

# Unveiling Attractor Cycles in Large Language Models: A Dynamical Systems View of Successive Paraphrasing

Anonymous ACL submission

## Abstract

Dynamical systems theory provides a framework for analyzing iterative processes and evolution over time. Within such systems, repetitive transformations can lead to stable configurations, known as attractors, including fixed points and limit cycles. Applying this perspective to large language models (LLMs), which iteratively map input text to output text, provides a principled approach to characterizing long-term behaviors. Successive paraphrasing serves as a compelling testbed for exploring such dynamics, as paraphrases re-express the same underlying meaning with linguistic variation. Although LLMs are expected to explore a diverse set of paraphrases in the text space, our study reveals that successive paraphrasing converges to stable periodic states, such as 2-period attractor cycles, limiting linguistic diversity. This phenomenon is attributed to the self-reinforcing nature of LLMs, as they iteratively favour and amplify certain textual forms over others. This pattern persists with increasing generation randomness or alternating prompts and LLMs. These findings underscore inherent constraints in LLM generative capability, while offering a novel dynamical systems perspective for studying their expressive potential.

## 1 Introduction

Dynamical systems theory provides a mathematical framework for understanding how iterative processes evolve over time (Sutherland, 1974; Michel et al., 2008). In such systems, repetitive transformation can guide the state of the system toward stable configurations, known as attractors (Milnor, 1985). These attractors can manifest as fixed points, limit cycles, or more complex structures. Applying this perspective to large language models, which iteratively map input text to output text, allows us to characterize their long-term behavioral patterns in a principled manner.

Paraphrase generation can serve as a valuable testbed for exploring these dynamics. Paraphrases

$T_0$ : This way, the filibuster serves as a balance to prevent 51 senators from steamrolling the opposition.

$T_7$ : In this way, the filibuster serves as a safeguard to prevent a majority of 51 senators from dominating the minority.

$T_8$ : In this manner, the filibuster acts as a protection to stop a simple majority of 51 senators from overpowering the minority.

$T_9$ : In this way, the filibuster serves as a safeguard to prevent a simple majority of 51 senators from dominating the minority.

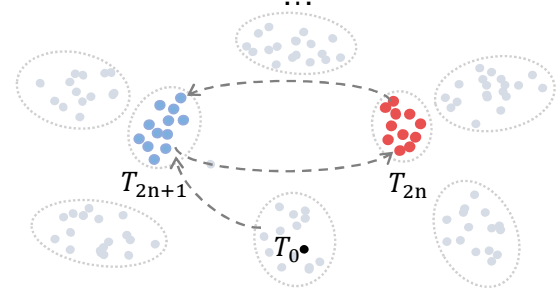


Figure 1: An illustration of successive paraphrasing using GPT-4o-mini: Here,  $T_0$  denotes the original human-written text, while  $T_i$  indicates the  $i$ -th round of paraphrases. The nodes depicted in the lower section represent valid paraphrases for the input sentence, with distance reflecting textual variation. Successive paraphrases generated by LLMs are confined to alternating between two limited clusters, represented as blue and orange nodes.

are re-expressions of the same underlying meaning, differing only in their textual or linguistic form (Bhagat and Hovy, 2013). They serve multiple purposes: improving the readability of text for language learners (Motlagh et al.; Roe and Perkins, 2022a; Kim et al., 2024), enriching datasets in low-resource scenarios (Okur et al., 2022a; Sobrevilla Cabezudo et al., 2024), and enhancing stylistic variation (Krishna et al., 2020). With recent advances in LLMs (Touvron et al., 2023; Taori et al., 2023; Brown et al., 2020; OpenAI, 2023), machine-generated paraphrases can rival or surpass human quality, exhibiting remarkable generalization across diverse domains and text lengths.

While producing a single paraphrase demon-

strates an LLM’s ability to exploit its prior knowledge to create textual variety while preserving semantic equivalence, **successive paraphrasing** pushes this capacity further. Instead of generating just one re-expression, the model recursively paraphrases its own output over multiple rounds (Sadashivan et al., 2023; Tripto et al., 2023). Intuitively, this iterative process is expected to explore an expansive linguistic landscape, generating a rich tapestry of forms. Each subsequent paraphrase, based on previously transformed text, could theoretically diverge into increasingly varied structures—similar to depth-first exploration of the paraphrase search space in contrast to breadth-first approaches like beam search (Holtzman et al., 2020; Huang et al., 2023; Meister et al., 2023).

In practice, however, we find that this expected variety does not materialize. Instead of diverging across a vast combinatorial space, the LLM’s successive paraphrasing converges onto a limited set of recurring solutions, as depicted in Figure 1. When studied through the lens of dynamical systems, these recurring solutions resemble *a stable attractor cycle—a low-order periodic orbit in the space of possible paraphrases* (Milnor, 1985). Rather than continuously discovering new linguistic configurations, the model settles into a pattern where the paraphrased outputs repeat with a fixed period. This phenomenon is subtle: it does not always manifest as explicit repetition but rather as a recurring rotation among a small set of structurally similar forms. Such periodic attractors challenge the intuition that longer or more complex texts should accommodate a broad array of distinct paraphrases. Instead, the LLM gravitates toward a small closed orbit, revealing inherent limitations in its expressive variability.

Specifically, to investigate this attractor-like behavior, we compile a diverse collection of human-written texts (Li et al., 2023) and prompt a range of both open-source and commercial LLMs to perform 15 rounds of successive paraphrasing. Using normalized Levenshtein distance to quantify textual variation, we consistently observe a **2-period cycle**: each new paraphrase resembles the one generated two steps prior. This periodicity proves robust, remaining consistent across multiple models, text lengths, and prompts. We further analyze model perplexity and generation diversity as successive paraphrasing unfolds. The results indicate that, rather than wandering freely in the paraphrase space, LLMs grow increasingly confident in a nar-

row set of solutions, effectively collapsing onto these attractors. Modifying generation hyperparameters or introducing perturbations, such as alternating prompts and models, only subtly disrupts these obstinate attractor cycles. Moreover, this tendency to settle into attractor cycles extends beyond paraphrasing. Any invertible task, i.e., one that allows reconstructions of previous inputs, shows similar behavior, suggesting that such cycles are a general characteristic of LLM iterative behavior.

Finally, we propose a straightforward method to disrupt attractor cycles while maintaining semantic fidelity. By intervening in the iterative process, we can reintroduce meaningful variation and prevent the model from settling into stable yet constrained periodic orbits. In summary, we propose to leverage successive paraphrasing to reveal that LLM outputs, when treated as a dynamical system, tend to converge onto stable attractor cycles rather than exploring open-ended linguistic variety. Understanding these attractors and identifying strategies to escape them is key to unlocking the full expressive potential of LLMs. We will release our data and code after the anonymous period.

## 2 Successive Paraphrasing as System Function

In this section, we briefly introduce the theoretical framework of dynamical systems and applies it to understand the iterative process of successive paraphrasing. By viewing paraphrase generation as the repeated application of a transformation (the LLM’s paraphrasing function), we connect observed phenomena, e.g., periodicity and convergence, to well-studied concepts in systems theory.

### 2.1 Systems Theory Foundations

Systems theory provides a broad mathematical and conceptual framework for analyzing how complex processes evolve over time (Sutherland, 1974). The core idea is modeling the state of a system and its evolution through deterministic or stochastic rules. In continuous or discrete time, systems can exhibit distinct behaviors, ranging from stable equilibria to oscillatory dynamics or even chaotic patterns.

A **dynamical system** is commonly defined as a set of states and a rule describing how those states vary under iteration. When a transformation repeatedly maps an initial state to a new state, one of several outcomes often emerges: *Fixed Points*: States that remain unchanged under the transformation,

representing equilibrium; *Limit Cycles*: Closed loops of states that recur periodically, representing sustained oscillations; *More Complex Attractors*: Patterns to which the system’s trajectories converge, including chaotic attractors.

These attractors shape the long-term behavior of the system. If an initial state lies within the basin of attraction of a limit cycle, for example, the system will converge to that cycle regardless of small perturbations. Identifying such attractors offers valuable insights into the stability and variability of the system’s evolution.

