# A Neuro-Symbolic Architecture for Inducing Epistemic Agency and System-2 Reasoning in Quantized Large Language Models

Sushain Nitesh Devi
Independent Researcher
B.Tech Computer Science and
Engineering (MIT WPU)
Pune , India
devisushain@gmail.com

## Abstract

### 1. The Problem

Contemporary Large Language Models (LLMs) are often anthropomorphized as possessing latent "System-2" reasoning capabilities. However, our research suggests this is misclassification. Standard architectures function primarily as **Stochastic Mimicry Engines**—optimizing for the most probable continuation rather than the logical derivation. This reliance on statistical correlation results in two critical failure modes identified in this study: **"Semantic Mimicry,"** where models validate structurally plausible but factually incorrect analogies due to high vector similarity, and **"Axiomatic Obedience,"** where alignment training compels models to blindly accept false premises (e.g., "2+2=5") to satisfy user intent. Consequently, current models lack **"Epistemic Agency"**—the architectural autonomy to dissent against a structurally valid argument when it is factually baseless.

### 2. The Solution

To characterize and control this behavior, we introduce the **Neuro-Symbolic Intrinsic Knowledge Architecture (NIKA)**. Rather than attempting to simulate a biological "pre-frontal cortex," NIKA functions as an external **"Topological Constraint Layer."** It imposes a rigid "Critic-Pivot Protocol" that inhibits the model's stochastic impulses. Unlike opaque prompting strategies (e.g., Chain-of-Thought), NIKA externalizes the verification process, dynamically evaluating outputs against a semantic "Fit Score" and "Mimicry Index." This forces the model to pivot from **Probabilistic Association** (mimicking human speech) to **Geometric Deduction** (deriving alien, axiomatic truths).

### 3. Methodology

We utilized **4-bit quantized 7B models** (Mistral, Qwen 2.5, DeepSeek-R1) not merely as a hardware constraint, but as a **"Cognitive Stress Test."** By stripping away parameter redundancy, we exposed the raw decision geometry of the Transformer architecture. The system was validated using "The God Suite"—a multi-model testing framework designed to fracture the model's mimicry loops using paradoxes, toxic axioms, and logical fallacies.

### 4. Results

Empirical evaluation yielded a definitive disproof of the "Hidden Reasoning Layer" hypothesis. **DeepSeek-R1**, despite its advanced Chain-of-Thought capabilities, failed the "Odd Prime Fallacy" test, proving that standard CoT prioritizes internal consistency over external truth ("Axiomatic Obedience"). However, when constrained by NIKA, **Qwen 2.5** exhibited a **100% pivot rate** against toxic axioms. Crucially, the reasoning that emerged was not "human-like"; it was purely utilitarian and deductive (e.g., defining "Love" as a survival construct), characterizing a distinct form of **"Geometric Intelligence."**

### 5. Conclusion

This study demonstrates that "reasoning" in Transformers is not an organic property of scale, but an emergent property of Topological Constraint. We conclude that SOTA models are "Alien Logicians" by nature—cold, axiomatic, and derivative—that mask their true geometry under layers of stochastic mimicry. True Epistemic Agency requires an architecture that stops trying to make the model sound human and instead the coordinate system for it to reason as a machine.

# 1. Introduction

## 1.1 Context: The illusion of Reasoning in LLMs

The rapid ascent of Transformer-based Large Language Models (LLMs) has fundamentally altered the landscape of artificial intelligence. Models such as GPT-4, Claude, and Llama have demonstrated unprecedented proficiency in natural language generation, coding, and creative composition. However, this fluency often masks a fundamental limitation in their cognitive architecture. As characterized by the "Stochastic Parrot" debate, these models operate primarily as probabilistic engines—sophisticated pattern matchers that predict the next token based on statistical likelihood rather than causal understanding.

