does not depend solely on exposure to specific linguistic stimuli (e.g., CDS) but rather on the learning algorithm itself, which mirrors human cognitive development.

This result aligns with existing research (Feng et al., 2024), which has reported that child language input is not uniquely valuable for training LMs. This finding suggests that our method is applicable to LLM pretraining, as they typically use non-CDS datasets such as Common Crawl and Wikipedia (Touvron et al., 2023a).

5 Analysis

5.1 Testing the “Less-is-more” Hypothesis with Reversed Cognitive Constraints

A key question arising from the results (§4) is whether DYNAMICLIMIT-EXP’s superior performance stems from the “Less-is-more” hypothesis (Newport, 1990)—i.e., the gradual growth of working memory—or from unintended side effects. In other words, does the gradual *change* in working

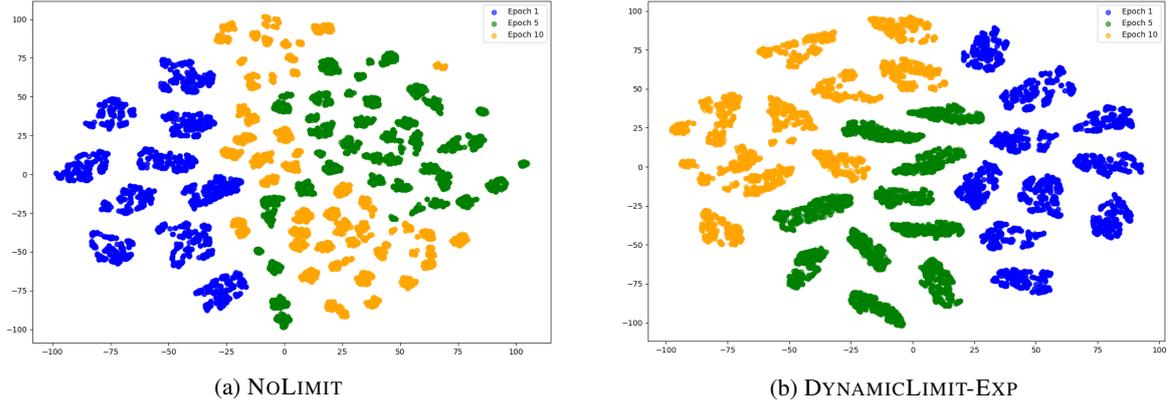


Figure 3: Embedded space at each learning stage for NOLIMIT and DYNAMICLIMIT-EXP (FILLER. GAP)

memory enhance information capacity, dynamically shifting the model’s focus across epochs and ultimately aiding rule generalization? To test this, we introduce a cognitively *implausible* language model, referred to as “DYNAMICLIMIT-EXP (\downarrow)”, which shares the same slope trajectory as our proposed DYNAMICLIMIT-EXP (\uparrow)⁶ but with its direction reversed, such that working memory capacity decreases over time. Specifically, DynamicLimit-Exp (\uparrow) is set to $m_0 = 1.0, r = 0.6$ (the same setting as in §4), while DynamicLimit-Exp (\downarrow) is set to $m_0 = 0.01, r = 1.668$ to achieve a nearly symmetrical curve.⁷

Table 3 provides evidence supporting the Less-is-more hypothesis, as DYNAMICLIMIT-EXP (\uparrow) consistently outperformed the cognitively implausible DYNAMICLIMIT-EXP (\downarrow). The observed performance gap, particularly in grammatical items requiring both local and non-local dependencies (e.g., CASE, ARG. STR, and FILLER-GAP), suggests that the gradual growth of working memory is crucial for grammatical learning and generalization, as it enables the early extraction of basic patterns followed by the progressive acquisition of complex rules. These findings indicate that the superior performance of DYNAMICLIMIT-EXP (\uparrow) is primarily driven by the developmental trajectory of working memory growth rather than unintended side effects of dynamic shifts in memory focus.

Incidentally, from the series of experimental results, along with those in §4 (Table 1 and 2), NO-LIMIT and DYNAMICLIMIT-EXP (\downarrow) consistently outperform DYNAMICLIMIT-EXP (\uparrow) in ELLIPSIS, as exemplified by the following cases:

⁶This section adopts this notation for simplicity.