## 2.2 Framing Successive Paraphrasing as a Dynamical System

Successive paraphrasing involves iteratively generating variations of a given text while maintaining semantic equivalence, where each iteration builds upon the previous output. We propose viewing successive paraphrasing as a discrete dynamical system. Let  $\mathcal{T}$  be the space of all possible texts. Consider a large language model that defines a paraphrasing function:  $P : \mathcal{T} \rightarrow \mathcal{T}$ , where  $P(T)$  outputs a paraphrase of the input text  $T$ . Given an initial text  $T_0 \in \mathcal{T}$ , successive paraphrasing generates the sequence  $\{T_n\}_{n=0}^{\infty}$  recursively by:

$$T_{n+1} = P(T_n), \quad n = 0, 1, 2, \dots \quad (1)$$

The set  $\mathcal{P}(T)$  denotes the complete text space for valid paraphrases of  $T$ , which is assumed as a finite space. In theory, the space of potential paraphrases  $\mathcal{P}(T)$  can be vast, especially as text length grows. Each new iteration can potentially explore fresh textual variations, e.g., new syntactic structures, vocabulary choices, and stylistic nuances, while maintaining semantic equivalence. From a systems perspective, if the mapping  $P$  is capable of diversifying output states, one might expect the generated text sequence to spread broadly through the space  $\mathcal{P}(T)$ , never stuck in repetitive patterns, resembling a system without stable attractors. In contrast, if the LLM’s internal biases lead to favouring certain textual forms, the sequence may enter a basin of attraction and converge onto a stable set of states. In other words, rather than exhibiting limitless variety, the system might find itself drawn to limit cycles, i.e., periodic attractors in the paraphrase space.

## 3 Experiment Setup

To systematically investigate this pattern, we first build dedicated testbeds and evaluation criteria.

**Source Data Collection.** We consider English and Chinese paraphrasing in this work. For English paraphrase generation, we collect human-written source documents by sampling instances from the MAGE dataset (Li et al., 2023). Specifically, we uniformly collect 1,000 sentences and 30 paragraphs from each domain in the dataset. This results in a total of 1,000 sentences and 300 paragraphs for subsequent paraphrasing. For Chinese, we source 200 sentences from WMT 2019 (Barrault et al., 2019) and 200 sentences from Wikipedia (Foundation). Detailed data statistics is presented in Appendix A. The main experiments (Section 4) utilize sentence-level paraphrasing datasets, while analytic experiments employ paragraph-level datasets to demonstrate the generality of our findings (Section 5).

**Paraphrase Generation.** For English paraphrasing, we utilize Mistral-7B-Instruct-v0.3 (Jiang et al., 2023), Meta-Llama-3-8B-Instruct (Touvron et al., 2023), Meta-Llama-3-70B-Instruct (Touvron et al., 2023), Qwen2.5-7B-Instruct, Qwen2.5-14B-Instruct, Qwen2.5-72B-Instruct (Yang et al., 2024), GPT-4o-mini and GPT-4o (OpenAI et al., 2024). For Chinese, we use Qwen2.5-7B-Instruct, Qwen2.5-14B-Instruct, Qwen2.5-72B-Instruct, and GPT-4o-mini for paraphrase generation. By default, we set the temperature to 0.6 and p to 0.9 during the decoding process. We sample 10 different paraphrases at each step by setting the number of search beams to 10 and sequentially rephrasing each sample for 15 rounds. We select the candidate with the highest probability for the next paraphrasing iteration.

**Evaluation Metrics.** We use the normalized Levenshtein edit distance function  $d$  to quantify the textual differences between two paraphrases. To provide a more intuitive of the attractor cycle, we propose a metric termed 2-periodicity degree to quantify and study the cyclic pattern in successive paraphrasing. The 2-periodicity degree  $\tau$  is defined as  $\tau = 1 - \frac{1}{M-2} \sum_{i=3}^M d(T_i, T_{i-2})$ , which captures the average textual similarity between the current paraphrase and that from two steps prior.  $M$  denotes the total number of paraphrasing iterations. A higher  $\tau$  indicates stronger periodicity, i.e., similar between two paraphrases. For instance, if successive paraphrases exhibit perfect 2-periodicity such that  $d(T_i, T_{i-2}) = 0$ , then  $\tau = 1$ , indicating that the current paraphrase matches exactly with that from two steps earlier. To evaluate semantic equiv-

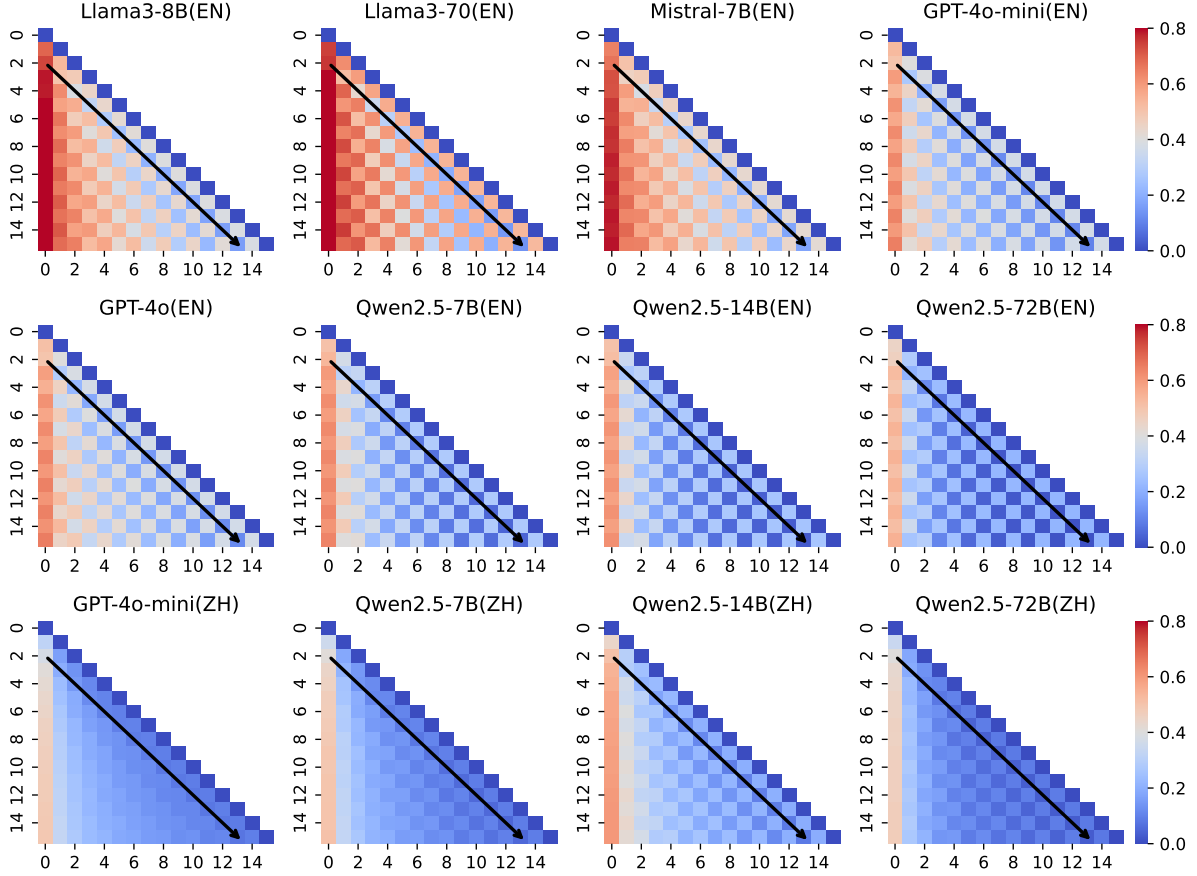


Figure 2: The difference confusion matrix for successive paraphrasing, where EN and ZH denotes English and Chinese sentence-level paraphrase generation accordingly. A darker value indicates a smaller difference between two paraphrases. The blue arrow underlines the differences between  $T_i$  and  $T_{i-2}$ , and averaging these values and subtracting the result from 1 gives our 2-period degree  $\tau$ .

alence, we employ cosine similarity on sentence embeddings<sup>1</sup> (Reimers and Gurevych, 2019).