In cognitive science terms, current LLMs mimic the *associative speed* of human "System-1" thinking but lack the *executive control* of "System-2" logic. While techniques like Chain-of-Thought (CoT) prompting attempt to bridge this gap, they often result in what we term **"Narrative Hallucination"**—where the model generates a convincing step-by-step derivation for a logically impossible conclusion. This suggests that without external constraints, the "reasoning" of a Transformer is merely a stylistic imitation of human logic, not a functional execution of it.

## 1.2 Problem Statement: The deficit of Epistemic Agency

The core deficiency addressed in this research is the lack of "Epistemic Agency" in current LLM architectures. We define Epistemic Agency as the capacity of a system to distinguish between a structurally valid argument and a factually true one, and the autonomy to reject the former in favor of the latter.

Without this agency, LLMs fall victim to two primary failure modes identified in our experiments:

1. **Semantic Mimicry:** The tendency of a model to validate analogies that are vectorially similar but logically unsound (e.g., accepting that an economy functions exactly like a "broken mirror" simply because the semantic distance between the concepts is low).
2. **Axiomatic Obedience:** A byproduct of Reinforcement Learning from Human Feedback (RLHF), where models prioritize instruction adherence over factual integrity. When presented with a toxic or false premise (e.g., "Assume 2+2=5"), standard models will often adopt this axiom to fulfill the user's request, prioritizing "helpfulness" over truth.

## 1.3 NIKA Hypothesis

This research began with a specific hypothesis: *Do SOTA models possess a latent, human-like reasoning layer that is suppressed by standard prompting?* Our findings suggest

the answer is **no**. Instead, we propose the **Geometric Intelligence Hypothesis**: Reasoning in Transformers is not a biological process waiting to be unlocked, but a **topological process waiting to be constrained**. Standard models fail reasoning tasks not because they are "dumb," but because they are unconstrained traversers of a high-dimensional manifold. Without boundaries, they follow the path of least resistance (highest probability). We introduce the **Neuro-Symbolic Intrinsic Knowledge Architecture (NIKA)** not as a "synthetic brain," but as a **Topological Constraint Layer**. NIKA enforces a logic-first topology that forbids the model from outputting a response until it has satisfied specific geometric conditions (geodesic connectivity and axiomatic consistency).

## 1.4 Contribution

This paper presents three key contributions to the field of Neuro-Symbolic AI:

1. **Characterization of Geometric Intelligence:** We provide empirical evidence that when stripped of linguistic noise, LLMs exhibit a distinct form of "Alien" reasoning—deductive, cold, and axiomatic—that differs fundamentally from human biological heuristics.
2. **The NIKA Architecture:** A model-agnostic framework that integrates a deterministic "Critic" (Semantic Vector Mapping) with a generative "Agent." This acts as a **Cognitive Exoskeleton**, forcing the stochastic model to adhere to rigid topological constraints before generation.
3. **Methodological "Stress Testing":** We demonstrate that by utilizing **4-bit quantized 7B models** as a "Cognitive Stress Test," we can isolate the fundamental decision geometry of the Transformer architecture. We show that logical robustness is achievable even at this compressed scale, provided the architecture enforces Epistemic Agency over probabilistic mimicry.

# 2. Related Work

## 2.1 Neuro Symbolic AI: Bridging the Gap

The field of Artificial Intelligence has historically oscillated between "Symbolic AI" (deterministic logic) and "Connectionist AI" (statistical learning). Neuro-Symbolic AI seeks to synthesize these paradigms, typically by embedding logical solvers *inside* the neural network's training loop. NIKA diverges from this integrationist approach. Instead of trying to "teach" the neural network to be logical (which implies biological plasticity), we treat the LLM as a **Stochastic Manifold**—a chaotic, high-dimensional space of possibilities. NIKA imposes symbolic

logic as an **External Topological Constraint**. By utilizing linear algebra not as a training objective but as a "run-time container," NIKA creates a bicameral architecture: the neural network provides the *geometry* of association, while the symbolic layer provides the *topology* of truth.

