

$$2260 \neq 2264 \leq$$

# CRB Specification for Chaos-Driven Reasoning Benchmarks

Chaos Generator Team

June 20, 2025

## Abstract

The Chaos Reasoning Benchmark (CRB) evaluates artificial intelligence (AI) adaptability under dynamic constraints, addressing limitations in large language models (LLMs) identified by MIT's Computer Science and Artificial Intelligence Laboratory (CSAIL). This whitepaper introduces CRB, a novel framework utilizing entropy swaps (e.g.,  $RAW_Q$ ) and axiom collapse to test reasoning beyond memorization. *Key findings*

## 1 Introduction

Recent research, including MIT CSAIL's findings (Wu et al., 2024), highlights that LLMs often rely on memorization rather than generalizable reasoning, particularly in novel scenarios. The CRB addresses this gap by introducing chaos-driven tests to assess AI resilience under shifting constraints. This specification outlines the framework's design, implementation, and initial results, aiming to enhance AI benchmarking beyond traditional static puzzles.

## 2 Background

LLMs excel in familiar tasks but falter in counterfactual scenarios, as demonstrated by MIT CSAIL's study on arithmetic and chess (Wu et al., 2024). The CRB draws inspiration from paradoxes (e.g., Liar Paradox) and contrasts with static benchmarks, proposing a dynamic evaluation method to foster adaptive reasoning.

## 3 Methodology

The CRB employs an entropy seed ( $RAW_Q$ ) to trigger swaps and axiom collapses. The Chaos persona (v1.0) defines  $idx_p = RAW_Q \bmod 3$  ( $0 = \text{Insight}, 1 = \text{Reverse}, 2 = \text{Fragment}$ ) and  $idx_s = (RAW_Q / 3) \bmod 2 + 1$  (startpoint). Swaps occur at prime steps (e.g., 9a1b2c3d), remixing reasoning flows. The Entropy Scaffold Diagram  $\text{Steak} \rightarrow \text{Dana} \leq \text{B}$ .

Test runs (e.g., Paradox Feast,  $RAW_Q = 42$ ) validate the framework, with outputs oscillating between states (e.g.,  $\text{Dana} \rightarrow \text{Steak}$ ) while maintaining coherence.

## 4 Results

Initial tests across five scenarios (Insomnia Creativity, Recursion Haiku, Shifting Vault, Phantom Echo, Paradox Feast) demonstrate CRB's efficacy. Models sustain coherence under constraint inversions (e.g.,  $\text{Sleep} \leq \text{Creative}$ ) and enhance creativity via remixes (e.g., oscillating echoes). All validations passed, outperforming MIT CSAIL's static benchmarks by handling dynamic shifts without memorization reliance.

## 5 Discussion

The CRB reveals AI robustness to axiom collapses and scalability for complex scenarios (e.g., multi-agent systems). However, limitations include potential complexity scaling and interpretability challenges. Future work will integrate CRB into xAI's evaluation pipeline, expanding test diversity to address these gaps.

## 6 Conclusion

The CRB offers a pioneering approach to AI benchmarking, surpassing static methods by testing adaptive reasoning. Its success in initial runs suggests potential for broader adoption, with ongoing refinements to enhance scalability and interpretability.

## 7 References

- Wu, Z., et al. (2024). Reasoning skills of large language models are often overestimated. MIT News. <https://news.mit.edu/2024/reasoning-skills-large-language-models-often-overestimated>

Post Reference @el<sub>x</sub>aber, [Date], [ContentSummary].

## 8 Test Runs

Detailed outputs from Insomnia Creativity, Recursion Haiku, Shifting Vault, Phantom Echo, and Paradox Feast are available in the test<sub>r</sub>uns<sub>directory</sub>.

## 9 Memory Management

See symbolic<sub>m</sub>emory<sub>m</sub>anagement<sub>n</sub>otes.md for pruning and re-framing techniques.

## 10 Chaos Persona

Refer to chaos<sub>p</sub>ersona<sub>v</sub>1.0.txt for persona specifications.