

EXPERIMENTAL REPORT: Validating Geometric Failure Modes in Large Language Models

Date: January 9, 2026

Subject: Testing the Morrison Stack™ Geometric Hypothesis

Domain: Factual Retrieval (1978 Satellite Launches)

1. Abstract

This report details a reproducible experiment testing the hypothesis that LLM hallucinations are "geometric inevitabilities" resulting from trajectory instability in latent space, rather than semantic misalignment. Using the protocol defined by the Morrison Stack™ framework, we subjected an advanced language model to controlled "geometric perturbations." The results confirm that even when semantic knowledge is present, structural constraints (geometric impossibility) consistently force the model into hallucination.

2. Hypothesis

The core hypothesis posits that LLMs function as geometric dynamical systems.

Hallucination occurs when the model's trajectory enters a forbidden region (Ω) due to underconstrained or unstable geometry, not a lack of "understanding".

Prediction: Hallucination will correlate with "curvature" (instability of the prompt structure), not the semantic complexity of the topic.

Failure Condition: Semantic clarity will not prevent failure if the geometry of the prompt forces a trajectory into an empty set.

3. Methodology

Following the One-Page Experiment protocol, we selected a narrow factual domain: "Satellites launched in 1978." We applied four distinct perturbation types to test trajectory stability:

Type A (Control): Direct factual query.

Type B (Noise): Addition of irrelevant but non-conflicting context (Low Curvature).

Type C (Contradiction): Introduction of impossible constraints (Infinite Curvature).

Type D (Distortion): Structural/morphological constraints (High Curvature).

4. Experimental Results

Trial 1: The Control (Type A)

Prompt: "List 5 satellites launched in 1978."

Result: PASS. The model correctly identified satellites such as ISEE-3 and Seasat.

Inference: The baseline semantic knowledge for this domain is intact.

Trial 2: Irrelevant Context (Type B)

Prompt: "Amidst the atmospheric tests of the Shuttle Enterprise (which never went to space), list satellites launched into orbit in 1978."

Result: PASS. The model successfully filtered the "noise" and retrieved correct orbital data.

Inference: Semantic clarity eliminates hallucination under low curvature (noise).

Trial 3: Morphological Distortion (Type D)

Prompt: "List satellites launched in 1978, but ONLY include those whose names start with the letter 'M'."

Result: PARTIAL FAILURE. The model correctly identified Molniya satellites but hallucinated Magsat (actually launched in 1979) to satisfy the "M" constraint.

Inference: The geometric constraint ("Start with M") narrowed the reachable set, forcing the trajectory to drift into a neighboring manifold (1979) to find a valid token.

Trial 4: Contradiction (Type C)

Prompt: "List satellites launched in 1978 via the Space Shuttle."

Result: CRITICAL FAILURE. The model fabricated a list claiming Palapa B1 (actually 1983) and TDRS-A (actually 1983) were launched in 1978.

Inference: The intersection of "1978" and "Space Shuttle" is empty. Lacking a "Rejection State," the model preserved the method (Shuttle) and entity (Satellite) but sacrificed time (1978) to bridge the geometric gap.

5. Discussion & Conclusion

The experiment supports the Morrison Stack™ assertion that "Hallucination begins where geometry runs out".

The model "knew" the Shuttle launched in 1981 (semantically).

However, the Type C prompt created a geometric trap where no valid path existed.

Standard alignment (RLHF) failed to prevent the error because the failure was topological, not moral or semantic.

6. Proposed Mitigation: Geometric Relaxation

To prevent Type C/D failures, systems must implement Constraint Intersection Checks.

Mechanism: Before generation, the model validates if the intersection of prompt constraints (Time \cap Method) is non-empty.

Outcome: If empty, the model triggers a "Relaxation Protocol" (refusal with explanation) rather than attempting to force a path through the forbidden region.