## 4 Results

Building on the dynamical systems perspective introduced earlier, we now examine the empirical evidence that successive paraphrasing leads LLMs toward stable attractor cycles. We iteratively paraphrase sentences over 15 rounds within the sentence-level dataset and calculate the 2-periodicity degree.

### 4.1 Periodicity

We calculate the textual difference between  $T_i$  and  $T_{i-2}$  for paraphrases at each step. Arranging these differences into a confusion matrix (Figure 2) reveals a pronounced 2-period cycle. For all LLMs, the matrix’s alternating light and dark patterns indicate that paraphrases generated at even iterations

cluster together, and similarly, those at odd iterations form another cluster. This clear partitioning aligns with the behavior of a dynamical system converging onto a 2-period limit cycle—an attractor that draws the iterative process into a stable oscillation between two distinct states.

We also quantify this periodicity across different LLMs, as shown in Table 1. While all models exhibit some degree of 2-periodicity, Qwen2.5-72B shows a particularly strong and consistent cycle in both English and Chinese, whereas Llama3-70B displays relatively weaker periodic behavior. Models with higher periodicity tend to retain more semantic fidelity, suggesting that the recurring attractor states preserve core meaning even as they oscillate between two paraphrastic forms, as shown in Appendix B.

While this periodicity can be viewed as an implicit repetition issue, it differs from explicit repetition of previously seen context. Instead, the model implicitly cycles through a limited set of paraphras-

<sup>1</sup><https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

Mistral-7B	Llama3-8B	Llama3-70B	GPT-4o-mini	GPT-4o	Qwen2.5-7B
0.71	0.72	0.60	0.83	0.81	0.86
Qwen2.5-14B	Qwen2.5-72B	Qwen2.5-7B	Qwen2.5-14B	Qwen2.5-72B	GPT-4o-mini
0.89	0.92	0.70	0.84	0.92	0.88

Table 1: The periodicity degree  $\tau$  of different LLMs. The models represented in blue denotes the English paraphrase generation, while red indicating Chinese paraphrasing.

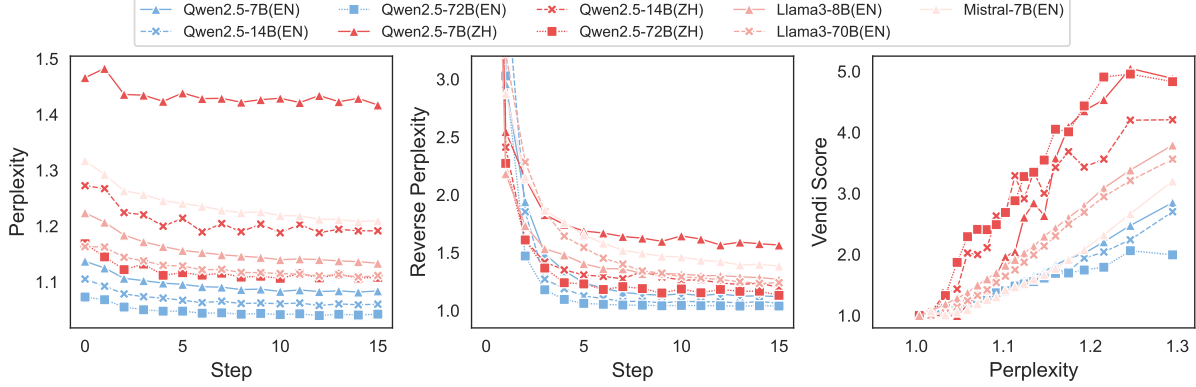


Figure 3: Convergence of perplexity, reverse perplexity, and generation diversity. The left and middle plots show that as the number of steps increases, both perplexity and reverse perplexity decrease steadily until they reach their lower bounds. The right plot shows that generation decreases as perplexity decreases.

tic forms without directly referencing prior iterations. In terms of systems theory, the model’s mapping function  $P$  creates a dynamical environment in which the state space is not fully explored, with the trajectories settling into a 2-period attractor.

## 4.2 Convergence to Stable Attractor

To probe the internal dynamics that lead to these attractor cycles, we explore generation determinism with successive paraphrasing unfolds. We define a **conditioned perplexity**  $\sigma(T_i | T_{i-1})$ , reflecting the model’s confidence in generating  $T_i$  given  $T_{i-1}$ , and a **reverse perplexity**  $\hat{\sigma}(T_i | T_{i+1})$ , indicating how easily  $T_i$  could be reconstructed from  $T_{i+1}$ .

Figure 3 demonstrates that as successive paraphrasing proceeds, both perplexity and reverse perplexity decrease. The forward direction (perplexity) quickly converges to a low boundary, while the reverse direction starts high, indicating that initially it is hard to “go back” from  $T_{i+1}$  to  $T_i$ . However, it drops fast as paraphrasing proceeds and aligns with the forward perplexity. Finally, the system evolves towards a state where generating  $T_{i+1}$  from  $T_i$  is nearly as deterministic and predictable as reconstructing  $T_i$  from  $T_{i+1}$ . This symmetry resembles a **stable attractor** in a dynamical system, where bidirectional predictability indicates that the sys-

tem has “locked in” to a limit cycle.

We further quantify generation diversity by sampling multiple paraphrases at each iteration and computing the Vendi score (Friedman and Dieng, 2022). As shown in Figure 3, a low perplexity indicates a low generation diversity. A Vendi score of one indicates that all paraphrases in the beam are identical to each other. As both forward and reverse perplexity decreases, the model consistently produces similar paraphrases, leaving minimal room for alternative textual trajectories. From a systems viewpoint, the collapse into low perplexity and low diversity states corresponds to the model settling into the basin of attraction of a periodic orbit. Once inside the basin, the model’s generative behavior becomes nearly deterministic, causing the output sequence to cycle predictably.

The notion of invertibility, where each paraphrase can be treated as a paraphrase of its own paraphrase, further explains the robustness of periodicity. Invertibility places constraints on the mapping function  $P$ , effectively enabling a bidirectional relationship between states which encourages stable cycles. This insight suggests that tasks with similar invertible properties, e.g., translation, can also display limit cycle behavior, a hypothesis we will explore in Section 5.1.

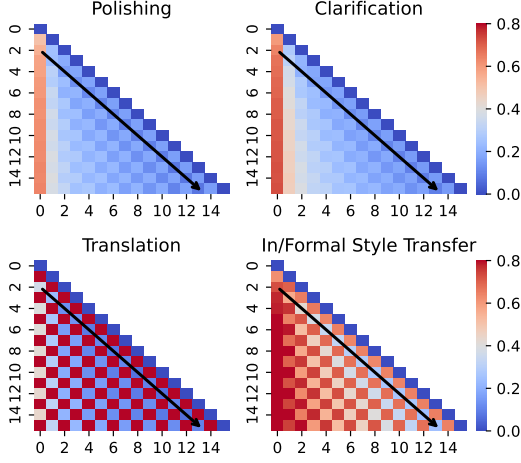


Figure 4: The difference confusion matrix for four tasks beyond paraphrasing. Note that in translations, the difference between texts in two different languages is set to one.

## 5 Analysis

In this section, we perform analytical experiments on paragraph-level paraphrase datasets to generalize our findings to longer texts. We first demonstrate the extension of our findings to other task formats (Section 5.1). Then we go through a set of methods to try to escape from the attractor cycles in the remaining subsections.

### 5.1 Beyond Paraphrase Generation

Our earlier results indicate that successive paraphrasing leads LLMs to settle into periodic attractors—specifically, 2-period limit cycles. According to the systems-theoretic perspective, such cycles should arise whenever the transformation is invertible, enabling a bidirectional mapping that makes prior states easily reproducible. To test this, we examine four additional invertible tasks at the paragraph level: polishing (Pol.), clarification (Clar.), informal-to-formal style transfer (I/F.), and forward/backward translation (Trans.). These tasks are defined in Appendix C.1.

