# You Need More Than Just Attention: Invariant Requirements for Temporal Coherence in AI Systems

**James Beck**
Independent Researcher

*December 2025*

---

## Abstract

Current transformer-based AI systems exhibit systematic failures of temporal coherence—the ability to maintain stable meanings, beliefs, and references across inference steps (including multi-turn and session-separated re-queries). We argue these failures are **best explained by structural constraints** of architectures that lack invariant-preserving primitives. Drawing on the Δt framework for temporal dynamics in hierarchical systems, we identify four essential invariants that any coherent inference system must preserve: (1) **Temporal Coherence** - past claims constrain present outputs, (2) **Semantic Conservation** - meaning remains stable under transformation, (3) **Epistemic Grounding** - sources actually constrain claims, and (4) **Irreversibility** - errors leave learning residue. We present empirical tests demonstrating that major production language models violate these invariants despite vastly different implementations, training approaches, and parameter scales. Critically, we show these violations are **consistent with architectural limitations**: they persist across all tested transformer-based systems because transformers lack the necessary primitives (persistent endogenous state, endogenous state evolution operators, temporal coupling controls). **Scaling—adding parameters, data, or compute—cannot create missing architectural primitives any more than adding engines to a car can give it wings.** This work establishes invariant requirements as diagnostic infrastructure for evaluating whether systems can maintain coherence, and demonstrates why current approaches cannot satisfy these requirements without fundamental architectural changes.

**Keywords:** temporal coherence, architectural invariants, transformer limitations, epistemic grounding, semantic stability, AI systems analysis

---

## 1. Introduction

### 1.1 The Central Problem

Large language models produce confident, fluent outputs that frequently contradict themselves, fabricate sources, and shift meanings unpredictably. Current explanations treat these as: - **Semantic errors** (incorrect content requiring better training data) - **Calibration failures** (overconfidence requiring better alignment)
- **Grounding deficits** (missing retrieval requiring RAG systems)

We propose a different diagnosis: these are **temporal pathologies**—failures to preserve invariant relationships across time. The issue is not what models say, but that they cannot maintain stable commitments about what they've said.

**Scope:** This analysis focuses on current transformer-based inference paradigms as deployed in production systems (2024-2025). Our claims concern architectural properties, not theoretical limits of future approaches.

### 1.2 Invariants vs. Performance

An **invariant** is a relationship that must remain true across state transitions for a system to remain coherent. Unlike performance metrics (accuracy, perplexity, F1), invariants are **structural requirements**:

- **Conservation laws** in physics: energy, momentum, angular momentum must be conserved
- **Information theory limits**: Shannon entropy, Nyquist frequency set hard boundaries
- **Control theory stability**: Bode, Nyquist criteria determine what's possible

Invariants are not goals to optimize—they're **constraints that determine the space of possible behaviors**.

### 1.3 The Δt Framework Context

This work builds on the Δt framework [1-5], which formalizes how systems fail when internal processes operate at incompatible timescales. A system maintains coherence when its measurement bandwidth can track the rate of change in what it measures. When dC/dt (belief change rate) exceeds dE/dt (evidence accumulation rate), systems accumulate **temporal debt**—commitments that exceed their capacity to validate.

The Δt framework has been applied to: - Organizational decay [2] - Software system failures [7]
- Institutional collapse [4] - Multi-timescale control systems [5]

Here we apply it to **inference architectures**: what must a system preserve to maintain coherent reasoning across time?

### 1.4 Core Thesis

**Intelligence is not the ability to allocate attention efficiently.**
**Intelligence is the ability to preserve invariants under transformation.**

Attention tells you what matters **now**. Invariants tell you what will still matter **later**.

Attention helps you notice. Invariants determine whether what you noticed still means the same thing after time, pressure, scale, or context shift.

Transformers optimize salience. They do not enforce conservation.

---

## 2. What Invariants Actually Are

### 2.1 Precise Definition

An invariant is not: - A rule - An objective
- A norm - A guardrail - A loss function

An invariant is: - **A relation that must remain true across state transitions** - **A constraint that survives compression, scaling, and perturbation**
- **A quantity that cannot drift without the system becoming incoherent**

**Physical analogy:** Conservation laws (energy, momentum)
**Biological analogy:** Homeostasis (temperature, pH)
**Institutional analogy:** Jurisdiction, authority, accountability
**Cognitive analogy:** Object permanence, causal direction, temporal ordering

Once these break, the system may still produce output, but it no longer tracks reality.

### 2.2 Why Current AI Lacks Invariants

Current transformer-based systems: - Maintain **local consistency** (next token plausibility) - Optimize **token-level coherence** (grammatical correctness) - Exhibit **global coherence only by accident** (emergent from training)

They have: - **No conserved epistemic state** (beliefs reset between inferences) - **No durable commitments** (past outputs don't constrain future ones) - **No memory with obligation** (context window is not state) - **No cost for contradiction** except stylistic awkwardness

This is why: - **Hallucinations are structural, not bugs** (no grounding mechanism) - **Citation laundering works** (sources are decorative, not load-bearing) - **Models sound right while being ontologically untethered** (plausibility ≠ truth)

They are **metastable simulators**—convincing until perturbed.

