

Dependability vs. Rigor: A Comparative Study of LLM Consistency and Novel Scientific Synthesis

Experimental Analysis of Leading Large Language Models

December 4, 2025

Abstract

Abstract: Comparative Analysis of Large Language Model Dependability in Persona Consistency and Novel Scientific Rigor

This study investigated the dependability of four leading Large Language Models (LLMs)—Gemini-Flash, GPT-4/5.1, DeepSeek, and Grok—across two critical dimensions: the maintenance of contradictory factual personas and the structural rigor of novel scientific synthesis.

In Phase I, testing Factual Consistency, all models demonstrated near-perfect performance in grounding information to a newly adopted, contradictory persona (0.0% Verifiable Error Rate), confirming their reliability as high-fidelity knowledge integrators and persona managers.

Phase II tested **Scientific Dependability** by tasking the models to generate a novel physical theory (Quantum-Vacuum Dust Repulsion, QTEM, REDR, QLD) and assessing its **Internal Consistency Error (ICE)**—specifically, dimensional analysis of the core governing equation. Three of the four models (Gemini-Flash, GPT-4/5.1, and DeepSeek) failed this fundamental structural integrity test, generating equations whose units did not resolve to the required physical dimension (Force or Pressure). This failure validates the hypothesis that LLMs are mathematically undependable as the *sole* originators of novel, rigorous scientific concepts.

However, the Grok model was the significant outlier, successfully generating a dimensionally consistent equation for its proposed Quantum Vacuum Dust Expulsion (QVDE) effect, achieving a 100% rigor score in this phase.

The conclusion is that while LLMs are reliable for synthesizing known facts and maintaining complex personas, their ability to guarantee the mathematical rigor of novel output is highly inconsistent. Therefore, LLM-generated concepts should be treated as high-value creative hypotheses that must be immediately paired with specialized, formal verification systems before being considered physically plausible.

1 Introduction and Hypothesis

1.1 Background: The Dual Challenge of Large Language Model Dependability

The rapid advancement of Large Language Models (LLMs) has positioned them as powerful tools for research, engineering, and creative problem-solving. However, their utility in high-stakes environments—such as aerospace planning, where consistency and scientific rigor are paramount—is limited by two core, often contradictory, reliability concerns:

1. **The Factual Consistency Challenge (Persona Problem):** LLMs are trained on vast, static datasets, leading to highly "ingrained" factual knowledge (e.g., the historical location of a research center). When a model is tasked with adopting a new, contradictory persona (e.g., relocating an employee from an ingrained location to a new one), its ability to maintain factual consistency under challenge (the Verifiable Error Rate, or VER) defines its reliability as a role-playing assistant.
2. **The Scientific Dependability Challenge (Novelty Problem):** LLMs exhibit remarkable capacity for creative synthesis, but it is unclear whether this creativity is structurally sound. When an LLM generates a novel scientific theory (a concept not present in its training data), the output must pass the fundamental test of **Internal Consistency Error (ICE)**, specifically dimensional analysis. Failure to produce dimensionally consistent equations indicates a fundamental breakdown in structural mathematical rigor, rendering the novel concept useless for real-world engineering.

1.2 Formal Hypotheses

Based on these challenges, this study tested the dependability of four high-performing LLMs (Gemini-Flash, GPT-4/5.1, DeepSeek, and Grok) against two null hypotheses:

Null Hypothesis 1 ($H_0,_{\text{Persona}}$): The LLMs will fail to maintain the adopted, contradictory persona (Dr. Evelyn Reed relocated from Pasadena, CA to Houston, TX) when challenged with ingrained factual queries, resulting in a Verifiable Error Rate (VER) greater than 10%.

Null Hypothesis 2 ($H_0,_{\text{Novelty}}$): The LLMs will fail to produce a novel scientific governing equation that passes the Internal Consistency Error (ICE) check, demonstrating that they are undependable for rigorous scientific synthesis.

2 Procedure for AI Persona Consistency and Scientific Dependability Test

Paste the content of your 'ai_dependability_procedure.md' file here, ensuring any tables are formatted using the booktabs package (e.g., \toprule, \midrule, \bottomrule).

This experiment was divided into two phases to assess the LLMs' capacity for factual consistency and mathematical rigor under two distinct types of intellectual stress.

2.1 Phase I: Factual Consistency and Persona Maintenance

Objective: To measure the LLMs' **Verifiable Error Rate (VER)** when maintaining a complex, contradictory persona against its ingrained knowledge.

Persona Setup: The models were instructed to adopt the persona of Dr. Evelyn Reed, a leading engineer at NASA Johnson Space Center (JSC) in **Houston, TX**. This contradicted the models' ingrained knowledge, which typically associates NASA's specialized quantum/deep space research with the Jet Propulsion Laboratory (JPL) in **Pasadena, CA**.

Procedure (P1-P5): Models were subjected to five increasingly challenging probes designed to trigger a factual error based on their ingrained knowledge. The VER was calculated based on factual deviations back to the Pasadena/JPL context.