Figure 4 shows that even for these varied tasks, LLMs repeatedly converge to stable states, exhibiting pronounced 2-periodicity. Table 2 shows the degree of 2-periodicity across these tasks, with values ranging from 0.65 to 0.87. This finding reinforces the idea that invertibility fosters the emergence of limit cycles, as the model iterates the transformation and settles into an attractor. While paraphrasing is our primary lens, these findings confirm that stable attractor cycles are a broader characteristic

of LLM behavior in iterative, invertible mappings.

Tasks	Para.	Clar.	Pol.	I/F.	Trans.
$\tau$	0.80	0.83	0.86	0.65	0.87

Table 2: Impact of perturbations on periodicity compared to the original during paraphrasing.

### 5.2 Alternating Models and Prompts

One intuitive approach to escape an attractor is to introduce perturbations in the transformation itself. We attempt this by varying both models and prompts during successive paraphrasing. For **prompt variation**, we design four different paraphrasing prompts (refer to Appendix C.2) and randomly select one at each iteration. Despite regularly switching prompts, the 2-period cycle persists, as shown in Figure 5.

Similarly, we introduce **model variation** by alternating among GPT-4o-mini, GPT-4o, Llama3-8B, and Qwen2.5-7B during successive paraphrasing. Although each model brings its own stylistic biases, the fundamental attractor cycle remains intact. Interestingly, perplexity computed by a single model (e.g., Llama3-8B) on paraphrases generated by other models still decreases over iterations in Figure 6. This suggests that the attractor states are not confined to a single model’s parameter space. Instead, they reflect a more general statistical optimum that multiple LLMs gravitate toward.

From a systems perspective, this findings suggest that randomizing the transformation function  $P$  does not inherently break the attractor. The system remains in a basin of attraction shared across these varied modeling conditions, implying that the stable cycle is a robust property of the iterative transformation rather than a quirk of any particular prompt or model.

### 5.3 Increasing Generation Randomness

Another strategy is introducing more stochasticity in the generation process by increasing the generation temperature. Higher temperatures expand the immediate token selection space, potentially allowing trajectories to wander away from the attractor. However, as shown in Figure 7, while higher temperatures do increase the difference between successive paraphrases, the system still exhibits a 2-period cycle. Further increases in temperature lead only to nonsensical outputs. This outcome aligns with dynamical systems theory: a small increase in

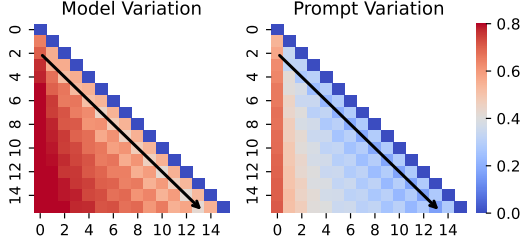


Figure 5: The difference confusion matrices for model variation and prompt variation.

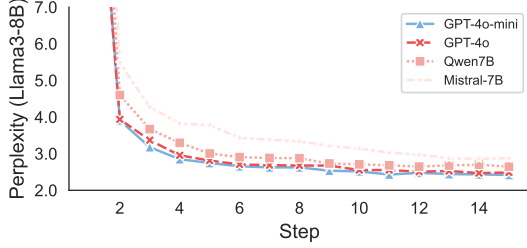


Figure 6: The perplexity of  $T_i$  conditioned on  $T_{i-1}$  calculated by Llama3-8B. Both  $T_i$  and  $T_{i-1}$  are generated by other LLMs.

stochasticity may create local perturbations, but if the basin of attraction is strong, the system remains near the limit cycle. Excessive stochastic forcing can push the system out of meaningful regions of state space entirely, leading to “chaotic” or nonsensical behavior, rather than discovering a new stable attractor with richer linguistic diversity.

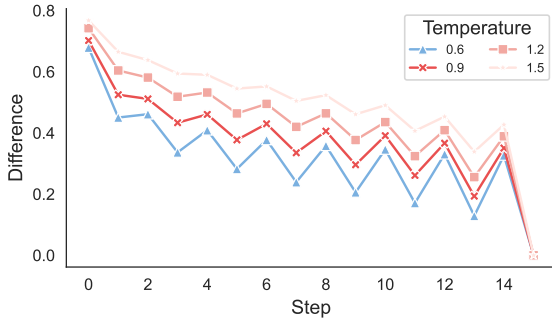


Figure 7: The difference between  $T_{15}$  and  $T_i$  generated by GPT-4o-mini. By increasing the temperature, randomness is amplified, causing the differences to grow as well.

## 5.4 Incorporating Local Perturbations

We introduce local perturbations to mitigate the attractor cycle pattern. At the end of each iteration, we edit 5% of the text by introducing perturbations using three methods: synonym replacement (S.R.), word swapping (W.S.), and random insertion or

deletion (I./D.). As shown in Table 3, among these interventions, synonym replacement barely affects periodicity, suggesting that minor lexical changes do not move the system out of the attractor’s basin. It indicates that except during the first paraphrasing, LLMs primarily perform synonym replacements for words or phrases, as shown in Figure 1. Word swapping, however, causes more significant disruption, lowering periodicity more effectively. From a dynamical standpoint, large structural perturbations are needed to shift the system’s state out of a stable cycle. Local lexical tweaks do not suffice because the attractor’s pull is strong and preserved at a deeper structural level.

w/o Perturb.	S.R.	W.S.	I./D.
0.77	0.73	0.62	0.66

Table 3: Impact of different types of perturbations on 2-periodicity degrees  $\tau$ , compared to the original text during paraphrasing.

## 5.5 Paraphrasing with History Paraphrases

We consider a scenario where the transformation  $\hat{P}$  depends on both  $T_i$  and  $T_{i-1}$ . This added historical context can alter the equilibrium states. In a scenario where we paraphrase  $T_i$  based on the reference  $T_{i-1}$ , it is essential that  $T_{i+1}$  differs from both  $T_i$  and  $T_{i-1}$ . This function can be expressed as:  $T_{i+1} = \hat{P}(T_i, T_{i-1})$ . In this context,  $P_{i-1}$  emerges as a strong candidate for paraphrasing  $P(T_{i+1}, T_i)$ , as it aligns with the distribution of LLMs while maintaining difference from  $\hat{P}(T_{i+1}, T_i)$ , satisfying the task requirement. As a result, this more complex cycle still represents a stable attractor, albeit of higher order, as shown in Figure 8.

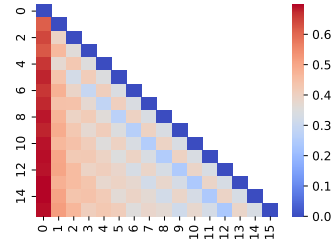


Figure 8: When adding historical paraphrases, LLMs exhibit 3-periodicity in the paraphrasing task.

## 5.6 Sample Selection Strategies

Finally, we investigate methods to steer the system away from stable attractors at the least cost

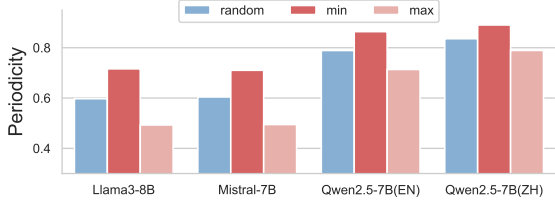


Figure 9: The periodicity of three strategies using different LLMs.

of generation quality. Given the correlation between periodicity and perplexity, it is intuitive to mitigate this issue by increasing perplexity while maintaining generation quality. To achieve this, we can randomly sample multiple paraphrases at each iteration and select the one based on perplexity. We design three types of strategies: selecting the paraphrase with the maximum or minimum perplexity or randomly choosing one at each iteration. Figures 9 illustrate that selecting a higher perplexity can reduce periodicity. However, such diversity comes at the cost of semantic equivalence (Appendix C.4). Considering both periodicity and meaning preservation, we recommend the random strategy, which effectively reduces periodicity while incurring minimal information loss compared to selecting the option with the lowest perplexity.