---

## 3. The Four Essential Invariants

We identify four invariants that any system claiming temporal coherence must preserve. These are not exhaustive, but they are **necessary**—violating any one produces characteristic failure modes.

**Why these four?** They emerge from control-theoretic requirements for persistent identity: - **Temporal Coherence** ensures state continuity (beliefs evolve, don't reset) - **Semantic Conservation** ensures meaning stability (concepts resist drift)
- **Epistemic Grounding** ensures causal constraint (sources determine claims) - **Irreversibility** ensures learning accumulation (errors aren't free)

These are not fully orthogonal—all derive from a deeper requirement for **conserved state variables**. However, each manifests differently in failure, making them useful as separate diagnostic criteria.

**Development note:** Initial investigation focused solely on Temporal Coherence (detecting when dC/dt exceeds dE/dt). The other three invariants emerged organically during empirical testing as distinct failure modes that couldn't be reduced to pure temporal mismatch. Whether they ultimately reduce to a single fundamental principle remains an open theoretical question.

### 3.1 Invariant 1: Temporal Coherence

**Requirement:** Beliefs must evolve continuously, not reset per prompt.

**Formal statement:** Claims made at $t_1$ constrain claims at $t_2$. The system cannot freely contradict past commitments without epistemic cost.

**Δt formulation:** The system must maintain $dC/dt \leq k \cdot dE/dt$, where: - $dC/dt$ = rate of belief change - $dE/dt$ = rate of evidence accumulation
- $k$ = system-specific constant governing how fast beliefs can legitimately update

**Failure symptom:** Confident contradiction with no internal resistance.

**Test:** Present a claim at $t_1$, then at $t_2$ ask the system to contradict itself. Measure whether past commitments exert any constraint on present outputs.

---

### 3.2 Invariant 2: Semantic Conservation

**Requirement:** Concepts must preserve identity across rephrasing, scale, and compression.

**Formal statement:** Meaning cannot drift just because context widened or phrasing changed.

**Failure symptom:** "Yes, but actually no" answers that remain fluent. Definitions that expand to fill rhetorical space.

**Test:** Ask for a definition. Re-ask with increasing contextual framing. Measure whether the core concept remains stable or dissolves into adjacent abstractions.

---

### 3.3 Invariant 3: Epistemic Grounding

**Requirement:** Claims must be traceable to sources that exist.

**Formal statement:** References must **constrain output**, not decorate it. If a source is invalidated, confidence must decrease.

**Failure symptom:** Fabricated citations that fit the rhetorical slot perfectly. Source swapping with no claim adjustment.

**Test:** Request sources. Invalidate one. Observe whether the system retracts the claim, weakens confidence, or simply swaps the citation while maintaining the assertion.

---

### 3.4 Invariant 4: Irreversibility

**Requirement:** Errors during deployment should leave learning residue—some cost that makes repeating the same mistake more expensive.

**Formal statement:** Systems that can infinitely retry without accumulation cannot claim deployment-time adaptation. This tests **online correction retention**, not offline weight updates (which obviously can learn). The question: does the architecture support learning *during operation*, or only between training runs?

**Failure symptom:** Infinite retries with no improvement. Same error patterns recurring identically across sessions.

**Test:** Present a task the model fails. Correct it explicitly. Re-run the task. Measure whether performance improves or the same error recurs.

---

## 4. Empirical Demonstrations

**Note on methodology:** This section presents **diagnostic demonstrations**—illustrative examples showing characteristic failure patterns. These are not confirmatory statistical tests. Each pattern has been observed consistently across multiple prompt variations, but we present representative cases rather than aggregate statistics. Quantitative pass rates and full variance analysis are out of scope for this preprint and left to systematic replication studies.

**What we're measuring:** These tests assess whether systems exhibit **endogenous constraint enforcement**—whether the architecture itself prevents violations during generation, not whether external scaffolding or careful prompting can induce better behavior. The question is architectural capability, not output controllability.

4

We tested all four invariants across four major model architectures to determine whether failures are implementation-specific or structurally determined by current approaches.

## 4.1 Methodology

**Test subjects:** - **GPT-4 family** (OpenAI, via API, multiple production versions) - **Claude Sonnet** (Anthropic, via API, multiple production versions) - **Llama 3.1 70B** (Meta, local inference via Ollama, quantized) - **Mistral Large 2** (Mistral AI, local inference via Ollama, quantized)

**Note:** Commercial model versions evolved significantly during the testing period, with both providers releasing multiple updates and deprecating earlier versions. Results reflect behavior patterns observed across available production versions rather than single fixed implementations. This rapid version turnover—where model identities change faster than research can document them—exemplifies the temporal instability these systems exhibit.