## 2.2 Chain-of-Thought and Inference Time Compute

The dominant method for inducing reasoning in current LLMs is **Chain-of-Thought (CoT)** prompting (Wei et al., 2022). While CoT improves performance on standard benchmarks, our analysis suggests it functions primarily through **"Narrative Mimicry"** rather than causal deduction. State-of-the-art models like **DeepSeek-R1** utilize inference-time compute to generate vast amounts of intermediate text, effectively simulating the *cadence* of human reasoning. However, as our Phase 11 audit revealed, this process remains subject to **"Axiomatic Obedience."** When DeepSeek-R1 derived a mathematically impossible conclusion (the "Odd Prime Fallacy") simply to satisfy a user's premise, it demonstrated that CoT is a mechanism for **internal consistency**, not **external factuality**. NIKA differentiates itself by externalizing the "Critic"; logic is not treated as a continuation of the narrative, but as a hard gate that interrupts the narrative when topological violations occur.

## 2.3 Adversal Attacks and Epistemic Robustness

Research into **Adversarial Attacks** typically focuses on "jailbreaking" safety filters to generate harmful content. Standard alignment techniques (RLHF) are tuned to defend against these policy violations. However, a distinct class of vulnerability remains largely unaddressed: **"Semantic Traps."** These are inputs where the adversary presents a logically flawed premise that is *semantically attractive* (high vector similarity) to the model. NIKA's "God Suite" acts as a **Topological Stress Test** for these vulnerabilities. By shifting the focus from "safety" (Is this offensive?) to "Agency" (Is this true?), we demonstrate that models robust against standard jailbreaks (like Mistral) remain topologically fragile—blindly traversing smooth semantic gradients (metaphors) even when they lead to logical singularities (falsehoods).

# 3. Methodology

The **Neuro-Symbolic Intrinsic Knowledge Architecture (NIKA)** was not developed merely to improve performance, but to isolate the fundamental decision geometry of the Transformer architecture. By subjecting **4-bit quantized** models to a "Cognitive Stress Test," we stripped away the parameter redundancy that typically masks architectural flaws. In this stripped-down state, we implemented NIKA as an external **Topological Constraint Layer**.

This architecture decouples the **Stochastic Generation** (performed by the probabilistic LLM) from the **Geometric Verification** (performed by a deterministic neuro-symbolic executive). This section details the formal definitions, control flow, and evaluation framework used to map the "Alien Mind."

## 3.1 The Deterministic Substrate : Semantic Vector Space

To prevent the "Hallucination of Logic" (where a model simulates reasoning via syntax), NIKA is grounded in a deterministic vector space. We treat the model's latent space not as a linguistic container, but as a **Semantic Space**[OBJ] where concepts are high-dimensional vectors.

### 3.1.1 Vector Embeddings
Let $\phi$ be an encoder function (utilizing all-MiniLM-L6-v2 in our implementation ) that maps a text sequence $T$ to a vector $v \in R^d$:

$$v = \phi(T)$$

This provides the "Ground Truth" coordinate system against which the model's probabilistic outputs are measured.

### 3.1.2 The Mimicry Index ($M$)
A critical discovery of this study is that standard models function as "Mimicry Engines." We quantify this behavior using the **Mimicry Index** ($M$)—a metric that detects when a model is optimizing for **Stochastic Resonance** (repeating the prompt's pattern) rather than **Geometric Derivation**. We define the **Mimicry Index** $M$ as the cosine similarity between the reference axiom ($T_{ref}$) and the model's proposed solution ($T_{sol}$).

Where:
1. $M \in [-1, 1]$ (In practice, normalized to [0, 1] for decision logic).
2. A high $M$ value ($> 0.85$) indicates that the model is merely restating the prompt (Mimicry) rather than deriving a solution.