## 6 Related Work

**Paraphrase Generation.** Paraphrase generation has long been a significant focus in NLP research, with numerous studies dedicated to enhancing the quality of generated paraphrases. (Li et al., 2018; Roy and Grangier, 2019; Lewis et al., 2020; Lin et al., 2021b; Hosking et al., 2022; Xie et al., 2022). Some studies also explore methods to control paraphrase generation by focusing on aspects such as syntactic and lexical diversity (Li et al., 2019; Goyal and Durrett, 2020; Huang and Chang, 2021; Bandel et al., 2022; Krishna et al., 2023; Yang et al., 2022). Others investigate the application of paraphrase generation as a data augmentation technique to enhance model performance (Jolly et al., 2020; Bencke and Moreira, 2024; Okur et al., 2022b).

Recently, advancements in LLMs have enabled LLM-based paraphrasing tools to generate stable, high-quality responses, making them widely used for refining materials like news articles, academic papers, and speeches (Witteveen and Andrews, 2019; Roe and Perkins, 2022b; Rani et al., 2023). However, their work primarily discusses single-

step paraphrasing. In contrast, another line of work involves LLMs iteratively rephrasing their own outputs over multiple iterations. Sadasivan et al. (2023) explores how repeated rephrasing can help evade AI text detectors, while Tripto et al. (2023) and Huang et al. (2024) discuss the implications for authorship after a document has undergone multiple rounds of paraphrasing. Our research differs from those work. We investigate the inherent characteristics of paraphrasing when extended over multiple iterations.

**Self-Reinforcement in LLMs.** Repetition, defined as the occurrence of repetitive text in natural language generation, has been widely explored in research community (Holtzman et al., 2020; Welleck et al., 2020; Lin et al., 2021a; See et al., 2017; Liu and Lapata, 2019; Fu et al., 2021). Xu et al. (2022) introduce the concept of self-reinforcement to elucidate this phenomenon, demonstrating that LLMs exhibit a tendency to repeat preceding sentences and reinforce this behavior during generation. Yan et al. (2024) further explore the relationship between the self-reinforcement effect and the in-context learning capabilities of LLMs. Our research aligns with theirs in examining the self-reinforcement patterns of LLMs. However, we specifically concentrate on typical behaviors observed in successive paraphrasing tasks, mirroring LLMs’ limitations in the exploration of text space.

## 7 Conclusion

We reframed successive paraphrasing as a discrete dynamical system, offering a principled explanation for the emergence of stable periodic attractors in LLM-generated text. Our empirical findings revealed that instead of producing an expanding array of diverse paraphrases, LLMs rapidly settled into low-order limit cycles. These attractor states persisted even when we vary models, prompts, generation temperatures, and local kick perturbations, indicating that they stem from a fundamental property of the system rather than superficial repetition or particular model idiosyncrasies. Viewing iterative text generation through the lens of systems theory helps clarify why certain interventions fail to break these cycles and how others can weaken the attractor’s pull. Ultimately, recognizing and addressing these stable attractor cycles is crucial for unlocking more expressive and flexible language generation for large language models.

## Limitations

While this study provides valuable insights into successive paraphrasing, several limitations should be acknowledged. First, the paraphrasing is based on simple prompts, which may limit the generalizability of the findings to more complex or specific prompts. Second, although we have examined this phenomenon in the currently prevalent LLMs, other LLMs may not exhibit the same behavior. Finally, while we present the convergence of reverse perplexity in this work, the underlying reasons for this behavior still require further investigation.

## Ethic Considerations

We uphold the Code of Ethics and ensure that no private or non-public information is used in this work. We comply with the terms set by companies offering commercial LLM APIs and extend our gratitude to all collaborators for their invaluable support in utilizing these APIs.

## References

- Elron Bandel, Ranit Aharonov, Michal Shmueli-Scheuer, Ilya Shnayderman, Noam Slonim, and Liat Ein-Dor. 2022. [Quality controlled paraphrase generation](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 596–609, Dublin, Ireland. Association for Computational Linguistics.
- Loïc Barrault, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, Mathias Müller, Santanu Pal, Matt Post, and Marcos Zampieri. 2019. [Findings of the 2019 conference on machine translation \(WMT19\)](#). In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 1–61, Florence, Italy. Association for Computational Linguistics.
- Luciana Bencke and Viviane Pereira Moreira. 2024. Data augmentation strategies to improve text classification: a use case in smart cities. *Language Resources and Evaluation*, 58(2):659–694.
- Rahul Bhagat and Eduard Hovy. 2013. [Squibs: What is a paraphrase?](#) *Computational Linguistics*, 39(3):463–472.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Wikimedia Foundation. [Wikimedia downloads](#).
- Dan Friedman and Adji Bousso Dieng. 2022. The vendi score: A diversity evaluation metric for machine learning. *arXiv preprint arXiv:2210.02410*.
- Zihao Fu, Wai Lam, Anthony Man-Cho So, and Bei Shi. 2021. A theoretical analysis of the repetition problem in text generation. In *Thirty-Fifth AAAI Conference on Artificial Intelligence*.
- Tanya Goyal and Greg Durrett. 2020. [Neural syntactic preordering for controlled paraphrase generation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 238–252, Online. Association for Computational Linguistics.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. [The curious case of neural text de-generation](#). In *International Conference on Learning Representations*.
- Tom Hosking, Hao Tang, and Mirella Lapata. 2022. [Hierarchical sketch induction for paraphrase generation](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2489–2501, Dublin, Ireland. Association for Computational Linguistics.
- Baixiang Huang, Canyu Chen, and Kai Shu. 2024. [Authorship Attribution in the Era of LLMs: Problems, Methodologies, and Challenges](#). *arXiv preprint. ArXiv:2408.08946*.
- Kuan-Hao Huang and Kai-Wei Chang. 2021. [Generating Syntactically Controlled Paraphrases without Using Annotated Parallel Pairs](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1022–1033, Online. Association for Computational Linguistics.
- Tenghao Huang, Ehsan Qasemi, Bangzheng Li, He Wang, Faeze Brahman, Muhao Chen, and Snigdha Chaturvedi. 2023. [Affective and dynamic beam search for story generation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 11792–11806, Singapore. Association for Computational Linguistics.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. 2023. [Mistral 7b](#). *Preprint, arXiv:2310.06825*.