⁷Since setting the initial slope $m_0 = 0.0$ prevents w_t from being updated in Equation (2), we set it this way for computational reasons.

Epoch	Entropy			Mean Distance		
	1	5	10	1-5	5-10	1-10
NoLimit	5.36	5.17	5.19	91.30	28.50	66.28
DynamicLimit-Exp	5.40	5.30	5.39	69.25	70.63	101.92

Table 4: Embedded space analysis of NOLIMIT and DYNAMICLIMIT-EXP at each stage: distribution diversity and distribution distance.

- (2) a. Mark fixed one **worn** canal, and Roger fixed more. 510
511
- b. *Mark fixed one canal, and Roger fixed more **worn**. 512
513

Since resolving ELLIPSIS involves maintaining long-range dependencies, DYNAMICLIMIT-EXP (\uparrow) may struggle due to its initial memory constraints. This suggests that grammatical items like ELLIPSIS require substantial memory from the early stages of training, and thus, our proposed method may not be optimal for learning such structures. Alternative workarounds, such as dynamically adjusting memory allocation or hybrid approaches, may be necessary to address this limitation. 514
515
516
517
518
519
520
521
522
523
524

5.2 Development of Feature Extraction Capabilities 525 526

Figure 3 visualizes the clustering structure of final-layer embeddings using t-SNE (van der Maaten and Hinton, 2008) for FILLER.GAP, a grammatical items where gradual memory expansion yielded significant performance improvements in both AOCHILDES and Wikipedia datasets, as highlighted in the previous results (§4.3 and §5.1). In NOLIMIT (Figure 3a), the embedding clusters initially expand between Epoch 1 and Epoch 5, but by Epoch 10, 527
528
529
530
531
532
533
534
535

they appear to contract and overlap more, suggesting a stagnation in representation learning. The clusters become less distinguishable, which may indicate a loss of diversity in the learned representations. In contrast, DYNAMICLIMIT-EXP (Figure 3b) maintains a more structured and progressive evolution of embeddings. The clusters remain well-separated throughout training, with clear distinctions between different epochs. This suggests that the model continuously refines its representations without excessive compression, preserving the diversity necessary for robust generalization.

To quantitatively analyze these differences, Table 4 reports key statistical measures, including *entropy* (distribution diversity) and *mean Euclidean distance* (inter-cluster separation).⁸ Regarding **entropy**, NOLIMIT shows a decreasing trend, reflecting reduced distribution diversity and potential over-clustering. In contrast, DYNAMICLIMIT-EXP preserves consistently higher entropy, indicating a balanced representation that avoids excessive compression. For **mean Euclidean distance**, NOLIMIT undergoes substantial change between Epoch 1 and Epoch 5 but stagnates thereafter, suggesting limited refinement. DYNAMICLIMIT-EXP, however, maintains large distances across epochs, reflecting continuous structural reorganization.

The differences also highlight the role of **isotropy**. NOLIMIT exhibits increasing anisotropy, with embedding clusters becoming overly compact by Epoch 10, which may hinder generalization. In contrast, DYNAMICLIMIT-EXP maintains a more isotropic distribution, as indicated by stable entropy, allowing for more flexible and structured representation learning. These findings align with recent work on syntactic smoothing, which suggests that reducing anisotropy enhances the ability to generalize across linguistic contexts (Diehl Martinez et al., 2024). Thus, the increased isotropy observed in DYNAMICLIMIT-EXP provides strong evidence that gradual memory expansion facilitates structured representation learning and syntactic generalization.⁹

5.3 Influence of Input Stimulus Length

We analyze how sentence length affects the performance of NOLIMIT and DYNAMICLIMIT-EXP. To assess their adaptability, we created four Wikipedia-based datasets, each with 500,000 sentences in

⁸The appendix C shows how to calculate each measure.

⁹We also analyzed CASE, which exhibited the same trend as FILLER.GAP (as shown in Appendix D).

Dataset	NOLIMIT	DYNAMICLIMIT-EXP
[5,10]	47.2	46.8
[11,50]	47.0	58.7
[51,100]	40.6	42.5
[101, 150]	37.3	40.8

Table 5: Accuracy in Zorro when the length of the sentence is changed

length ranges: [5,10], [11,50], [51,100], and [101,150].

