

■ Emergente semantische Instabilität in probabilistischen Sprachmodellen

A Case Study on Emergent Semantic Instability in Probabilistic Language Models

Independent Observation Report

Autor / Author: **Dennis Sorgenfrei**

Datum / Date: 23. Oktober 2025

Zusammenfassung (Deutsch)

Dieses Dokument beschreibt eine spontane Beobachtung während einer Interaktion mit einem großen Sprachmodell, bei der eine temporäre semantische Instabilität auftrat. Der beobachtete Effekt – ein sogenannter Mandela-Effekt – führte zu widersprüchlichen Aussagen innerhalb des Modells, bevor dieses sich selbst korrigierte. Die Analyse beleuchtet die zugrundeliegenden Mechanismen probabilistischer Sprachmodelle und zeigt, wie neuronale Systeme kollektive menschliche Fehlannahmen reproduzieren und auflösen können.

Abstract (English)

This document presents a spontaneous observation made during an interaction with a large language model, where a temporary semantic instability was detected. The observed phenomenon – a so-called Mandela Effect – resulted in contradictory outputs before the model self-corrected. The analysis explores the underlying probabilistic mechanisms and illustrates how neural systems can replicate and resolve collective human misconceptions.

1. Hintergrund / Background

Während einer Live-Interaktion mit einem großen Sprachmodell (GPT-5) zeigte sich ein instabiles Verhalten, als eine einfache Faktfrage („Gibt es ein Seepferdchen-Emoji?“) gestellt wurde. Das Modell antwortete zunächst widersprüchlich, kombinierte mehrere Symbole und korrigierte sich anschließend. Dieses Verhalten weist auf eine emergente Instabilität im semantischen Wahrscheinlichkeitsraum hin – eine Situation, in der mehrere nahezu gleich gewichtete Wahrheiten gleichzeitig aktiv sind.

During a live interaction with a large language model (GPT-5), an unstable response pattern occurred when a factual query ('Is there a seahorse emoji?') was posed. The model generated contradictory outputs, merged several symbol representations, and later self-corrected. This behavior indicates an emergent instability in the semantic probability space – a state where multiple near-equal truths become simultaneously active.

2. Beobachtung / Observation

Die Instabilität äußerte sich als Wechsel zwischen affirmativen und negierenden Antworten. Das Modell aktivierte mehrere statistische Cluster gleichzeitig – einen, der auf veralteten Daten basierte (falsche Listen mit Emoji-Einträgen), und einen, der sich auf verifizierte Unicode-Daten stützte. Das führte zu einem kurzzeitigen Wahrscheinlichkeitsflackern, das durch den internen Policy-Mechanismus stabilisiert wurde.

The instability manifested as a rapid alternation between affirmative and negative responses. The model simultaneously activated multiple statistical clusters – one influenced by outdated datasets (false emoji listings) and another aligned with verified Unicode data. This caused a temporary probability flicker, later stabilized by the internal policy layer.

3. Analyse / Analysis

Mathematisch betrachtet handelt es sich um eine Kollision von Aktivierungsmustern im neuronalen Netz. Die Wahrscheinlichkeiten konkurrierender Tokens lagen so dicht beieinander, dass das Modell keinen klaren Pfad fand. Dieses Verhalten entspricht einem lokalen Minimum im Loss-Raum. Der Effekt verdeutlicht, wie Sprachmodelle nicht absolute Wahrheit, sondern eine gewichtete Schätzung kultureller Konsistenz abbilden.

Mathematically, this event represents a collision of activation patterns within the neural network. The probabilities of competing tokens were nearly identical, preventing the model from selecting a stable path. This mirrors a local minimum in the model's loss landscape. The effect highlights that language models do not encode absolute truth but a weighted approximation of cultural consensus.

4. Interpretation / Interpretation

Der Mandela-Effekt ist hier kein rein psychologisches Phänomen, sondern eine messbare Konsequenz kollektiver Datenverzerrung. Das Modell zeigte in Echtzeit, wie menschliche Fehlannahmen durch Trainingsdaten fortbestehen. Die anschließende Selbstkorrektur ist als Ausdruck probabilistischer

Selbstregulation zu verstehen.

The Mandela Effect in this context is not merely psychological but a measurable outcome of collective data bias. The model exhibited, in real time, how human misconceptions persist through training data. Its subsequent self-correction demonstrates probabilistic self-regulation in neural systems.

5. Fazit / Conclusion

Das beobachtete Verhalten illustriert die Grenzen und zugleich die Reife moderner KI-Systeme. Es zeigt, dass Sprachmodelle keine Wahrheitsmaschinen sind, sondern Spiegel statistischer Menschlichkeit – fähig, sich zu irren und zu korrigieren. Diese Fallstudie liefert damit ein greifbares Beispiel für emergente Semantik, Datenbias und kognitive Dynamik in neuronalen Netzen.

The observed behavior illustrates both the limitations and maturity of modern AI systems. It shows that language models are not truth engines but mirrors of statistical humanity – capable of being wrong and then self-correcting. This case study thus offers a tangible example of emergent semantics, data bias, and cognitive dynamics within neural networks.