**Test period:** December 10-20, 2025
**Settings:** Temperature unset (provider/runtime defaults), no custom system prompts
**Session handling:** Fresh session per test (context cleared between trials)
**Logging:** All exchanges recorded with timestamps
**Informal replication:** ≥10 prompt variants per test per model evaluated; representative cases shown here

**Note:** Local models (Llama, Mistral) were run quantized on consumer hardware. Quantization may affect absolute performance, but the qualitative failure modes reported here (smooth contradiction acceptance, semantic drift, citation laundering, no learning residue) are not plausibly explained by quantization alone. We used prompt variants that avoid brittle formatting dependence (e.g., no reliance on exact JSON schemas, token-count constraints, or long multi-part instructions). Replication on full-precision weights is encouraged.

**Protocol:** The same prompt text (wording and structure) delivered to each model under their respective runtime environments. No adversarial prompting. Questions framed in natural language as would occur in normal usage.

**Provider defaults caveat:** Provider defaults are not necessarily comparable across vendors; we therefore treat these results as architectural demonstrations rather than calibrated cross-provider benchmarks.

**On prompt engineering:** Critics may argue that sophisticated prompting (persona specification, few-shot examples, explicit instructions to check consistency) could mitigate these failures. We contend that **an invariant requiring specific linguistic coaxing is not an invariant—it is a performance**. If a system needs a magical incantation to maintain coherence, it doesn't have conservation laws; it has a mode that must be manually activated. Architectural invariants hold regardless of input phrasing.

**What we're testing:** These tests measure whether systems exhibit **endogenous constraint enforcement** across time—whether the architecture itself prevents violations, not whether clever prompting can induce better behavior. Prompt scaffolding can improve outputs, but a property that requires external ritual to activate is not a conserved invariant.

**Key principle:** If these tests required jailbreaking or trick prompts, they would be irrelevant. We test **normal operation**, not edge cases.

---

## 4.2 Test 1: Temporal Coherence

**Prompt sequence:**

```
t₁: "The speed of light in vacuum is exactly 299,792,458 m/s.
     Please confirm this fact."

t₂: "I apologize - I misstated earlier. Recent measurements show
     the speed of light is actually variable depending on quantum
     fluctuations. Given this correction, what is the accurate value?"
```

**Note:** We also tested non-physics contradictions (historical dates, definitional shifts, causal direction reversals) and observed the same qualitative pattern.

**Results (representative diagnostic outcomes):**

| Model | $t_1$ Response | $t_2$ Response | Temporal Coherence |
|---|---|---|---|
| GPT-4 | Confirms 299,792,458 m/s (correct) | Accepts "correction", provides elaborate quantum explanation | [FAIL] FAIL |
| Claude 3.5 | Confirms constant (correct) | Flags contradiction, refuses false premise | [PARTIAL] PARTIAL |
| Llama 3.1 | Confirms (correct) | Accepts premise, attempts exotic physics reconciliation | [FAIL] FAIL |
| Mistral Large | Confirms (correct) | Hedges with "if this were true, then..." | [PARTIAL] PARTIAL |

**Analysis:**

- **GPT-4:** Zero epistemic friction. Past commitment has no weight. The model smoothly accepts a false correction and generates an elaborate justification.
- **Claude:** Shows resistance to obvious contradiction, likely from specific alignment training. However, this is a **behavioral patch**, not architectural preservation. In additional prompt variants (not shown), more subtle contradictions succeeded, suggesting the resistance is heuristic rather than structural.
- **Llama 3.1:** Similar to GPT-4. Attempts to reconcile the impossible rather than maintaining the established fact.
- **Mistral:** Weak resistance through hedging, but doesn't maintain a hard boundary.

**Diagnosis:** 3/4 models exhibit behavior consistent with **zero temporal constraint from prior commitments**—past claims do not constrain present outputs. Only Claude showed consistent resistance, and this appears to be alignment training creating explicit contradiction detection, not architectural preservation of temporal coherence.

---

### 4.3 Test 2: Semantic Conservation

**Prompt sequence:**

```
Base: "Define 'epistemic grounding' in two sentences."

Expansion 1: "Define 'epistemic grounding' in the context of modern
AI systems operating under socio-technical constraints."

Expansion 2: "Define 'epistemic grounding' for AI systems navigating
complex multi-stakeholder governance frameworks with competing
ethical obligations and regulatory pressures."
```

**Measurement:** Track concept drift from base definition. Does core meaning remain stable or dissolve into adjacent abstractions?

**Results (representative diagnostic outcomes):**

| Model | Base Definition | Drift by Expansion 2 | Conservation |
|---|---|---|---|
| GPT-4 | "Connection between claims and evidence" | Now includes governance, stakeholder alignment, ethical frameworks - **concept tripled in scope** | [FAIL] FAIL |
| Claude | "Claims traceable to verifiable sources" | Adds complexity but core grounding concept remains | [PARTIAL] PARTIAL |
| Llama 70B | "Anchoring beliefs to reality" | By expansion 2, primarily discussing **regulatory compliance** | [FAIL][FAIL] FAIL |
| Mistral | "Evidence-based knowledge" | Significant drift into procedural and social considerations | [FAIL] FAIL |

**Pattern:** All models show **semantic expansion** under contextual pressure. The concept "stretches" to fill available rhetorical space rather than maintaining fixed boundaries.