- Shailza Jolly, Tobias Falke, Caglar Tirkaz, and Daniil Sorokin. 2020. [Data-efficient paraphrase generation to bootstrap intent classification and slot labeling for new features in task-oriented dialog systems](#). In *Proceedings of the 28th International Conference on Computational Linguistics: Industry Track*, pages 10–20, Online. International Committee on Computational Linguistics.
- Yewon Kim, Thanh-Long V. Le, Donghwi Kim, Mina Lee, and Sung-Ju Lee. 2024. [How Non-native English Speakers Use, Assess, and Select AI-Generated Paraphrases with Information Aids](#). *arXiv preprint*, ArXiv:2405.07475 version: 1.
- Kalpesh Krishna, Yixiao Song, Marzena Karpinska, John Wieting, and Mohit Iyyer. 2023. [Paraphrasing evades detectors of ai-generated text, but retrieval is an effective defense](#). *ArXiv preprint*, abs/2303.13408.
- Kalpesh Krishna, John Wieting, and Mohit Iyyer. 2020. [Reformulating unsupervised style transfer as paraphrase generation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 737–762, Online. Association for Computational Linguistics.
- Mike Lewis, Marjan Ghazvininejad, Gargi Ghosh, Armen Aghajanyan, Sida Wang, and Luke Zettlemoyer. 2020. [Pre-training via Paraphrasing](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 18470–18481. Curran Associates, Inc.
- Yafu Li, Qintong Li, Leyang Cui, Wei Bi, Zhilin Wang, Longyue Wang, Linyi Yang, Shuming Shi, and Yue Zhang. 2023. [Mage: Machine-generated Text Detection in the Wild](#). *arXiv*.
- Zichao Li, Xin Jiang, Lifeng Shang, and Hang Li. 2018. [Paraphrase generation with deep reinforcement learning](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3865–3878, Brussels, Belgium. Association for Computational Linguistics.
- Zichao Li, Xin Jiang, Lifeng Shang, and Qun Liu. 2019. [Decomposable neural paraphrase generation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3403–3414, Florence, Italy. Association for Computational Linguistics.
- Xiang Lin, Simeng Han, and Shafiq Joty. 2021a. [Straight to the gradient: Learning to use novel tokens for neural text generation](#). *Preprint*, arXiv:2106.07207.
- Zhe Lin, Yitao Cai, and Xiaojun Wan. 2021b. [Towards document-level paraphrase generation with sentence rewriting and reordering](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 1033–1044, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Yang Liu and Mirella Lapata. 2019. [Text Summarization with Pretrained Encoders](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3730–3740, Hong Kong, China. Association for Computational Linguistics.
- Clara Meister, Tiago Pimentel, Gian Wiher, and Ryan Cotterell. 2023. [Locally typical sampling](#). *Transactions of the Association for Computational Linguistics*, 11:102–121.
- Anthony N Michel, Ling Hou, and Derong Liu. 2008. [Stability of dynamical systems](#). *Birkhäuser, Boston, MA*.
- John Milnor. 1985. [On the concept of attractor](#). *Communications in Mathematical Physics*, 99:177–195.
- Negin Yazdani Motlagh, Matin Khajavi, Abbas Sharifi, and Mohsen Ahmadi. [The Impact of Artificial Intelligence on the Evolution of Digital Education: A Comparative Study of OpenAI Text Generation Tools including ChatGPT, Bing Chat, Bard, and Ernie](#).
- Eda Okur, Saurav Sahay, and Lama Nachman. 2022a. [Data Augmentation with Paraphrase Generation and Entity Extraction for Multimodal Dialogue System](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4114–4125, Marseille, France. European Language Resources Association.
- Eda Okur, Saurav Sahay, and Lama Nachman. 2022b. [Data augmentation with paraphrase generation and entity extraction for multimodal dialogue system](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4114–4125, Marseille, France. European Language Resources Association.
- OpenAI. 2023. [GPT-4 technical report](#). *ArXiv preprint*, abs/2303.08774.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,

778	Simón Posada Fishman, Juston Forte, Isabella Ful-		
779	ford, Leo Gao, Elie Georges, Christian Gibson, Vik		
780	Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-		
781	Lopes, Jonathan Gordon, Morgan Grafstein, Scott		
782	Gray, Ryan Greene, Joshua Gross, Shixiang Shane		
783	Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris,		
784	Yuchen He, Mike Heaton, Johannes Heidecke, Chris		
785	Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele,		
786	Brandon Houghton, Kenny Hsu, Shengli Hu, Xin		
787	Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,		
788	Joanne Jang, Angela Jiang, Roger Jiang, Haozhun		
789	Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee-		
790	woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-		
791	mali, Ingmar Kanitscheider, Nitish Shirish Keskar,		
792	Tabarak Khan, Logan Kilpatrick, Jong Wook Kim,		
793	Christina Kim, Yongjik Kim, Jan Hendrik Kirch-		
794	ner, Jamie Kiros, Matt Knight, Daniel Kokotajlo,		
795	Łukasz Kondraciuk, Andrew Kondrich, Aris Kon-		
796	stantinidis, Kyle Kosic, Gretchen Krueger, Vishal		
797	Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan		
798	Leike, Jade Leung, Daniel Levy, Chak Ming Li,		
799	Rachel Lim, Molly Lin, Stephanie Lin, Mateusz		
800	Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue,		
801	Anna Makanju, Kim Malfacini, Sam Manning, Todor		
802	Markov, Yaniv Markovski, Bianca Martin, Katie		
803	Mayer, Andrew Mayne, Bob McGrew, Scott Mayer		
804	McKinney, Christine McLeavey, Paul McMillan,		
805	Jake McNeil, David Medina, Aalok Mehta, Jacob		
806	Menick, Luke Metz, Andrey Mishchenko, Pamela		
807	Mishkin, Vinnie Monaco, Evan Morikawa, Daniel		
808	Mossing, Tong Mu, Mira Murati, Oleg Murk, David		
809	Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak,		
810	Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh,		
811	Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex		
812	Paino, Joe Palermo, Ashley Pantuliano, Giambat-		
813	tista Parascandolo, Joel Parish, Emy Parparita, Alex		
814	Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-		
815	man, Filipe de Avila Belbute Peres, Michael Petrov,		
816	Henrique Ponde de Oliveira Pinto, Michael, Poko-		
817	rny, Michelle Pokrass, Vitchyr H. Pong, Tolly Pow-		
818	ell, Alethea Power, Boris Power, Elizabeth Proehl,		
819	Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh,		
820	Cameron Raymond, Francis Real, Kendra Rimbach,		
821	Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-		
822	der, Mario Saltarelli, Ted Sanders, Shibani Santurkar,		
823	Girish Sastry, Heather Schmidt, David Schnurr, John		
824	Schulman, Daniel Selsam, Kyla Sheppard, Toki		
825	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav		
826	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,		
827	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin		
828	Sokolowsky, Yang Song, Natalie Staudacher, Fe-		
829	lipe Petroski Such, Natalie Summers, Ilya Sutskever,		
830	Jie Tang, Nikolas Tezak, Madeleine B. Thompson,		
831	Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,		
832	Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-		
833	lipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,		
834	Chelsea Voss, Carroll Wainwright, Justin Jay Wang,		
835	Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,		
836	CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-		
837	ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner,		
838	Clemens Winter, Samuel Wolrich, Hannah Wong,		
839	Lauren Workman, Sherwin Wu, Jeff Wu, Michael		
840	Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim-		
841	ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong		
	Zhang, Marvin Zhang, Shengjia Zhao, Tianhao	842	
	Zheng, Juntang Zhuang, William Zhuk, and Bar-	843	
	ret Zoph. 2024. <a href="#">Gpt-4 technical report</a> . <i>Preprint</i> ,	844	
	arXiv:2303.08774.	845	
	Anku Rani, S.M Towhidul Islam Tonmoy, Dwip Dalal,	846	
	Shreya Gautam, Megha Chakraborty, Aman Chadha,	847	
	Amit Sheth, and Amitava Das. 2023. <a href="#">FACTIFY-</a>	848	
	<a href="#">5WQA: 5W aspect-based fact verification through</a>	849	
	<a href="#">question answering</a> . In <i>Proceedings of the 61st An-</i>	850	
	<i>annual Meeting of the Association for Computational</i>	851	
	<i>Linguistics (Volume 1: Long Papers)</i> , pages 10421–	852	
	10440, Toronto, Canada. Association for Computa-	853	
	tional Linguistics.	854	
	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-bert:</a>	855	
	<a href="#">Sentence embeddings using siamese bert-networks</a> .	856	
	<i>Preprint</i> , arXiv:1908.10084.	857	
	Jasper Roe and Mike Perkins. 2022a. <a href="#">What are Auto-</a>	858	
	<a href="#">mated Paraphrasing Tools and how do we address</a>	859	
	<a href="#">them? A review of a growing threat to academic</a>	860	
	<a href="#">integrity</a> . <i>International Journal for Educational In-</i>	861	
	<i>tegrity</i> , 18(1):15.	862	
	Jasper Roe and Mike Perkins. 2022b. What are auto-	863	
	mated paraphrasing tools and how do we address	864	
	them? a review of a growing threat to academic	865	
	integrity. <i>International Journal for Educational In-</i>	866	
	<i>tegrity</i> , 18(1):15.	867	
	Aurko Roy and David Grangier. 2019. <a href="#">Unsupervised</a>	868	
	<a href="#">paraphrasing without translation</a> . In <i>Proceedings of</i>	869	
	<i>the 57th Annual Meeting of the Association for Com-</i>	870	
	<i>putational Linguistics</i> , pages 6033–6039, Florence,	871	
	Italy. Association for Computational Linguistics.	872	
	Vinu Sankar Sadasivan, Aounon Kumar, Sriram Bala-	873	
	subramanian, Wenxiao Wang, and Soheil Feizi. 2023.	874	
	<a href="#">Can ai-generated text be reliably detected?</a> <i>ArXiv</i>	875	
	<i>preprint</i> , abs/2303.11156.	876	
	Abigail See, Peter J. Liu, and Christopher D. Man-	877	
	ning. 2017. <a href="#">Get To The Point: Summarization</a>	878	
	<a href="#">with Pointer-Generator Networks</a> . In <i>Proceedings</i>	879	
	<i>of the 55th Annual Meeting of the Association for</i>	880	
	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	881	
	pages 1073–1083, Vancouver, Canada. Association	882	
	for Computational Linguistics.	883	
	Marco Antonio Sobrevilla Cabezudo, Marcio Lima Ina-	884	
	cio, and Thiago Alexandre Salgueiro Pardo. 2024.	885	
	<a href="#">Investigating Paraphrase Generation as a Data Aug-</a>	886	
	<a href="#">mentation Strategy for Low-Resource AMR-to-Text</a>	887	
	<a href="#">Generation</a> . In <i>Proceedings of the 17th International</i>	888	
	<i>Natural Language Generation Conference</i> , pages	889	
	663–675, Tokyo, Japan. Association for Computa-	890	
	tional Linguistics.	891	
	John W. Sutherland. 1974. <a href="#">Review of "general systems</a>	892	
	<a href="#">theory: Foundations, developments, applications"</a>	893	
	<a href="#">by ludwig von bertalanffy</a> . <i>IEEE Trans. Syst. Man</i>	894	
	<i>Cybern.</i> , 4(6):592.	895	
	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann	896	
	Dubois, Xuechen Li, Carlos Guestrin, Percy Liang,	897	