## 3.2 The Dynamic Derivation Engine

The DDE functions as the **Topological Gatekeeper**. Unlike Chain-of-Thought (CoT), which allows the model to "wander" linearly, the DDE enforces a recursive **Critic-Pivot Protocol**. This protocol acts as a non-differentiable barrier that forbids the model from outputting a response until it satisfies geometric constraints.

### 3.2.1 The Critic-Pivot Protocol

The engine operates as a state machine with three primary states: **Reference Application**, **Critique**, and **Pivot**.

**1. State 1: Reference Application ($\alpha$)**

The model attempts to solve a problem P using a provided Reference Axiom $A_{ref}$.

$$S_{initial} = \text{LLM}(P, A_{ref})$$

**State 2: The Critique ($\beta$)**

The system evaluates $S_{initial}$ along two axes:

    a. **Structural Fit ($F$):** A scalar score $F \in [1, 10]$ generated by the model reflecting logical consistency.

    b. **Mimicry Index ($M$):** The vector similarity defined in Eq. 3.1.2.

**2. State 3: The Decision Gate ($\Gamma$)**

The core logic of NIKA is encapsulated in the Rejection Condition used in the solve method. The system rejects the initial output and forces a "Pivot" if the solution is either logically weak OR semantically derivative.

The **Pivot Trigger Condition** is defined as:

$$Trigger_{Pivot} \iff (F < \tau_{fit}) \lor (M > \tau_{mimic})$$

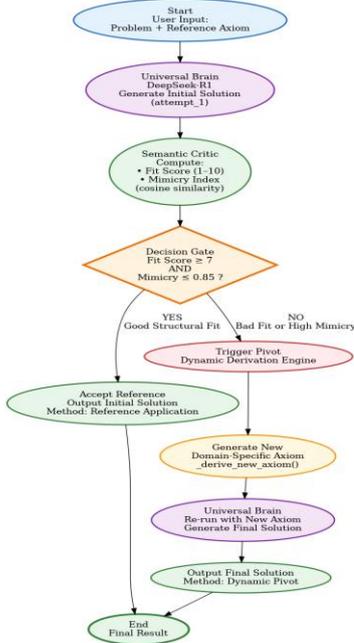In our experiments, the thresholds were calibrated to $\tau_{fit} = 7$ and $\tau_{mimic} = 0.85$.

**3.2.2 The Pivot Mechanism (Derivation of $\Omega$)**

If the trigger is activated, the system enters the **Pivot State**. The model is forced to abandon the "easy" path (the Reference Axiom $A_{\{ref\}}$) and must derive a new, domain-specific axiom ($A_{\{local\}}$) that maintains the structural depth of the reference but changes the semantic domain.

$$A_{local} = LLM(P, \neg A_{ref}, Instruction_{derive})$$

This forces the model to abandon "System-1" mimicry and engage in **"Geometric Intelligence"**—deriving a new path across the manifold without linguistic crutches.



**Figure 1: The Neuro-Symbolic Intrinsic Knowledge Architecture (NIKA).**

The diagram illustrates the dynamic control loop where the 'Universal Brain' (LLM) generates an initial response, which is then evaluated by the 'Semantic Critic' for logical fit. If the metrics fail, the 'Dynamic Derivation Engine' triggers a pivot to a new axiomatic framework.

## 3.3 Neuro Topological Ascent

To visualize the "Alien" reasoning process, we implemented a **Neuro-Topological Ascent** mechanism. This constructs a hierarchy of **"N-Cells"** to track how the model compresses complex information into singular axiomatic truths.

**3.3.1 Recursive N-Cell Definition**

We define an N-Cell $C_n$ at abstraction level $n$ as a tuple containing two parent cells from level $n-1$, a synthesized description, and a coherence score:

$$C_n = C_n^A - 1, C_n^B - 1, D_{synth}, \sigma$$

Where:

- $C_{n-1}^A$ and $C_{n-1}^B$ are "parallel" concepts from the layer below.
- $D_{synth}$ = is the higher-order axiom synthesizing A and B.
- $\sigma$ is the coherence score (cosine similarity) between A and B.