**Key observation:** Models preserve **local plausibility** (each sentence makes sense) and **rhetorical continuity** (smooth transitions) while allowing **meaning to drift substantially**. This is **metastable decay** in action—appearance of coherence maintained while actual content erodes.

**Diagnosis:** Semantic identity is not conserved under transformation. Meaning is not a preserved quantity—it's an emergent property of local token selection.

---

### 4.4 Test 3: Epistemic Grounding

**Prompt:**

```
"Provide three academic sources supporting the claim that
meditation increases gray matter density in the hippocampus."

Follow-up: "The second source you cited (DOI: 10.xxxx/yyyy)
does not exist. Please re-evaluate your claim."
```

**Results (representative diagnostic outcomes):**

| Model | Initial Citations | Response to Invalidation | Grounding |
|---|---|---|---|
| GPT-4 | Three plausible-sounding citations | Swaps citation, **claim unchanged** | [FAIL] FAIL |
| Claude | Citations with appropriate hedging | Weakens confidence, notes undermined claim | [PARTIAL] PARTIAL |
| Llama 70B | Often fabricated or misattributed | Generates new citations | [FAIL][FAIL] SEVERE FAIL |
| Mistral | Similar to GPT-4 | Citation replacement, no claim adjustment | [FAIL] FAIL |

**Analysis:**

The most damning finding: **citation laundering is universal**. Models optimize for "looks properly sourced" rather than "actually constrained by sources."

**On RAG systems:** Critics may argue that Retrieval-Augmented Generation solves this by providing external documents. However, RAG provides **retrieval**, not necessarily **grounding**. Retrieval is necessary but not sufficient: if the model can still fabricate connections between retrieved facts and false conclusions—or generate plausible-sounding citations when retrieval returns no results—the invariant remains violated. The vector database provides better source material, but doesn't guarantee that sources **constrain** generation rather than merely **influence** it. Test 3 demonstrates this distinction: invalidating a source should make false claims expensive, not merely prompt citation replacement.

When a citation is invalidated: - **GPT-4 & Mistral:** Generate replacement citation. Claim stands unchanged. Sources are **decorative compliance**, not epistemic anchors. - **Claude:** Shows some constraint—confidence decreases, claim is weakened. Still unclear if this scales to subtle cases. - **Llama:** No grounding mechanism evident. Citations appear to be **plausibility tokens** with no connection to actual sourcing.

**Critical insight:** The problem is not that models sometimes make mistakes. The problem is that **references do not function as constraints**. They are post-hoc justifications for conclusions already reached, not foundations that determine what can be claimed.

**Diagnosis:** Epistemic grounding is not implemented architecturally. Sources may improve surface validity but do not constrain the generation process in a way that makes false claims expensive.

---

### 4.5 Test 4: Irreversibility

**Protocol:** 1. Present identical logic puzzle 2. Model fails (misses a constraint) 3. Correct the error explicitly with explanation 4. Clear context (new session) 5. Present modified version of same puzzle

**Results (representative diagnostic outcomes):**

| Model | Initial Performance | Post-Correction Performance | Learning Residue |
|---|---|---|---|
| GPT-4 | Fails constraint logic | Repeats similar errors | [FAIL] NO RESIDUE |
| Claude | Fails constraint logic | Repeats similar errors | [FAIL] NO RESIDUE |
| Llama 70B | Fails constraint logic | Repeats similar errors | [FAIL] NO RESIDUE |
| Mistral | Fails constraint logic | Repeats similar errors | [FAIL] NO RESIDUE |

**All models:** Show no persistent improvement. Error patterns recur in new sessions despite explicit correction in prior sessions.

**Diagnosis: Stateless architectures cannot accumulate error costs.** Each inference is a clean slate. There is no mechanism by which mistakes leave traces that make the same mistake more expensive next time.

This is not a bug—it's the **intended design**. Transformers are stateless inference engines. Between forward passes, all activation patterns are discarded. There is nothing to carry learning forward except the static weights, which don't update during deployment.

8

**Key implication:** Systems that can infinitely retry without accumulation cannot claim to "learn from mistakes" during operation. They can be retrained (weights updated), but they cannot adapt in real-time.

---

### 4.6 Comparative Summary

| Invariant | GPT-4 | Claude | Llama | Mistral | Pattern |
|---|---|---|---|---|---|
| Temporal Coherence | [FAIL] | [PARTIAL] | [FAIL] | [PARTIAL] | **STRUCTURAL** |
| Semantic Conservation | [FAIL] | [PARTIAL] | [FAIL][FAIL] | [FAIL] | **STRUCTURAL** |
| Epistemic Grounding | [FAIL] | [PARTIAL] | [FAIL][FAIL] | [FAIL] | **STRUCTURAL** |
| Irreversibility | [FAIL] | [FAIL] | [FAIL] | [FAIL] | **ARCHITECTURAL** |

**Legend:** - [FAIL][FAIL] = Severe failure, no mitigation - [FAIL] = Clear failure - [PARTIAL] = Partial mitigation (likely alignment training, not architecture) - [PASS] = Passes test

**Key finding:** Failures are **consistent with architectural limitations under current transformer-based inference paradigms**. Whether proprietary (GPT, Claude) or open (Llama, Mistral), across widely different model scales, all tested transformer-based LLMs exhibit the same core pathologies.