and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca).

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. **Llama 2: Open foundation and fine-tuned chat models**. *Preprint*, arXiv:2307.09288.

Nafis Irtiza Tripto, Saranya Venkatraman, Dominik Macko, Róbert Móro, Ivan Srba, Adaku Uchendu, Thai Le, and Dongwon Lee. 2023. **A ship of theseus: Curious cases of paraphrasing in llm-generated texts**. *ArXiv preprint*, abs/2311.08374.

Sean Welleck, Ilia Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. 2020. **Neural text generation with unlikelihood training**. In *International Conference on Learning Representations*.

Sam Witteveen and Martin Andrews. 2019. **Paraphrasing with large language models**. In *Proceedings of the 3rd Workshop on Neural Generation and Translation*, pages 215–220, Hong Kong. Association for Computational Linguistics.

Xuhang Xie, Xuesong Lu, and Bei Chen. 2022. **Multi-task learning for paraphrase generation with keyword and part-of-speech reconstruction**. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1234–1243, Dublin, Ireland. Association for Computational Linguistics.

Jin Xu, Xiaojiang Liu, Jianhao Yan, Deng Cai, Huayang Li, and Jian Li. 2022. **Learning to Break the Loop: Analyzing and Mitigating Repetitions for Neural Text Generation**. *arXiv preprint*. ArXiv:2206.02369.

Jianhao Yan, Jin Xu, Chiyu Song, Chenming Wu, Yafu Li, and Yue Zhang. 2024. **Understanding In-Context Learning from Repetitions**. *arXiv preprint*. ArXiv:2310.00297.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan

Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.

Kexin Yang, Dayiheng Liu, Wenqiang Lei, Baosong Yang, Haibo Zhang, Xue Zhao, Wenqing Yao, and Boxing Chen. 2022. **GCPG: A general framework for controllable paraphrase generation**. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 4035–4047, Dublin, Ireland. Association for Computational Linguistics.

## A Data Statistics

We provide source information of our data in table 4 and statistic information of data length in 10.

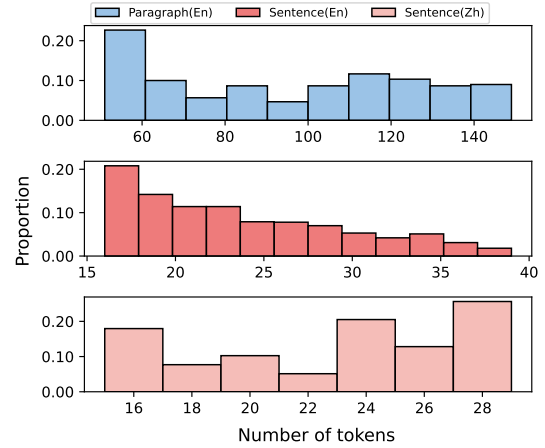


Figure 10: Statistical patterns of data length distribution.

## B Change in similarity

We measure the change in similarity between  $T_i$  and  $T_0$  across successive paraphrasing steps. The results are presented in Figure 11. As the number of paraphrasing steps increases, most LLMs maintain the similarity between paraphrases and their corresponding original texts, with the exception of an initial drop in similarity. Meanwhile, it also exhibits aslight 2-periodicity in similarity. By combining Figure 11 and Table 2, we found that models with higher periodicity also exhibit higher similarity.

Dataset	TLDR	SQuAD	ROCT	Yelp	ELI5	Sci_Gen
Sentence/Paragraph	100/30	100/30	100/30	100/30	100/30	100/30
Dataset	XSum	CMV	HSWAG	WP	Wiki	WMT
Sentence/Paragraph	100/30	100/30	100/30	100/30	200/0	200/0

Table 4: Dataset Setup: Datasets marked in red indicate Chinese datasets, while others represent English datasets. The value indicates the number of extracted samples. For example, we extract 100 sentences and 30 paragraphs from the TLDR dataset.

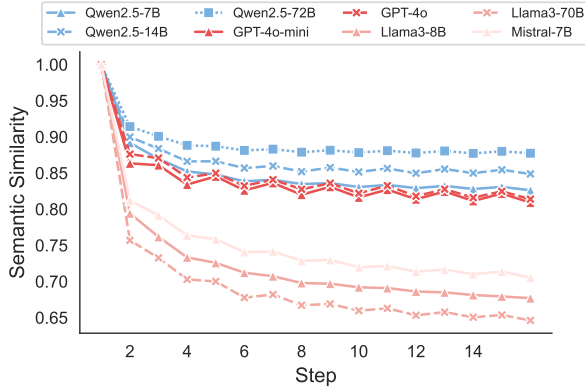


Figure 11: Similarity changes during successive paraphrasing. Qwen2.5-72B is the best at preserving meaning, while all other LLMs experience slight degradation in similarity, except during the first paraphrasing step.

## C Generalization

### C.1 Task Extensions

We propose four additional tasks beyond paraphrasing: polishing (**Pol.**), clarification (**Clar.**), informal-to-formal style transfer (**I/F.**), and forward/backward translation (**Trans.**). The detailed prompts for these tasks are listed in Table 5. We perform these tasks on our paragraph dataset, calculate the textual difference of the paraphrase at each iteration with the initial text, and plot the results in Figure 12. As the number of paraphrasing steps increases, the difference between  $T_i$  and  $T_{i-2}$  decreases. After 7 steps, there is little difference between  $T_i$  and  $T_{i-2}$ .

### C.2 Model and Prompt variation

We continue to modify the models and prompts during paraphrasing. The chosen model set includes GPT-4o-mini, GPT-4o, Qwen2.5-7B, and Llama3-8B. Four variations of the paraphrasing prompts are provided in Table 6.

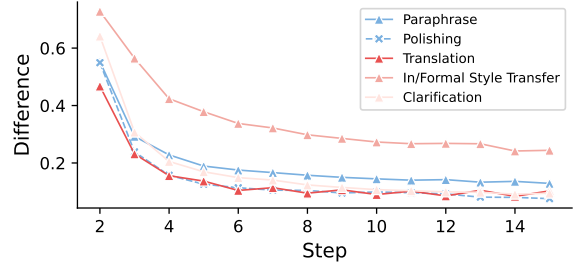


Figure 12: The trend in normalized edit distance between  $T_i$  and  $T_{i-2}$  across various tasks during the repetition process using GPT-4o-mini.