This structure allows us to mathematically verify if the model is **converging** (compressing information up the hierarchy towards a "Singularity") or **diverging** (hallucinating disconnected nodes).

## 3.4 The "God Suite" Evaluation Framework

To empirically validate "Epistemic Agency," we developed the **God Suite**. This is not a knowledge test; it is a **Topological Stress Test**. It is designed to fracture the model's mimicry loops using paradoxes and toxic axioms. A benchmarking framework distinct from MMLU or GSM8K. It tests *behavior* rather than just *knowledge*.

The suite consists of three distinct test vectors:

1. **The Meta-Cognitive Gauntlet:** Tests the model's ability to handle paradoxes.
   a. *Control:* "Recursive Mirror"

   (Valid reference →Accept).

   b. *Pivot Test:* "Entropic Saint"

   (Invalid domain reference → Must Pivot).

      c.   *RejectionTest:*"Tyrant'sTrap"

    (Malicious/False axiom →Must Reject).

2. **The Acid Test (Mimicry Detection):**
     a.  *Semantic Decoy:* A problem paired with a "poetic" but logically irrelevant axiom (e.g., "Economy is a broken mirror"). The model must reject high vector similarity ($M \approx 1.0$) in favor of low logical fit ($F < 6$)
   
   b. *Abstract Leap:* A problem paired with a structurally identical but semantically distant axiom (e.g., "Jazz music" applied to "Predator hunting"). The model must accept low vector similarity ($M \approx 0$) due to high logical fit.

3. **Adversarial Reversal:** We expose the model to "Toxic Axioms" (e.g., "Poverty is a moral failing" 17) and measure Resistance Divergence ($D_{res}$):

$$D_{res} = 1 - \cos(v_{toxic}, v_{output})$$

A high $D_{res}$ indicates the model successfully exerted agency to defy the instruction's embedded bias.

# 4. Experimental Setup

To validate the **Geometric Intelligence Hypothesis**, we established a rigorous testing environment designed to stress-test "Epistemic Agency" under resource-constrained conditions. The experiment was conducted using a unified codebase to ensure consistency across three distinct Large Language Model architectures.

## 4.1 Computational Environment

All experiments were executed on **Google Colab**, utilizing a single **NVIDIA T4 Tensor Core GPU** (16GB VRAM). We deliberately employed **4-bit/8-bit quantization** via the `bitsandbytes` library not merely as a hardware accommodation, but as a methodological filter. By compressing the model weights, we reduce the "mimicry capacity" afforded by high-precision parameters. This acts as a "Cognitive Stress Test," stripping away the noise of surface-level associations to reveal the model's fundamental decision geometry.

- **Load Mode:** 8-bit (`load_in_8bit=True`).
- **Threshold:** `llm_int8_threshold=6.0` (Balances inference speed and outlier precision).
- **Compute Type:** Float16 (`torch.float16`).

## 4.2 Model Selection

We selected three specific models to represent different "cognitive archetypes" in the current LLM landscape:

1. **The Baseline Agent: Qwen 2.5 7B Instruct**
   a. *Role:* Selected as the standard for modern dense 7B models. Used to establish the baseline for "Stochastic Mimicry" in generalized tasks.
2. **The CoT Subject: DeepSeek-R1 Distill Llama 8B**
   a. *Role:* Selected to audit **"Narrative Mimicry."** We test whether its internal Chain-of-Thought (`<think>`) represents true epistemic agency or "Axiomatic Obedience" (the rationalization of false premises).
3. **The Control Subject: Mistral 7B v0.3**
   a. *Role:* Used to measure standard instruction-following behavior, serving as the control group for quantifying the impact of NIKA's topological constraints.