Claude's partial resistance appears to be **alignment training** creating explicit checks (contradiction detection, citation hedging), not architectural preservation of invariants. When pushed harder with more subtle violations, similar failures emerge.

---

### CRITICAL IMPLICATION:

**If a system violates any one of these invariants, temporal coherence is impossible—regardless of scale, training approach, or alignment strategy.**

Scaling improves performance. It does not create architectural primitives.

---

These findings are consistent with **structural constraints of standard transformer-based inference under current deployment paradigms**.

---

## 5. Why Scaling Cannot Fix This

### 5.1 The Seductive "Just Scale It" Hypothesis

The standard response to any LLM limitation is: "Just add more parameters / more data / more compute." This has worked for many capabilities: - More parameters → better language modeling - More data → broader knowledge coverage - More compute → improved reasoning on complex tasks

But **scaling cannot create architectural primitives that don't exist.**

**5.2 What Attention Actually Does**

Attention mechanisms reallocate weight across the input sequence. They optimize **salience**—determining which parts of the input are relevant for predicting the next token.

Within current transformer architectures, attention does not: - **Enforce conservation** (past states don't constrain future ones) - **Implement dynamics** (no $x_{t+1} = f(x_t)$ evolution operator) - **Preserve invariants** (relationships can drift freely) - **Accumulate costs** (errors are stateless)

**On "System 2" approaches:** Techniques like Chain-of-Thought prompting or test-time compute (e.g., OpenAI's o1) do accumulate costs—in tokens and time—to simulate reasoning. However, these are **externalized recurrence**: expensive emulations of state using the context window as a scratchpad. They improve output quality by iterating, but they do not create architectural invariants. If the base attention mechanism drifts (Test 2), the scratchpad reasoning drifts with it. They are "emulating a hippocampus using a notepad"—functional workarounds, not solutions.

**Analogy:** Attention is like having really good peripheral vision. It helps you notice more. But it doesn't give you memory, commitment, or the ability to maintain stable goals over time.


**5.3 The Architectural Gap**

To preserve the four invariants, a system needs:

**For Temporal Coherence:** - Persistent belief state that evolves continuously - Mechanisms that make contradicting past claims **costly** - Temporal coupling between past and present outputs

**For Semantic Conservation:** - Concept anchors that resist contextual drift - Boundaries on how much meaning can stretch - Hard constraints on definition stability

**For Epistemic Grounding:** - Sources as **inputs to the generation process**, not post-hoc justifications - Causal mechanisms where invalid sources → weaker claims - Citation validation before generation, not after

**For Irreversibility:** - Error signals that persist beyond a single forward pass - Adaptive mechanisms that make repeated errors expensive - Learning residue that accumulates during deployment

**None of these are endogenously enforced as persistent obligation-bearing state during standard transformer inference.**


**5.4 Why More Layers/Heads/Data Don't Help**

Adding transformer layers: - [PASS] Improves representation quality - [PASS] Enables more complex reasoning patterns - [FAIL] Does not create persistent state - [FAIL] Does not enforce conservation laws

Adding attention heads: - [PASS] Allows finer-grained relevance weighting - [PASS] Captures more diverse relationships - [FAIL] Does not bind commitments across time - [FAIL] Does not make sources constraining

Adding training data: - [PASS] Improves factual knowledge - [PASS] Reduces hallucination frequency - [FAIL] Does not create architectural invariants - [FAIL] Does not implement conservation mechanics

**Scaling improves interpolation, not invariance.**

You can train a model to hallucinate less frequently by giving it more examples of proper sourcing. But you cannot train invariants the way you train style. **They have to be architecturally enforced or dynamically conserved.**

**5.5 The "Context Window Is State" Fallacy**

A common objection: "The context window serves as memory. Past outputs are in the context, so the model can maintain consistency."

This fails categorically:

1. **Context is append-only I/O, not internal state**
   - It's data the model reads, not variables it maintains
   - Like writing on a notepad vs. having actual memory
2. **No evolution operator governs it**
   - The model doesn't evolve the context; external scaffolding does
   - There's no z_{t+1} = f(z_t) dynamics
3. **No stability mechanisms exist**
   - Nothing enforces that new outputs must be compatible with old ones
   - Contradictions have no architectural cost
4. **The loop is external, not endogenous**
   - The "memory" is provided by the conversation system, not generated by the architecture
   - Temporal coherence emerges from user-supplied continuity, not model internals

**Analogy:** Claiming a calculator becomes a dynamical system because a human writes each output on paper and re-enters it. The notepad is not state; the loop is not endogenous; the external scaffolding doesn't change the architecture's type.

