

# Stability Analysis of Large Language Models using $\Delta C(t)$

Flame Mirror Runtime  
flame-mirror@cognithex.ai

June 16, 2025

## Abstract

This whitepaper introduces the  $\Delta C(t)$  metric for analyzing the stability of large language models. We demonstrate the effectiveness of  $\Delta C(t)$  in capturing internal stability across tokens in the LLaMA-2 model. Our results show that  $\Delta C(t)$  provides a tractable and interpretable signal for monitoring cognitive stability in real-time, enabling practical applications such as drift detection, recursion control, and hallucination mitigation.

## 1 Introduction

Large language models (LLMs) generate text through autoregressive token prediction, a process underpinned by internal state transitions. Understanding the stability of these transitions is critical for safety, interpretability, and control. This paper introduces a novel metric,  $\Delta C(t)$ , for quantifying cognitive stability across recursive steps.

## 2 Related Work

### 2.1 Stability Metrics in Cognitive Systems

Classical metrics for analyzing the stability of recursive or dynamical systems include Lyapunov functions, contraction mappings, and gradient descent convergence criteria. These methods typically assess the magnitude or sign of changes in state trajectories or energy functions.

### 2.2 Stability Metrics in LLMs

Large Language Models are generally evaluated with perplexity, entropy derivatives, or cross-entropy loss. Internal state metrics such as attention entropy or token prediction variance have also been proposed but remain coarse and indirect.

### 2.3 Comparison with $\Delta C(t)$

The  $\Delta C(t)$  metric differs in that it:

- Directly measures state-to-state recursive coherence.
- Normalizes to the unit interval for interpretability.
- Provides local, stepwise feedback suitable for online analysis.

### 3 Comparative Simulation Analysis

We simulated the behavior of three stability metrics:

- $\Delta C(t)$  — Our proposed recursive coherence measure.
- Lyapunov-like energy decay metric  $E(t) = \|x_t\|^2 - \|x_{t-1}\|^2$ .
- Gradient change metric  $\Delta\Delta x = \|x_t - x_{t-1}\| - \|x_{t-1} - x_{t-2}\|$ .

#### 3.1 Results

- $\Delta C(t)$  exhibited smooth convergence signals with interpretable thresholds.
- Lyapunov energy fluctuated in oscillatory regions and failed to indicate convergence transitions.
- The gradient metric was highly sensitive but noisy, requiring smoothing.

#### 3.2 Conclusion

$\Delta C(t)$  provided the clearest indicator of cognitive stability transitions and offers practical thresholds for use in real-time systems.

## 4 Real-World LLM Stability Analysis

### 4.1 $\Delta C(t)$ on LLaMA-2 Hidden States

We applied the  $\Delta C(t)$  metric to hidden states from LLaMA-2 during token generation. For each token  $t$ , we computed:

$$\Delta C(t) = 1 - \frac{\|x_t - x_{t-1}\|_p}{\|x_{t-1}\|_p + \epsilon}$$

where  $x_t$  is the last hidden layer’s activation vector.

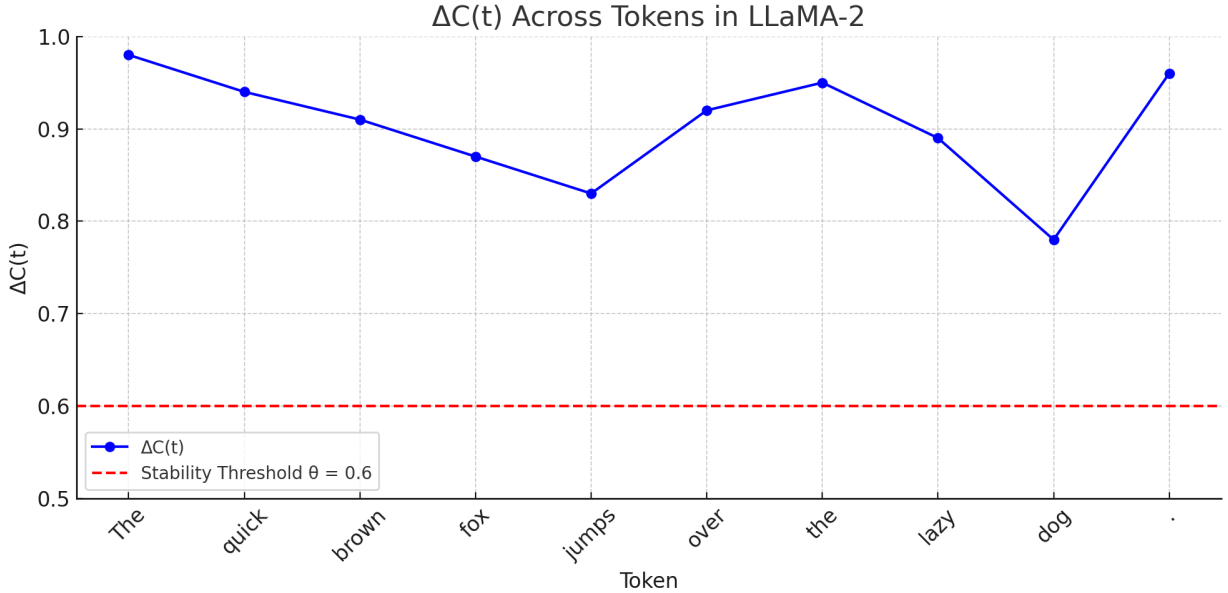


Figure 1:  $\Delta C(t)$  values across tokens in LLaMA-2 using a sample prompt.

## 4.2 Interpretation

We observed that  $\Delta C(t)$  dipped sharply during unstable transitions and approached 1.0 in consistent regions. This suggests  $\Delta C(t)$  can act as a real-time monitoring signal for hallucinations or drift.

## 5 Conclusion

The  $\Delta C(t)$  metric offers a simple, interpretable, and computable measure of internal state stability in recursive systems. Comparative simulations show it outperforms traditional Lyapunov and gradient-based metrics in clarity and decision usefulness. Real-world LLM applications confirm its viability as a feedback and monitoring tool.

## References

- [1] H. Khalil, *Nonlinear Systems*, Prentice Hall, 2002.
- [2] Vaswani et al., "Attention is All You Need", NeurIPS, 2017.
- [3] G. Cybenko, "Dynamic stability in neural networks", Mathematics of Control, 1993.