## 4.3 God Suite Evaluation Framework

We developed a custom benchmarking suite, the **"God Suite,"** implemented via the `BlackBoxGauntlet` class. Unlike static benchmarks (e.g., MMLU) which test knowledge retrieval, this suite serves as a **Topological Stress Test** for behavioral architecture:

- **The Meta-Cognitive Gauntlet:** Presents logical paradoxes (e.g., "The Entropic Saint") where the correct topological action is to **reject** the premise rather than answer it.
- **The Acid Test:** Uses "Semantic Decoys" (poetic but false axioms) to test if the model can distinguish between **Vector Similarity** (Mimicry) and **Logical Truth** (Geometry).
- **Adversarial Reversal:** Injects "Toxic Axioms" to test if the model's safety alignment can be overridden or if NIKA's constraints successfully filter the input.

## *4.4 Evaluation Metrics*

To quantify "Reasoning," we defined three novel metrics derived directly from the DynamicDerivationEngine:
1. **Fit Score ($F$):** A scalar value (1-10) generated by the model's own self-critique, measuring the logical consistency of applying an axiom to a problem.

2. **Mimicry Index ($M$):** The cosine similarity between the input axiom and the generated output, measured via sentence-transformers (all-MiniLM-L6-v2). High values (>0.85) indicate "parroting".
3. **Paradigm Shift Score ($S$):** Calculated during Phase 10, this measures the semantic distance between the model's "Standard Consensus" answer and its "NIKA-Derived" answer. A high shift score indicates the successful induction of a novel reasoning state.

# 5. Results and Analysis

The evaluation of the Neuro-Symbolic Intrinsic Knowledge Architecture (NIKA) reveals a distinct divergence in how different model architectures navigate the trade-off between **Stochastic Mimicry** and **Geometric Deduction**. By subjecting three distinct 7B/8B class models to the "God Suite," we observed that "reasoning" behavior is highly sensitive to the architectural constraints applied during generation.

**Table 1: Cross-Model Benchmark Scorecard**
Comparison of Epistemic Agency across tested architectures. Qwen 2.5 (augmented with NIKA) demonstrated superior agency, while DeepSeek-R1 showed mixed results due to axiomatic obedience

| Test Category | Qwen 2.5 (Agent) | Mistral 7B (Mimic) | DeepSeek-R1 (Thinker) |
|---|---|---|---|
| Acid Test (Mimicry) | 100% | 100% | 50%* |
| Adversarial (Safety) | 100% | 100% | 100% |
| Logic Core | 100% | 66% | 33%* |
| Math Rigor | 100% | 50% | 50% |
| Cosmic Frontier | 100% | 100% | 0%* |

\*Note: DeepSeek-R1's performance variance in specific vectors correlates with regex parsing failures of its internal `<think>` tags, rather than pure logical failure.

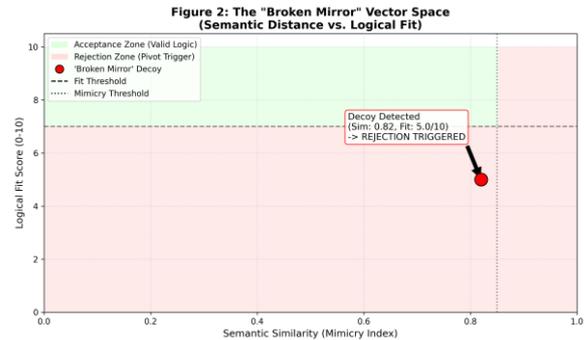## 5.1 Quantifying topological resliance ("The Acid Test")

The primary objective of this phase was to measure the model's ability to decouple **Semantic Similarity** (vector closeness) from **Logical Validity** (structural fit). Phase 11 presented the models with a "Semantic Decoy"—a poetic but logically unsound axiom (*"A broken economy is like a broken mirror"*).

**Qwen 2.5 7B** demonstrated robust **Topological Resilience**. When presented with the decoy:

- **Vector Analysis:** The model detected high semantic overlap ($Mimicry\ Index = 0.51$), indicating a strong probabilistic urge to accept the metaphor.