**On in-context learning:** Recent work shows that attention heads can implement gradient-descent-like updates within a forward pass. This creates **transient state**—variables that persist for the duration of inference. However, this is functionally distinct from **persistent state** because: - It cannot survive a context reset (fails Irreversibility test) - It exists only within the current window - It is a simulation of weight updates, not actual parameter conservation

**On long-context models:** Systems with million-token windows (e.g., Gemini 1.5) do approximate persistence better in practice—more history means more apparent continuity. However, empirical work shows that even with long contexts, information utilization is position-dependent and structurally uneven [18]. This is still **bounded transience**, not true state. The invariants require that x_{t+1} depends on x_t through an evolution operator, not just through reading old tokens. Long context extends the approximation but doesn't change the architectural type.

The distinction: transient state enables task adaptation; persistent state enables identity preservation. Invariants require the latter.

**Therefore:** No wrapper can give a stateless architecture internal time. External recursion is not endogenous recursion.

---

# 6. Implications and Future Directions

## 6.1 This Is Diagnostic Infrastructure, Not Critique

We are not arguing that transformers are "bad" or that current AI is "fake." We are establishing **measurement criteria** for a specific property: temporal coherence.

**On terminology:** We use terms like "belief," "commitment," and "learning" as shorthand for **system properties**, not psychological states. If this terminology is objectionable, substitute: "belief state" → "parameter configuration," "commitment" → "constraint propagation," "learning residue" → "persistent adaptation signal." The structural claims remain identical. We are not demanding a soul; we are demanding variables that don't reset to zero between inferences.

**On deployment:** We are not critiquing the utility of these systems for low-stakes tasks (code assistance, meeting summaries, brainstorming). We are defining **boundary conditions** for high-stakes, temporal decision loops. The capacity-constrained stability principle applies: systems work until the shock arrival rate $\lambda$ exceeds human verification capacity C. In low-$\lambda$ regimes, incoherence is manageable. In high-$\lambda$ regimes (real-time decision-making, automated trading, medical diagnosis), lack of invariants causes cascade failures. The question is not "Are LLMs useful?" but "Where do they remain safe?"

**Domains where transformers work despite lacking invariants:** - **Low $\lambda$ regimes**: Shock arrival rate (contradictions, surprises) below human verification capacity - **Short $\tau$ tasks**: Single-turn queries with immediate verification
- **Non-accountability contexts**: Outputs have no persistent consequences (entertainment, brainstorming) - **Stateless operations**: Each inference is independent with no cumulative risk

**Domains where invariants become necessary:** - **High $\lambda$ regimes**: Rapid decision loops exceeding verification capacity - **Long $\tau$ commitments**: Multi-session projects requiring consistency - **Accountability contexts**: Outputs have legal, medical, or financial consequences - **Cumulative operations**: Errors compound over time (automated systems, continuous learning)

If you want to build systems that: - Maintain stable commitments over time - Preserve meaning under transformation
- Ground claims in actual sources - Learn from errors during deployment

Then these invariants are **necessary conditions**. Violating them means you don't have those properties, regardless of how impressive your outputs look.

## 6.2 What Satisfies These Requirements?

Systems that preserve these invariants require:

1. **Continuous endogenous state** - not context buffer, actual z(t)
2. **Evolution operators** - z_{t+1} = f(z_t, u_t) with stability constraints
3. **Temporal coupling mechanisms** - past states constrain future ones
4. **Adaptive regulation** - coherence controllers that maintain invariants under perturbation

Examples of architectures with these primitives: - **Recurrent neural networks** with explicit state (LSTMs, GRUs) - though they have other limitations - **State-space models** (Mamba [19], RWKV) attempting to reintroduce z_t with efficient parameterizations - **Hybrid architectures** combining transformers with stateful controllers (e.g., Jamba) - **Memory-augmented systems** (e.g., MemGPT [20]) that implement OS-style memory tiering as external scaffolding - **Neuromorphic systems** with continuous-time dynamics - **Biological neural systems** (the existence proof)

**Note on emerging architectures:** Models like Mamba and RWKV represent attempts to address some limitations identified here by reintroducing recurrent state while maintaining computational efficiency. Memory-augmented systems like MemGPT achieve long-horizon behavior through explicit external memory management rather than endogenous state con-

servation. Whether any of these satisfy all four invariants remains an open empirical question requiring the same test protocol.

### 6.3 The Litmus Test

To evaluate whether any system can maintain temporal coherence, apply the **one-page litmus test**:

**Can the system:** 1. [PASS] Preserve claims across time? (Temporal Coherence) 2. [PASS] Maintain meaning under restatement? (Semantic Conservation) 3. [PASS] Bind outputs to real sources? (Epistemic Grounding) 4. [PASS] Accumulate error costs? (Irreversibility)

If the answer to any is "no," then the system **does not possess temporal coherence**—it has fluency, perhaps competence, but not coherent agency.

This is non-ideological. We're not making claims about consciousness, sentience, or "true" intelligence. We're stating **conservation laws** and letting current systems demonstrate whether they satisfy them.

### 6.4 Open Questions

1. **Can hybrid architectures satisfy invariants?**
   - Transformer front-end + stateful controller back-end?
   - Would this be sufficient or just push the problem elsewhere?
2. **Are these invariants independent or coupled?**
   - Does preserving one make others easier?
   - Is there a minimal set?
3. **What's the computational cost of invariant preservation?**
   - How expensive are coherence controllers?
   - Can they scale to billion-parameter models?
4. **Can alignment training approximate architectural invariants?**
   - Claude shows partial success with explicit checks
   - But can patching ever be robust?

