

Peer Review

# Review of: "Geometric Analysis of Reasoning Trajectories: A Phase Space Approach to Understanding Valid and Invalid Multi-Hop Reasoning in LLMs"

Alexandru Tăbușcă<sup>1</sup>

1. Computer Science for Business Management, Romanian-American University, Romania

This article introduces an innovative framework for analyzing multi-hop reasoning in large language models (LLMs) through Hamiltonian mechanics and differential geometry. Reasoning processes are represented as trajectories within embedding spaces, with kinetic energy modeling the progression of reasoning and potential energy encoding the alignment between reasoning steps and question relevance. The author formalizes a “reasoning Hamiltonian,” explores canonical transformations, and employs Frenet-Serret geometry to characterize curvature and torsion of reasoning trajectories. Using BERT embeddings and the OpenBookQA dataset, the study empirically compares valid and invalid reasoning chains, analyzing their “energetic” profiles and geometric dynamics.

The central finding is that valid reasoning corresponds to “lower Hamiltonian energies” and “smoother geometric trajectories,” implying cognitive efficiency and stability. Invalid reasoning tends to exhibit higher energy variance and irregular geometric patterns. The author concludes that the Hamiltonian formalism provides a novel mathematical language for understanding reasoning, interpretability, and potential conservation laws within cognitive or AI systems.

- Ethics and Integrity

The article is original and adheres to ethical standards. It relies on publicly available data (OBQA, QASC) and standard open models (BERT-base). The article’s transparency in methods and the reproducibility of its experiments are commendable. Citations are broad and interdisciplinary, though a few mathematical derivations would benefit from clearer attribution to classical mechanics sources (e.g., Goldstein 1980; Arnold 1989). No ethical, data-handling, or integrity issues are apparent.

- Quality

This is a highly technical and conceptually ambitious article, merging physics-inspired formalism with natural language reasoning analysis. The theoretical framework—defining a Hamiltonian for reasoning and applying canonical transformations—is mathematically sound and internally coherent. The author effectively bridges multiple domains: theoretical physics, AI interpretability, and cognitive modeling.

The empirical section complements theory with statistical rigor: t-tests, MANOVA, and PCA-based geometric visualization are used appropriately to compare reasoning trajectories. However, the empirical evidence remains preliminary—the dataset (OpenBookQA) is limited in scale and complexity, and results, while suggestive, do not yet confirm the universality of the framework.

I consider that a major strength lies in the mathematical transparency and the integration of symbolic reasoning with geometric interpretation, while a key limitation is the lack of ablation or comparative baselines (e.g., with simpler metrics such as cosine similarity or entropy of reasoning chains). The connection between physical conservation laws and cognitive processes is elegant but still metaphorical.

- Novelty

The work demonstrates exceptional originality. Modeling reasoning chains via Hamiltonian dynamics and associating reasoning “energy” with logical validity represents a strikingly fresh conceptual leap (at least to my knowledge, this is the first article of this length and complexity focused on exactly this topic). The introduction of phase-space reasoning analysis, canonical transformations in embedding space, and the use of Frenet-Serret curvature to describe cognitive flexibility are unprecedented in AI interpretability literature.

This originality is twofold:

- Theoretical novelty: introducing conservation laws and symmetries into reasoning dynamics (via Noether-like analogies).
- Methodological novelty: quantifying reasoning trajectories geometrically and statistically.

The author situates this framework in a lineage of physics-inspired AI models but expands it significantly. It positions reasoning not as discrete symbolic steps but as continuous trajectories constrained by dynamic invariants—a paradigm potentially transformative for understanding LLM cognition.

- Impact

If further validated, this framework could redefine how reasoning quality, coherence, and interpretability are measured. It could lead to new diagnostic tools for evaluating reasoning in AI systems—especially in explainability and bias detection. The interdisciplinary impact extends to computational cognitive science, philosophy of AI, and physics-inspired computation.

However, the practical impact is currently limited by the lack of demonstrable improvements in model performance or reasoning generation. The work is primarily diagnostic and theoretical. Future efforts should test whether manipulating Hamiltonian parameters (or enforcing conservation constraints) can improve reasoning stability or steer model outputs.

Still, as a conceptual contribution, the article offers paradigm-level innovation, connecting the geometric formalism of physics to the abstract reasoning processes of LLMs—a potentially foundational step toward the physics of cognition.

- Language and Organization

The text is quite dense but articulate. The mathematical sections are rigorous and carefully written, though they will really challenge readers without backgrounds in classical mechanics or differential geometry. The structure (Introduction – Theoretical Foundations – Framework – Methodology – Results – Discussion) is clear and logical, and the figures (especially Figures 2-18) greatly enhance understanding.

Minor language refinements could further improve flow:

- Simplify extended mathematical explanations by moving detailed derivations to an appendix.
- Clarify the physical-cognitive analogies (e.g., distinguishing metaphorical from literal uses of “energy,” “momentum,” or “symmetry”).
- Enhance accessibility by summarizing mathematical results qualitatively for interdisciplinary readers.

Overall, the writing reflects a great level of scholarly maturity, precision, and interdisciplinary fluency.

- Recommendations for Improvement

- Broaden empirical evaluation—apply the Hamiltonian framework to multiple reasoning benchmarks (e.g., HotpotQA, StrategyQA) to test robustness.
- Introduce baselines—compare Hamiltonian energy metrics to simpler alternatives (entropy, cosine coherence) to validate unique explanatory power.

- Clarify metaphoric scope—explicitly define which physical analogies (Hamiltonian, Noether symmetries) are metaphorical versus directly computable.
- Visualization refinement—enhance figures with clearer legends and consistent color schemes distinguishing valid/invalid trajectories.
- Extend discussion—reflect more deeply on philosophical implications: does reasoning obey conservation laws, or are these emergent regularities?
- Potential for real-time application—speculate on how the framework could be embedded into reasoning-monitoring modules for LLMs.

Addressing these points would elevate the current work from conceptual innovation to empirical significance.

## Declarations

**Potential competing interests:** No potential competing interests to declare.