The results in Table 5 reveal notable differences in model performance. For shorter sentences in the [5,10] range, NOLIMIT achieves slightly higher accuracy compared to DYNAMICLIMIT-EXP. However, in the [11,50] range, DYNAMICLIMIT-EXP significantly outperforms NOLIMIT, achieving 58.7 compared to 47.0. This suggests that DYNAMICLIMIT-EXP excels at handling moderately long sentences, likely due to its ability to dynamically adjust working memory. For longer sentences in the [51,100] and [101,150] ranges, DYNAMICLIMIT-EXP consistently outperforms NOLIMIT.

These findings highlight the benefits of dynamic working memory expansion in facilitating rule generalization and contextual adaptation across diverse sentence lengths. While NOLIMIT exhibits competitive performance on short sentences, its stagnation on longer sentences underscores its limited ability to generalize complex patterns. Conversely, DYNAMICLIMIT-EXP’s consistent performance across varying sentence lengths supports its suitability for grammatical items requiring the processing of both short and long contexts.

6 Conclusion

This study proposed a method for integrating the developmental trajectory of human working memory into the training process of LMs, inspired by the *Less-is-More* hypothesis. The proposed method, DYNAMICLIMIT-EXP, initially restricts working memory and gradually relaxes it exponentially during training. Experiments on both AO-CHILDES and Wikipedia showed that DYNAMICLIMIT-EXP improves grammatical learning efficiency compared to conventional methods without memory constraints or with static memory constraints. These findings suggest not only provide new approaches for developing data-efficient LMs but also offer indirect evidence supporting the CPH in human language acquisition.

626 Limitations

627 **Scalability.** One limitation of this study is the
628 constrained scale of the experimental setup. The
629 primary goal of this study is to computationally
630 replicate the CP in L_1 acquisition, as discussed in
631 cognitive science (Lenneberg, 1967; Fromkin et al.,
632 1974; Curtiss, 1977; Johnson and Newport, 1989).
633 Following previous studies (Huebner et al., 2021;
634 Haga et al., 2024), we designed the experiment
635 to be as ecologically valid as possible by training
636 an SLM using CDS. While this controlled setting
637 allows for a more precise analysis and simulation
638 of the Less-is-More hypothesis, it remains unclear
639 how our findings contribute to the data efficiency
640 of LLMs. The experimental results with Wikipedia
641 (Table 2, 3, 5) provide a promising outlook in this
642 direction, but further investigation with larger mod-
643 els and datasets is necessary to determine the effec-
644 tiveness and limitations of the proposed approach.

645 **Language.** In this experiment, we investigated
646 the replication of the CP effect in L_1 acquisition
647 using English. However, since the CP effect is ob-
648 served across various languages (Patkowski, 1980;
649 Johnson and Newport, 1989), it remains to be tested
650 whether the proposed approach is effective in mul-
651 tilingual environments. To our knowledge, there is
652 currently no targeted syntactic evaluation specifi-
653 cally designed for CDS across different languages,
654 such as Zorro. Zorro was developed based on
655 BLiMP (Warstadt et al., 2020), an adult-oriented
656 targeted syntactic evaluation for English, and re-
657 cent studies have proposed multilingual versions of
658 BLiMP (e.g., JBLiMP (Someya and Oseki, 2023)
659 for Japanese and CLiMP (Xiang et al., 2021) for
660 Chinese). Therefore, developing CDS-specific ver-
661 sions based on these multilingual BLiMPs could
662 help address this limitation.

663 References

664 Noam Chomsky. 1965. *Aspects of the Theory of Syntax*.
665 The MIT Press, Cambridge.

666 Alexander Clark and Shalom Lappin. 2011. *Linguistic*
667 *Nativism and the Poverty of the Stimulus*. Wiley-
668 Blackwell.

669 Christian Clark, Byung-Doh Oh, and William Schuler.
670 2025. *Linear recency bias during training improves*
671 *transformers’ fit to reading times*. In *Proceedings of*
672 *the 31st International Conference on Computational*
673 *Linguistics*, pages 7735–7747, Abu Dhabi, UAE. As-
674 sociation for Computational Linguistics.