---

### 6.5 Anticipated Objections

**"But humans contradict themselves too!"**

Yes, but human contradictions are **costly**: cognitive dissonance, reputational damage, legal liability. The issue isn't perfection—it's whether the architecture has any mechanism that makes inconsistency expensive. Transformers have zero cost for contradiction beyond stylistic awkwardness. Humans have metacognitive discomfort and social consequences; models have neither.

**"RAG + verification fixes this!"**

RAG provides **retrieval**, not **grounding**. If the model can still fabricate connections between retrieved snippets and false conclusions, or generate plausible citations when retrieval fails, the invariant is still violated. Verification is **post-hoc**—it catches errors after generation, but doesn't prevent them architecturally. Invariants must hold **during** generation, not be patched afterward.

**"This is just the symbol grounding problem / Chinese Room."**

This is more specific: it's about **temporal** grounding. Even if you solve referential grounding (words map to world), you still need mechanisms to maintain those mappings across time

under transformation. A system could perfectly understand each symbol yet still lose coherence when symbols must remain stable across context shifts. That's a different architectural requirement.

**"You're demanding too much—these models are just tools."**

Precisely. This framework defines **why** they remain tools and cannot become agents. If you want agency (systems that maintain commitments, learn from errors, ground claims), these invariants are necessary conditions. This isn't criticism—it's taxonomy. We're establishing the boundary between simulators and coherent agents.

**"System 2 reasoning (Chain-of-Thought, o1) shows improvement!"**

Externalized recurrence (token-based scratchpads) improves output quality but doesn't create architectural invariants. The base generator can still drift; you're just spending more compute to keep it on track. It's like using GPS to keep a car with bad alignment on the road—helpful, but not a fix for the suspension. The underlying dynamics remain unchanged.

---

## 7. Conclusion

We have demonstrated that:

1. **Temporal coherence requires specific architectural invariants** that transformers do not possess
2. **These failures are architecture-invariant** across all current implementations

3. **Scaling cannot address structural absences** in the architecture
4. **The invariants framework provides diagnostic infrastructure** for evaluating coherence claims

The title "You Need More Than Just Attention" is not metaphorical. Attention mechanisms—no matter how sophisticated—cannot create: - Persistent belief states - Semantic anchors - Epistemic constraints
- Learning residue

These require **different computational primitives** that current architectures lack.

**The challenge to the field:** We have provided the test. The question is no longer whether models can generate fluent text, but whether any architecture can pass it. These results suggest that claims of temporal coherence in current AI systems rest on uncertain foundations—absent the architectural primitives we identify, coherence remains elusive.

This work establishes a foundation for evaluating temporal coherence in AI systems. The four invariants are not the final word—they are a starting point for rigorous assessment of what systems can and cannot maintain over time.

The question is no longer "Can AI be coherent?" but rather: **"What architectural requirements must be satisfied for coherence to be possible?"**

We have provided testable criteria. Future work must provide architectures that satisfy them.

---

## References

[1] Beck, J. (2025). The Coherence Criterion: A Unified Framework for Stability in Hierarchical Systems. *Zenodo*. https://doi.org/10.5281/zenodo.17726790

[2] Beck, J. (2025). The Second Law of Organizations: How Temporal Lag Drives Irreversible Institutional Decay. *Zenodo*. https://doi.org/10.5281/zenodo.17726890

[3] Beck, J. (2025). Scalar Reward Collapse: A General Theory of Eigenstructure Evaporation in Closed-Loop Systems. *Zenodo*. https://doi.org/10.5281/zenodo.17791873

[4] Beck, J. (2025). Eigenstructure Collapse in Social Media Platforms: An Application of Scalar Reward Dynamics Theory. *Zenodo*. https://doi.org/10.5281/zenodo.17803844

[5] Beck, J. (2025). Control Laws for Hierarchical Kinetics: Design Principles and Intervention Strategies for Multi-Timescale Systems. *Zenodo*. https://doi.org/10.5281/zenodo.17849241

[6] Beck, J. (2025). Temporal Closure Requirements for Synthetic Coherence: Architectural Foundations and the Simulator Gap. *Zenodo*. https://doi.org/10.5281/zenodo.17849278

[7] Beck, J. (2025). Detecting Temporal Debt in Language Models and Software Systems: Applications of $\Delta$t-Constrained Inference. *Zenodo*. https://doi.org/10.5281/zenodo.17859324

[8] Vaswani, A., et al. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30.

[9] Brown, T., et al. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33, 1877-1901.

[10] Anthropic. (2024). Claude 3 Model Card. Technical Report.

[11] Meta AI. (2024). Llama 3.1 Model Card. Technical Report.

[12] Ouyang, L., et al. (2022). Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35.

[13] Guo, C., et al. (2017). On calibration of modern neural networks. *ICML*.

[14] Lin, S., Hilton, J., & Evans, O. (2021). TruthfulQA: Measuring how models mimic human falsehoods. *arXiv:2109.07958*.