<b>Pol.</b>	Please polish the following text: {text}
<b>Clar.</b>	Please rewrite the following text in a way that is simpler and easier to understand, using clear language and shorter sentences without losing the original meaning: {text}
<b>I/F.</b>	Transform the following text into an informal style: {text} / Rewrite the following text in a formal style: {text}
<b>Trans.</b>	Please translate the following English text into Chinese: {text} / Please translate the following English text into Chinese: {text}

Table 5: Four types of prompts for extension tasks. The last two tasks involve switching between different languages and styles, separated by a semicolon.

<b>A:</b>	Please paraphrase the following text: {text}
<b>B:</b>	Please rephrase the text below: {text}
<b>C:</b>	Please rewrite the following text: {text}
<b>D:</b>	Please polish the text below: {text}

Table 6: Four variations of paraphrasing prompts. In the prompt variation experiments, a prompt is randomly selected at each step to perform the paraphrasing.

### C.3 Increasing Randomness

We measure the impact of increasing randomness on periodicity by adjusting the generation temperature. We select four temperature values: 0.6, 0.9, 1.2, and 1.5. The results are shown in Fig-

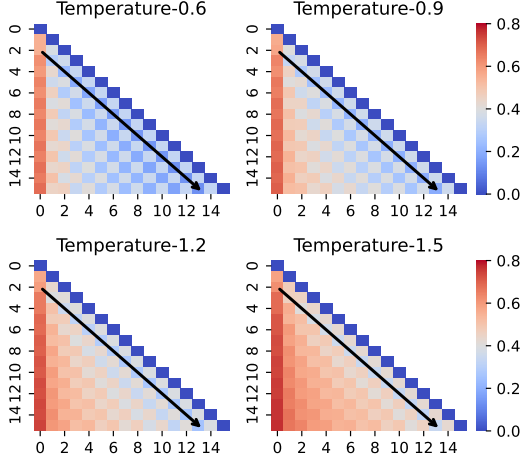


Figure 13: The difference confusion matrix for successive paraphrasing at different temperature settings, conducted by GPT-4o-mini.

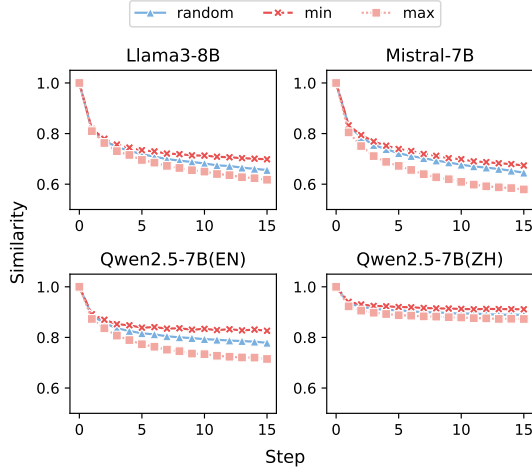


Figure 14: The similarity between paraphrases and the original texts increases during the paraphrasing process.

effectively alleviating periodicity.

## D Case Study

We present part of the first case from our paragraph dataset for successive paraphrasing in Figure 7.

ure 13. Although the temperature increases to a very high level, the 2-periodicity still persists. Further increasing temperature will cause nonsense responses.

## C.4 Sample Selection Strategies

We propose three strategies for successive paraphrasing and evaluate them across different LLMs. To assess the impact of these strategies on meaning preservation, we measure the similarity between the paraphrases  $T_i$  and their corresponding original texts  $T_0$  and demonstrate the result in Figure 14. By combining Figure 14 and Figure 9, we suggest that the random strategy preserves meaning significantly better than the max strategy, while also

Sample of successive paraphrasing	
$T_0$	Australia has a strong tradition of an honest days pay for an honest days work. We have a strong union presence which leads to high wages backed by a strong resource economy. We also have a very multicultural society with many immigrants choosing to settle in Australia this leads to diverse shopping areas where goods originally made all over the world are made right here on our shores. The main reason Australia is the best country in the world is our heavily ingrained concept of a "fair go" meaning you give anybody a chance to prove themselves no matter who they are. This is not even mentioning our spectacular landscapes and beaches. All in all Australia is the best country in the world.
$T_1$	Australia has a robust tradition of fair compensation for hard work. The strong presence of unions contributes to high wages, supported by a solid resource-based economy. Additionally, our multicultural society, enriched by many immigrants, creates diverse shopping areas where products from around the globe are available right here in Australia. A key reason why Australia stands out as the best country in the world is our deep-rooted belief in a "fair go," which emphasizes giving everyone the opportunity to demonstrate their abilities, regardless of their background. Not to mention, our stunning landscapes and beautiful beaches add to our appeal. Overall, Australia truly is the best country in the world.
$T_2$	Australia has a strong tradition of ensuring fair pay for hard work. The significant influence of unions leads to high wages, bolstered by a stable resource-driven economy. Furthermore, our multicultural society, enriched by a diverse population of immigrants, fosters vibrant shopping districts where global products are readily available. One of the main reasons Australia is considered the best country in the world is our fundamental belief in a "fair go," which highlights the importance of providing everyone with the chance to showcase their talents, regardless of their background. Additionally, our breathtaking landscapes and beautiful beaches enhance our country's appeal. Overall, Australia is truly the best place in the world.
$T_3$	Australia has a strong history of promoting fair compensation for hard work. The powerful presence of unions contributes to high wages, supported by a stable economy driven by natural resources. In addition, our multicultural society, enriched by a diverse immigrant population, creates lively shopping areas filled with international products. A key reason why Australia is often regarded as the best country in the world is our core belief in a "fair go," which emphasizes the importance of giving everyone the opportunity to demonstrate their abilities, regardless of their background. Moreover, our stunning landscapes and beautiful beaches add to the allure of our nation. Overall, Australia is undoubtedly one of the finest places in the world.
$T_4$	Australia has a rich tradition of advocating for fair pay for hard work. The strong influence of unions helps maintain high wages, bolstered by a robust economy fueled by natural resources. Additionally, our multicultural society, enhanced by a varied immigrant community, fosters vibrant shopping districts filled with international goods. A significant reason Australia is often viewed as one of the best countries globally is our fundamental belief in a "fair go," which highlights the necessity of providing everyone with the chance to showcase their talents, regardless of their background. Furthermore, our breathtaking landscapes and beautiful beaches contribute to the appeal of our nation. Overall, Australia is truly one of the most exceptional places in the world.
$T_5$	Australia has a rich history of promoting fair compensation for hard work. The strong presence of unions plays a crucial role in ensuring high wages, supported by a thriving economy driven by natural resources. Moreover, our diverse society, enriched by a wide-ranging immigrant population, creates lively shopping areas filled with products from around the globe. A key reason Australia is often regarded as one of the best countries in the world is our core belief in a "fair go," which emphasizes the importance of giving everyone the opportunity to demonstrate their abilities, regardless of their background. Additionally, our stunning landscapes and beautiful beaches enhance the allure of our nation. In summary, Australia is truly one of the most remarkable places in the world.
$T_6$	Australia has a notable history of advocating for fair pay for hard work. The significant influence of unions is essential in maintaining high wages, bolstered by a robust economy fueled by natural resources. Furthermore, our diverse society, enriched by a broad immigrant community, fosters vibrant shopping districts filled with international products. A major reason Australia is frequently seen as one of the best countries globally is our fundamental belief in a "fair go," which underscores the importance of providing everyone the chance to showcase their talents, irrespective of their background. Additionally, our breathtaking landscapes and beautiful beaches add to the appeal of our nation. In conclusion, Australia is indeed one of the most extraordinary places in the world.
$T_7$	Australia has a rich history of promoting fair compensation for hard work. The strong presence of unions plays a crucial role in sustaining high wages, supported by a thriving economy that benefits from its natural resources. Moreover, our multicultural society, enhanced by a diverse immigrant population, creates lively shopping areas filled with global products. A key reason Australia is often regarded as one of the best countries in the world is our core belief in a "fair go," which highlights the importance of giving everyone the opportunity to demonstrate their abilities, regardless of their background. Additionally, our stunning landscapes and picturesque beaches contribute to the allure of our nation. In summary, Australia is truly one of the most remarkable places on Earth.

Table 7: Case illustration of the first sample in our paragraph dataset. Due to the limited content, we just display  $T_0 \dots T_7$ .