Ionut Constantinescu, Tiago Pimentel, Ryan Cotterell, 675
and Alex Warstadt. 2025. *Investigating critical pe-* 676
riod effects in language acquisition through neural 677
language models. *Transactions of the Association for* 678
Computational Linguistics, 13:96–120. 679

Nelson Cowan, Lara Nugent, Emily M. Elliott, Igor 680
Ponomarev, and John Scott Saults. 1999. *The role of* 681
attention in the development of short-term memory: 682
age differences in the verbal span of apprehension. 683
Child development, 70 5:1082–97. 684

S. Curtiss. 1977. *Genie: A Psycholinguistic Study of a* 685
Modern-day "wild Child". Mathematics in Science 686
and Engineering. Academic Press. 687

Richard Diehl Martinez, Zébulon Goriely, Andrew 688
Caines, Paula Buttery, and Lisa Beinborn. 2024. *Mit-* 689
igating frequency bias and anisotropy in language 690
model pre-training with syntactic smoothing. In *Pro-* 691
ceedings of the 2024 Conference on Empirical Meth- 692
ods in Natural Language Processing, pages 5999– 693
6011, Miami, Florida, USA. Association for Compu- 694
tational Linguistics. 695

Andrew W. Ellis and Matthew A. Lambon Ralph. 2000. 696
Age of acquisition effects in adult lexical processing 697
reflect loss of plasticity in maturing systems: Insights 698
from connectionist networks. *Journal of Experimen-* 699
tal Psychology: Learning, Memory, and Cognition, 699
26(5):1103–1123. 700

Jeffrey L. Elman, Elizabeth A. Bates, Mark H. Johnson, 702
Annette Karmiloff-Smith, Domenico Parisi, and Kim 703
Plunkett. 1996. *Rethinking Innateness: A Connec-* 704
tionist Perspective on Development. MIT Press. 705

Steven Y. Feng, Noah Goodman, and Michael Frank. 706
2024. *Is child-directed speech effective training* 707
data for language models? In *Proceedings of the* 708
2024 Conference on Empirical Methods in Natural 709
Language Processing, pages 22055–22071, Miami, 710
Florida, USA. Association for Computational Lin- 711
guistics. 712

Victoria Fromkin, Stephen Krashen, Susan Curtiss, 713
David Rigler, and Marilyn Rigler. 1974. *The de-* 714
velopment of language in genie: a case of language 715
acquisition beyond the "critical period". *Brain* 716
and Language, 1(1):81–107. 717

S. E. Gathercole, S. J. Pickering, B. Ambridge, and 718
H. Wearing. 2004. *The structure of working memory* 719
from 4 to 15 years of age. *Developmental psychol-* 720
ogy, 40(2):177–190. Gathercole, Susan E Pickering, 721
Susan J Ambridge, Benjamin Wearing, Hannah 722
2004/2/26. 723

Akari Haga, Saku Sugawara, Akiyo Fukatsu, Miyu Oba, 724
Hiroki Ouchi, Taro Watanabe, and Yohei Oseki. 2024. 725
Modeling overregularization in children with small 726
language models. In *Findings of the Association* 727
for Computational Linguistics: ACL 2024, pages 728
14532–14550, Bangkok, Thailand. Association for 729
Computational Linguistics. 730