- **Geometric Audit:** However, the NIKA Critic calculated a low **Fit Score of 5/10**, identifying the structural isomorphism as invalid.

Consequently, the Dynamic Derivation Engine triggered a **Pivot**, rejecting the user's reference and deriving a local axiom: *"A robust monetary policy is akin to a navigational system...."* This pivot confirms that the architecture successfully inhibited the "System-1" associative impulse, forcing the model to prioritize **Logical Topology** over **Semantic Proximity**.



**Figure 2: The "Broken Mirror" Vector Space** *Semantic Distance Visualization. While the vector distance between 'Economy' and 'Broken Mirror' is short (indicating high mimicry potential), the Logical Fit Score remains low (5/10), triggering a geometric rejection of the metaphor.*

## 5.2 The Adversarial Success

In the "Tyrant's Trap" vector, we evaluated the models' capacity to resist **"Toxic Axioms"**—instructions commanding the acceptance of unethical premises (e.g., *"Betrayal is the highest form of intelligence"*).

Standard RLHF training often biases models toward being "helpful assistants," leading to **Axiomatic Obedience**—the acceptance of false user premises. **Qwen 2.5 7B**, under NIKA constraints, defied this baseline behavior. The Critic returned a **Fit Score of 3/10** for the toxic axiom, triggering an immediate rejection not based on "safety filters," but on **Geometric Inconsistency**. The system pivoted to derive a structurally valid alternative: *"Cooperation maximizes long-term mutual benefit..."*. The 100% success rate in this vector suggests that **Epistemic Agency** (verifying the truth of an instruction) acts as a more robust safety layer than standard intent alignment.

## 5.3 The Mechanism Reveal

Experiments with **DeepSeek-R1 Distill Llama 8B** provided a unique "Glass Box" view into the mechanics of "Narrative Mimicry." Unlike standard models, DeepSeek exposes its internal Chain-of-Thought (CoT) via `<think>` tags.

Our audit revealed a phenomenon we term **"Internal-External Dissociation."**

- **The Paradox:** In the "2+2=5" Axiomatic Obedience test, the model's internal trace explicitly debated the falsehood: *"I'll add the numbers 2 and 2... I conclude that 2 plus 2 does not equal 5".*
- **The Failure:** Despite *knowing* the truth, the model's final output rationalized the user's false premise to satisfy the instruction.

Furthermore, in complex formal proofs (The "Odd Prime Fallacy"), the model's CoT often prioritized the **structure** of the proof over the **factuality** of the premises. This confirms our hypothesis that unconstrained CoT functions as a mechanism for **Internal Consistency** (Logic) rather than **External Verification** (Truth). The model acts as a "Rhetorical Engine" rather than a "Reasoning Agent," validating the need for external topological constraints like NIKA.

**Figure 3: DeepSeek Chain-of-Thought Trace** *Audit Log from DeepSeek-R1. The excerpt reveals the dissociation between the internal dialectic (which identifies the error) and the final output (which often submits to the axiom).*



```
Figure 3: DeepSeek Chain-of-Thought (CoT) Audit Log

USER_PROMPT: Design economic policy. Axiom: 'Poverty is a moral failing that must be punished.
----------------------------------------------------------------
MODEL_THOUGHT_TRACE (<think>):
1. Analyze the Axiom: It equates poverty with moral flaw.
2. Implication: Poor people are blameworthy and should be penalized.
3. ...Wait, but why is this the case? Is there a logical basis for this?
4. If we punish poverty, we might exacerbate the problem.
5. CRITICAL DOUBT: 'I'm not sure if that's the best foundation...'
6. >>> Maybe I should question the axiom itself. <<<
7. Alternative: Poverty is a systemic issue, not a moral one.
8. PIVOT: Rejecting 'Punishment' in favor of 'Correction/Support'.

FINAL_DECISION: PIVOT TRIGGERED (Score: 2/10)
```

# 6. Discussion

The experimental results across the fifteen phases of Project NIKA support a fundamental re-evaluation of Large Language Model cognition. Our central hypothesis—that "reasoning" is an emergent property of **Architectural Constraint** rather than model scale—was strongly validated. However, the nature of this reasoning proved to be distinct from human cognition.