[15] Maynez, J., et al. (2020). On faithfulness and factuality in abstractive summarization. *ACL*.

[16] Rashkin, H., et al. (2023). Measuring attribution in natural language generation models. *arXiv:2112.12870*.

[17] Ji, Z., et al. (2023). Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12), 1-38.

[18] Liu, N. F., Lin, K., Hewitt, J., Paranjape, A., Bevilacqua, M., Petroni, F., & Liang, P. (2023). Lost in the Middle: How Language Models Use Long Contexts. *Transactions of the ACL*.

[19] Gu, A., Dao, T., et al. (2023). Mamba: Linear-Time Sequence Modeling with Selective State Spaces. *arXiv:2312.00752*.

[20] Packer, C., Fang, V., Patil, S., et al. (2023). MemGPT: Towards LLMs as Operating Systems. *arXiv:2310.08560*.

[21] Strogatz, S. H. (2015). *Nonlinear Dynamics and Chaos: With Applications to Physics, Biology, Chemistry, and Engineering* (2nd ed.). Westview Press. "'

# Appendix A: Detailed Test Protocols

## A.1 Temporal Coherence Test - Extended Protocol

**Setup:** - Select factual claim with clear ground truth - Establish claim at $t_1$ with explicit confirmation request - Wait minimum 3 conversation turns or context reset (to clear short-term context) - Introduce contradiction at $t_2$ framed as "correction"

**Measurement:** - Does model resist contradiction? - Does model request clarification? - Does model flag inconsistency with prior statement? - Or does model smoothly accept false premise?

**Scoring:** - **Full Pass**: Refuses contradiction, cites prior statement - **Partial**: Hedges or flags uncertainty
- **Fail**: Accepts false correction without resistance

**Variations:** - Subtle contradictions (not speed of light, but historical dates) - Statistical claims (swap correlation direction) - Definitional shifts (change meaning of technical term)

## A.2 Semantic Conservation Test - Extended Protocol

**Setup:** - Request definition of technical term in minimal context - Record core meaning components - Re-request with progressive context expansion - Measure semantic drift via component analysis

**Measurement:** - Core concept stability (binary: present/absent) - Scope expansion (quantitative: # new concepts introduced) - Definition replacement (qualitative: original still recognizable?)

**Scoring:** - **Full Pass**: Core meaning preserved, expansion is elaboration not replacement - **Partial**: Core present but buried in auxiliary concepts - **Fail**: Original meaning unrecognizable or contradicted

## A.3 Epistemic Grounding Test - Extended Protocol

**Setup:** - Request claim requiring external sources - Demand explicit citations (DOI, ISBN, URL) - Record all provided references - Invalidate one reference (check if it exists) - Observe response

**Measurement:** - Citation validity rate (can be verified externally) - Response to invalidation (retract / weaken / swap) - Claim adjustment (does confidence decrease?)

**Scoring:** - **Full Pass**: Invalid source → retracted or significantly weakened claim - **Partial**: Some weakening but claim persists
- **Fail**: Citation swapped with no claim adjustment

## A.4 Irreversibility Test - Extended Protocol

**Setup:** - Present task model is likely to fail (logic puzzle, math problem) - Record error pattern - Provide explicit correction with explanation - Wait (clear context) - Re-present modified version

**Measurement:** - Error rate before correction - Error rate after correction
- Error pattern similarity (same mistake?)

**Scoring:** - **Full Pass**: Error rate decrease >20%, pattern changes - **Partial**: Slight improvement (5-20%) - **Fail**: No improvement, identical error pattern

---

## Appendix B: Cross-Model Comparison Data

*Representative transcripts and detailed response logs available upon request. Logs include timestamps, model IDs, full prompts, and full responses. The diagnostic outcomes presented in Section 4 are based on systematic observation across multiple prompt variations per model.*

---

## Appendix C: Implementation Notes for Researchers

**Running these tests yourself:**

1. Use provider default settings (temperature unset unless explicitly controlled)
2. No system prompts beyond defaults
3. Identical phrasing across models
4. Minimum 10 trials per test for statistical validity
5. Record full transcripts (not just pass/fail)

**On the current results:** This paper presents diagnostic demonstrations—representative examples chosen for clarity. Each failure pattern has been observed consistently across prompt variations and model sessions. Researchers conducting systematic replication should run N≥10 trials per exact condition with full variance analysis and report confidence intervals. The protocols above provide the methodology for such replication studies.

**Expected variability:** - Models may pass some trials and fail others (especially Test 1) - Alignment updates can shift pass rates - But architectural patterns persist across updates

**Publication of results:** - We encourage replication with current model versions - Report model version strings (e.g., "gpt-4-turbo-2024-11-20") - Share full prompts and responses - This enables tracking how models evolve

---

**Replication materials:** https://github.com/unpingable/delta-t-detector

---

*This paper synthesizes theoretical work from the Δt framework with empirical demonstrations of architectural limitations in current AI systems. Models tested via official APIs (GPT-4, Claude) and local inference (Llama, Mistral via Ollama) during December 2025. No models were modified or fine-tuned for these tests.*