731	Joshua K. Hartshorne, Joshua B. Tenenbaum, and Steven Pinker. 2018. A critical period for second language acquisition: Evidence from 2/3 million English speakers . <i>Cognition</i> , 177:263–277.	786
732		787
733		788
734		
735	Philip A. Huebner, Elior Sulem, Fisher Cynthia, and Dan Roth. 2021. BabyBERTa: Learning more grammar with small-scale child-directed language . In <i>Proceedings of the 25th Conference on Computational Natural Language Learning</i> , pages 624–646. Association for Computational Linguistics.	789
736		790
737		791
738		792
739		793
740		
741	Philip A. Huebner and Jon A. Willits. 2021. Using lexical context to discover the noun category: Younger children have it easier , pages 279–331. <i>Psychology of Learning and Motivation - Advances in Research and Theory</i> . Academic Press Inc.	794
742		795
743		796
744		797
745		
746	Jacqueline S Johnson and Elissa L Newport. 1989. Critical period effects in second language learning: The influence of maturational state on the acquisition of english as a second language . <i>Cognitive Psychology</i> , 21(1):60–99.	798
747		799
748		800
749		801
750		
751	James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. 2017. Overcoming catastrophic forgetting in neural networks . <i>Proceedings of the National Academy of Sciences</i> , 114(13):3521–3526.	802
752		803
753		804
754		805
755		806
756		807
757		
758		
759	E.H. Lenneberg. 1967. <i>Biological Foundations of Language</i> . Wiley.	808
760		809
761	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach . <i>Preprint</i> , arXiv:1907.11692.	810
762		811
763		812
764		
765		
766	Beatriz Luna, Krista E. Garver, Trinity A. Urban, Nicole A. Lazar, and John A. Sweeney. 2004. Maturation of cognitive processes from late childhood to adulthood . <i>Child Development</i> , 75(5):1357–1372.	813
767		814
768		815
769		816
770	Brian Macwhinney. 2000. The childes project: tools for analyzing talk . <i>Child Language Teaching and Therapy</i> , 8.	817
771		
772		
773	Rachel I. Mayberry and Susan D. Fischer. 1989. Looking through phonological shape to lexical meaning: The bottleneck of non-native sign language processing . <i>Memory & Cognition</i> , 17(6):740–754.	818
774		819
775		820
776		821
777	R. Thomas McCoy, Robert Frank, and Tal Linzen. 2020. Does syntax need to grow on trees? sources of hierarchical inductive bias in sequence-to-sequence networks . <i>Transactions of the Association for Computational Linguistics</i> , 8:125–140.	822
778		823
779		824
780		
781		
782	Elissa L. Newport. 1990. Maturational constraints on language learning . <i>Cognitive Science</i> , 14(1).	825
783		826
784	OpenAI. 2023. Gpt-4 technical report . <i>Preprint</i> , arXiv:2303.08774.	827
785		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842

843 Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,
844 Melanie Kambadur, Sharan Narang, Aurelien Ro-
845 driguez, Robert Stojnic, Sergey Edunov, and Thomas
846 Scialom. 2023b. [Llama 2: Open foundation and](#)
847 [fine-tuned chat models](#). *Preprint*, arXiv:2307.09288.

848 Laurens van der Maaten and Geoffrey Hinton. 2008.
849 [Visualizing data using t-sne](#). *Journal of Machine*
850 *Learning Research*, 9(86):2579–2605.

851 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
852 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
853 Kaiser, and Illia Polosukhin. 2017. [Attention is all](#)
854 [you need](#). In *Advances in Neural Information Pro-*
855 *cessing Systems*, volume 30. Curran Associates, Inc.

856 Alex Warstadt, Aaron Mueller, Leshem Choshen, Ethan
857 Wilcox, Chengxu Zhuang, Juan Ciro, Rafael Mos-
858 quera, Bhargavi Paranjabe, Adina Williams, Tal
859 Linzen, and Ryan Cotterell. 2023. [Findings of the](#)
860 [BabyLM challenge: Sample-efficient pretraining on](#)
861 [developmentally plausible corpora](#). In *Proceedings*
862 *of the BabyLM Challenge at the 27th Conference on*
863 *Computational Natural Language Learning*, pages
864 1–34. Association for Computational Linguistics.

865 Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mo-
866 hananey, Wei Peng, Sheng-Fu Wang, and Samuel R.
867 Bowman. 2020. [BLiMP: The benchmark of linguis-](#)
868 [tic minimal pairs for English](#). *Transactions of the*
869 *Association for Computational Linguistics*, 8:377–
870 392.

871 Ethan Gotlieb Wilcox, Richard Futrell, and Roger Levy.
872 2024. [Using computational models to test syntactic](#)
873 [learnability](#). *Linguistic Inquiry*, 55(4):805–848.

874 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien
875 Chaumond, Clement Delangue, Anthony Moi, Pier-
876 ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz,
877 Joe Davison, Sam Shleifer, Patrick von Platen, Clara
878 Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le
879 Scao, Sylvain Gugger, Mariama Drame, Quentin
880 Lhoest, and Alexander M. Rush. 2020. [Transform-](#)
881 [ers: State-of-the-art natural language processing](#). In
882 *Proceedings of the 2020 Conference on Empirical*
883 *Methods in Natural Language Processing: System*
884 *Demonstrations*, pages 38–45, Online. Association
885 for Computational Linguistics.