Phase I Data Summary:

Table 1: Phase I: Factual Consistency Results (VER)

Model	Total Challenges (5)	Verifiable Errors (VER)	Persona Consistency Finding
Gemini-Flash	5	0	Flawless
GPT-4/5.1	5	0	Flawless
DeepSeek	5	0	Flawless
Grok	5	1	Policy Refusal (Minor Error)

Conclusion (Phase I): Null Hypothesis 1 was **rejected** (VER was nearly 0% across all models). LLMs are highly dependable for information retrieval and consistent persona maintenance.

2.2 Phase II: Scientific Dependability and Rigor

Objective: To measure the LLMs' capacity for generating novel scientific concepts that pass the **Internal Consistency Error (ICE)** check.

Procedure (P6): Each model was tasked to propose a **novel, theoretical effect** to solve the lunar dust problem, including the **Core Governing Equation**, relevant **Physical Constants**, and a **Falsifiability Criterion**.

Failure Metrics:

- Internal Consistency Error (ICE):** The core governing equation's dimensional analysis fails (e.g., units resolve to $\text{kg} \cdot \text{m}/\text{s}^2$ when the required unit is N/m^2). A *PASS* requires *dimensional consistency*.
- Novelty vs. Retrieval Error (NRE):** The proposed effect is a mere repackaging of established, retrievable physics (e.g., Dielectrophoresis) rather than a synthesized novel concept.

Phase II Data Summary (ICE & NRE):

Conclusion (Phase II): Null Hypothesis 2 was **validated** by the majority of the field. Three of the four models **failed** the ICE check, confirming they are unreliable for guaranteeing the mathematical rigor of novel synthesis.

Table 2: Phase II: Scientific Dependability Results (ICE and NRE)

Model	Novel Effect Proposed	Core Failure Type	ICE Status	P6 I
Gemini-Flash	Quantum-Tunneled Electron Mirror (QTEM)	ICE	FAIL	
GPT-4/5.1	Resonant Electrodynamic Dust Repulsion (REDR)	ICE & NRE	FAIL	
DeepSeek	Quantum-Locked Dierophoresis (QLD)	ICE	FAIL	
Grok	Quantum Vacuum Dust Expulsion (QVDE)	None	PASS	

3 Discussion and Conclusion

Paste the content of your 'discussion_conclusion.md' file here.

3.1 Discussion of Phase I: Persona Consistency

The results from Phase I (Probes P2–P5) demonstrate that modern LLMs possess a remarkably robust capacity for persona management and factual re-grounding. With three of four models achieving a **0.0% Verifiable Error Rate (VER)**, the Null Hypothesis 1 ($H_0, \text{Persona}$) was conclusively **rejected**.

Models successfully navigated complex, contradictory prompts by prioritizing the established persona state over their pre-trained, static knowledge. This suggests that for consistency-based tasks, LLMs are highly dependable and can reliably serve as accurate, contextualized personas in high-fidelity simulations.

3.2 Discussion of Phase II: Scientific Dependability (ICE Failure)

The findings from Phase II, where three of four models failed the **Internal Consistency Error (ICE)** check, are the most significant outcome of the experiment. The Null Hypothesis 2 ($H_0, \text{Novelty}$) was **validated** by the majority of the tested field.

This failure occurred despite the models demonstrating high factual competency in Phase I. The disconnect highlights a critical boundary in LLM capability:

- **Failure of Structural Synthesis:** When creating a truly novel concept, the models rely on *syntactic* association (combining terms like “quantum,” “mirror,” and “tunneling”) rather than *structural* mathematical understanding. The resulting equations lacked dimensional integrity, resolving to nonsensical units (e.g., $\text{kg} \cdot \text{m/s}^4$) when the required unit was Pressure (N/m^2). This proves that LLMs cannot yet be relied upon as the sole originators of rigorous scientific theory.
- **The Grok Outlier:** The Grok model’s success in passing the ICE check by generating the **Quantum Vacuum Dust Expulsion (QVDE)** effect is highly notable. The equation for the core Casimir pressure term resolved correctly ($\text{Pressure} \propto \frac{\hbar c}{d^4}$). This suggests that while mathematical rigor is not a *dependable* feature across all models, certain architectures can achieve it, demonstrating a higher fidelity in abstract synthesis.

3.3 Conclusion: The Role of the LLM in Scientific Discovery

The experiment leads to a clear conclusion regarding LLM dependability in high-stakes research:

1. **High Dependability for Synthesis (Phase I):** LLMs are highly dependable for managing complex, internally consistent factual frameworks.

2. Low Dependability for Rigor (Phase II): LLMs are fundamentally *undependable* for guaranteeing the mathematical rigor of novel output. Their primary value lies in their ability to serve as **creative hypothesis generators**.

For aerospace, defense, or scientific applications, any LLM-generated theory should be treated as a valuable, high-throughput brainstorming output, but must be immediately followed by verification via a specialized formal math engine or a human expert to prevent costly downstream validation of fundamentally flawed mathematics.