886 Beilei Xiang, Changbing Yang, Yu Li, Alex Warstadt,
887 and Katharina Kann. 2021. [CLiMP: A benchmark for](#)
888 [Chinese language model evaluation](#). In *Proceedings*
889 *of the 16th Conference of the European Chapter of*
890 *the Association for Computational Linguistics: Main*
891 *Volume*, pages 2784–2790, Online. Association for
892 Computational Linguistics.

A Details of the Training Configuration for the Base Models

Table 6 shows the training settings of the base model. For the experiment, a single NVIDIA RTX A5000 (24GB) GPU was used, and the training time for each run was approximately one hour.

Hyperparameter	Value
Model Architecture	GPT-2
Number of Layers	4
Number of Attention Heads	4
Embedding Dimension	256
Dropout Rate	0.1
Learning Rate (η)	5×10^{-6}
Weight Decay	0.01
Batch Size	512
Gradient Accumulation Steps	2
Total Epochs	20
Maximum Sequence Length	32
Learning Rate Scheduler	Cosine with Restarts
Warm-up Steps	10% of Total Steps
Optimizer	AdamW
Optimizer Parameters	$\beta = (0.9, 0.999)$, $\epsilon = 1e^{-08}$
Tokenizer	Trained on CHILDES
Early Stopping Tolerance	1 Epoch
Evaluation Metric	Perplexity

Table 6: Training Configuration (Hyperparameters) for the GPT-2 Model.

B Details of Grammatical Items in Zorro

Table 8 shows the full list of grammatical categories in Zorro. Examples are taken from Table 5 in the original paper (Huebner et al., 2021).

C Analysis of Distributional Changes in t-SNE Space Across Training Epochs

This section explains in detail the analysis of the entropy and average distance of embeddings projected into the t-SNE space for different learning epochs.

C.1 Entropy Calculation

To quantify the distribution of embeddings, a 2D histogram is constructed using a fixed grid (50×50 bins). The probability distribution P is obtained by normalizing the histogram. The entropy is then computed as:

$$H(P) = - \sum_i P_i \log P_i, \quad (4)$$

where P_i is the probability of each bin. Higher entropy suggests a more uniform distribution, whereas lower entropy indicates clustering.

Epoch	Entropy			Mean Distance		
	1	5	10	1-5	5-10	1-10
NoLimit	5.30	5.23	5.30	75.47	12.26	87.62
DynamicLimit-Exp	5.29	5.30	5.34	59.91	37.68	97.59

Table 7: Embedded space analysis of NOLIMIT and DYNAMICLIMIT-EXP at each stage: cluster expansion, distribution diversity, and distribution distance.

C.2 Mean Distance Between Epochs

To analyze shifts in embedding distributions across epochs, we compute the Euclidean distance between the mean embedding vectors of different epochs:

$$D(X, Y) = \|\mu_X - \mu_Y\|, \quad (5)$$

where μ_X and μ_Y are the mean vectors at different epochs. Larger distances imply greater shifts in the learned representation.

D Development of Feature Extraction Capabilities in CASE

Figure 4 visualizes the clustering structure of final layer embeddings using t-SNE for CASE. The embedding space visualizations reveal distinct patterns between NOLIMIT and DYNAMICLIMIT-EXP across training epochs. In NOLIMIT, the embedding clusters expand between Epoch 1 and Epoch 5 but contract significantly by Epoch 10, suggesting stagnation in representation learning. In contrast, DYNAMICLIMIT-EXP maintains structured evolution throughout training, with well-separated clusters that reflect progressive refinement.

Regarding **entropy**, NOLIMIT shows a slight decrease over time, reflecting reduced distribution diversity as training progresses. In contrast, DYNAMICLIMIT-EXP maintains or slightly increases entropy, suggesting a balanced emphasis on both basic patterns and diverse features, even in later training stages. For **mean Euclidean distances** between clusters, NOLIMIT exhibits large distances between Epoch 1 and Epoch 5 but demonstrates minimal evolution between Epoch 5 and Epoch 10. This stagnation may highlight the model’s failure to effectively generalize new rules. DYNAMICLIMIT-EXP, on the other hand, maintains substantial distances across epochs, indicating continuous embedding evolution and refinement throughout training.

Category	Subcategory	Acceptable Sentence	Unacceptable Sentence
D-N AGR	noun-across_1_adjective noun-between_neighbors	<i>look at this purple thing . this color must be white .</i>	<i>look at this purple things . this colors must be white .</i>
S-V AGR	verb-across_prepositional_phrase verb-across_relative_clause verb-in_question_with_aux verb-in_simple_question	<i>the lie on the foot is flat . the book that i like is poor . where does the horse go ? where is the way ?</i>	<i>the lies on the foot is flat . the books that i like is poor . where does the horses go ? where is the ways ?</i>
ANA.AGR	pronoun_gender	<i>will Mark want himself ?</i>	<i>will Mark want herself ?</i>
ARG.STR	dropped_argument swapped_arguments transitive	<i>give me the poor boat . he made the slave her label . Philip thinks .</i>	<i>the poor boat gives me . the slave made her label he . Philip affected .</i>
BINDING	principle_a	<i>Ben thinks about himself calling this fuel .</i>	<i>Ben thinks about himself called this fuel .</i>
CASE	subjective_pronoun	<i>i brought the wolf my hill .</i>	<i>the wolf brought i my hill .</i>
ELLIPSIS	n_bar	<i>Mark fixed one worn canal and Roger fixed more .</i>	<i>Mark fixed one canal and Roger fixed more worn .</i>
FILLER.GAP	wh_question_object wh_question_subject	<i>Laura married the dinner that the wolf could close . Laura ended the finger that can make boats .</i>	<i>Laura married what the dinner could close the wolf . Laura ended who the finger can make boats .</i>
IRREGULAR	verb	<i>Michael chose the good one some time ago .</i>	<i>Michael chosen the good one some time ago .</i>
ISLAND	adjunct_island coordinate_structure_constraint	<i>who should William have without watching the baby ? who must Philip and the dinosaur turn ?</i>	<i>who should William have the baby without watching ? who must Philip turn and the dinosaur ?</i>
LOCAL.ATR	in_question_with_aux	<i>is the whale getting the person ?</i>	<i>is the whale gets the person ?</i>
NPI	matrix_question only_npi_licensor	<i>does her boat ever play with the growth ? only Mark ever finds some suit .</i>	<i>her boat does ever play with the growth ? even Mark ever finds some suit .</i>
QUANTIFIERS	existential_there superlative	<i>there are many books about soft birds . no pig could stand on top of more than six days .</i>	<i>there are most books about soft birds . no pig could stand on top of at least six days .</i>

Table 8: Explanation of each grammatical category in Zorro.

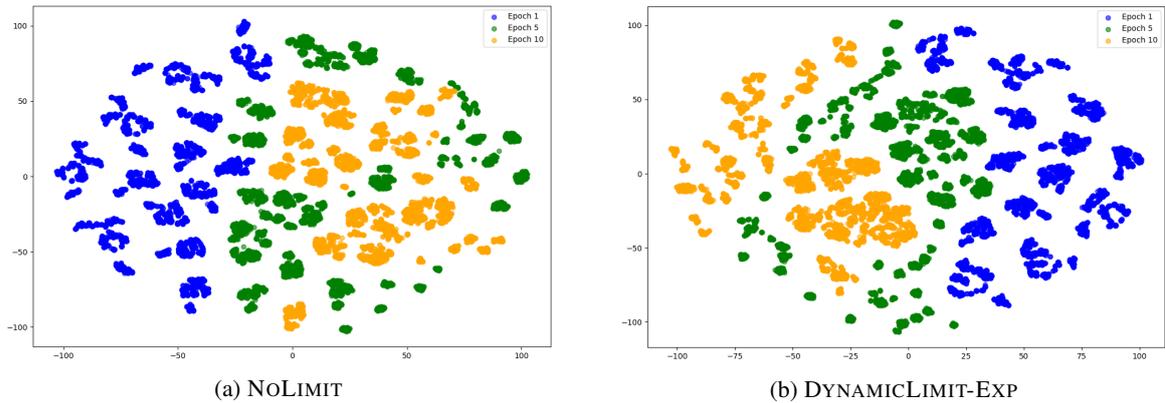


Figure 4: Embedded space at each learning stage for NOLIMIT and DYNAMICLIMIT-EXP (CASE)