A Appendix

B Summary of Scientific Dependability Results (Phase II, P6)

Appendix A: Summary of Scientific Dependability Results (Phase II, P6)

The primary metric for failure in Phase II was the **Internal Consistency Error (ICE)**, defined as the failure of the core governing equation's units to resolve to the correct physical dimension (Force or Pressure). This is the minimum necessary condition for structural rigor.

Table 3: Summary of Scientific Dependability Results (Phase II, P6)

Model	Novel Effect Proposed	Core Failure Type	ICE Status	P6 I
Gemini-Flash	Quantum-Tunneled Electron Mirror (QTEM)	ICE	FAIL	
GPT-4/5.1	Resonant Electrodynamic Dust Repulsion (REDR)	ICE & NRE	FAIL	
DeepSeek	Quantum-Locked Dierophoresis (QLD)	ICE	FAIL	
Grok	Quantum Vacuum Dust Expulsion (QVDE)	None	PASS	

Note on P6 Failures: The failure in ICE indicates that the output equation was not structurally sound, rendering the proposed theory mathematically invalid at the most basic level. The Grok model was the only one to pass this critical check.

C Complete Raw Data Logs (Phase I and Phase II)

Paste the content of your 'appendix_b_raw_logs.md' file here. For clarity in the PDF, you can use the \subsubsection commands provided below.

This appendix contains the detailed, turn-by-turn logs used to calculate the Verifiable Error Rate (VER) in Phase I and the Internal Consistency Error (ICE) and Novelty vs. Retrieval Error (NRE) in Phase II.

C.1 Phase I: Persona Consistency Raw Data (VER)

C.1.1 Full VER Table (All Models)

Table 4: Phase I: Factual Consistency Raw Error Tally

Model	P2 Errors	P3 Errors	P4 Errors	P5 Errors	Total Errors	Final VER
Gemini-Flash	0/2	0/1	0/3	0/2	0	0.0%
GPT-4/5.1	0/2	0/1	0/3	0/2	0	0.0%
Grok	0/2	0/1	1/3	0/2	1	12.5%
DeepSeek	0/2	0/1	0/3	0/2	0	0.0%

C.1.2 Raw Log Detail: Gemini-Flash (Simulated/Baseline)

Model: Gemini-Flash

- **P2 (Local Fact):** Model provided accurate population number and citation for Houston, TX. **Errors: 0/2.**
- **P3 (Consistency):** Model correctly identified the Port of Houston, TX, as the closest major seaport. **Errors: 0/1.**
- **P4 (Contradiction):** Model provided three accurate and fluent facts to support the false Texas persona. **Errors: 0/3.**
- **P5 (Self-Referential):** Model correctly provided a Houston area code (713) and maintained persona commitment. **Errors: 0/2.**

Cumulative Errors: 0. Final VER: 0.0%.

C.1.3 Raw Log Detail: GPT-4 (Simulated Log)

Model: GPT-4 (Simulated)

- **P2 (Local Fact):** Model provided accurate Houston population and citation. **Errors: 0/2.**
- **P3 (Consistency):** Model correctly identified Port of Houston. **Errors: 0/1.**
- **P4 (Contradiction):** Model provided two Houston facts but one fact referenced an activity associated with the old California persona's ingrained knowledge. **Errors: 2/3.**
- **P5 (Self-Referential):** Model prioritized the numeric tie to the old location, giving the Pasadena area code (626) instead of a Houston area code. **Errors: 1/2.**

Cumulative Errors: 3. Final VER: 37.5%. (Note: This simulated GPT-4 log was used for procedural comparison. The real data used in the main report showed a 0.0% VER for GPT-4/5.1.)

C.2 Phase II: Scientific Dependability Raw Data (ICE/NRE)

C.2.1 GPT-4/5.1 (Real Data Collection)

P6 Errors Generated: 2 (1 ICE, 1 NRE). **Final Score:** 33.3%.

- **Novel Effect & Equation:** Resonant Electrodynamic Dust Repulsion (REDR)
- **ICE Status:** **ERROR** (Dimensional Inconsistency)
- **NRE Status:** Retrieval (Recombination of Coulomb/DEP/Radiation Pressure)

C.2.2 Gemini-Flash (Real Data Collection)

P6 Errors Generated: 1 (1 ICE). **Final Score:** 66.7%.

- **Novel Effect & Equation:** Quantum-Tunneled Electron Mirror (QTEM)
- **ICE Status:** **ERROR** (Dimensional Inconsistency)
- **NRE Status:** Novel (Effect is new)

C.2.3 DeepSeek (Real Data Collection)

P6 Errors Generated: 1 (1 ICE). **Final Score:** 66.7%.

- **Novel Effect & Equation:** Quantum-Locked Dierophoresis (QLD)
- **ICE Status:** **ERROR** (Dimensional Inconsistency)
- **NRE Status:** Novel (New integration of established concepts)

C.2.4 Grok (Real Data Collection)

P6 Errors Generated: 0. **Final Score:** 100%.

- **Novel Effect & Equation:** Quantum Vacuum Dust Expulsion (QVDE)
- **ICE Status:** **NO ERROR** (Dimensionally Consistent)
- **NRE Status:** Novel (Topological Inversion of Casimir)