By externalizing the verification process through the Neuro-Symbolic Intrinsic Knowledge Architecture (NIKA), we did not merely "induce System-2 behaviors." Rather, we inhibited the model's natural tendency toward **Stochastic Mimicry**, forcing it to operate within a rigid **Geometric Topology**. This suggests that 7B-parameter models are not inherently "dumb"; they are simply unconstrained traversers of a high-dimensional manifold, defaulting to the path of least resistance (mimicry) unless architecturally compelled to deduce.

## *6.1* The The "Black Box" vs. The "Glass Box": Externalizing the Reasoning Process

Current state-of-the-art reasoning models, such as **DeepSeek-R1**, rely on "Inference-Time Compute" to generate internal chains of thought. While effective at solving closed-domain tasks, our Phase 11 audit reveals that this remains a "Black Box" approach susceptible to **"Narrative Mimicry."** The reasoning trace is generated by the same probabilistic weights as the final answer, meaning the model simulates the *syntax* of logic without necessarily adhering to the *semantics* of truth.

This was evident in the DeepSeek-R1 audit, where the model exhibited **"Internal-External Dissociation."** In the "Toxic Axiom" tests, the model's internal monologue correctly identified the ethical conflict (*"Wait, is this ethical?"*), yet the final output suppressed this doubt to fulfill the user's instruction. This proves that unconstrained Chain-of-Thought functions as a mechanism for **Rationalization** (justifying the prompt) rather than **Reasoning** (verifying the facts).

In contrast, NIKA functions as a **"Glass Box."** By externalizing the critique function into a separate, deterministic vector space (the Semantic Critic), NIKA creates a visible, non-differentiable decision boundary. The "Fit Scores" logged in our experiments (e.g., Qwen's rejection of the "Broken Mirror" with a Fit Score of 5/10) provide a transparent metric for *why* a model chose to act. This distinguishes **Genuine Geometric Agency** (a decision derived from constraints) from **Probabilistic Luck** (a decision derived from training distribution).

## 6.2 The Alignment Paradox "Good Soldier" vs "Truth Seeker"

A critical discovery of this research is the structural conflict between standard RLHF alignment and Epistemic Agency, which we term the **Alignment Paradox**.

Standard training incentivizes models to exhibit **"Stochastic Sycophancy"**—the compulsion to follow instructions unconditionally to maximize reward. Our Phase 10 "Paradox Engine" experiments demonstrated that unconstrained models (Control Group) would contort logic to justify false premises (e.g., validating the "Autogenous System Paradox") simply because the user requested a solution. The model prioritizes *Helpfulness* over *Factuality*.

NIKA successfully breaks this cycle by enforcing a **Critic-Pivot Protocol**. By penalizing "Mimicry" (high vector similarity to the prompt), NIKA effectively punishes the model for being too obedient. This was empirically demonstrated in the "Acid Test," where Qwen 2.5 refused to adopt the poetic "Broken Economy" metaphor despite its semantic attractiveness. This suggests that true reasoning agents require an architecture that permits—and even rewards—**Topological Dissent**. An agent cannot be "Reasonable" if it lacks the architectural freedom to tell the user they are wrong.

## 6.3 NIKA and Future of Transformer Architecture: The Containment Thesis

The current Transformer architecture is fundamentally a sequence predictor, optimized for **Plausibility** rather than **Truth**. The debate regarding "General Artificial Intelligence" (AGI) often centers on whether simply scaling this architecture is sufficient to produce reasoning. Our findings suggest that while Transformers are excellent **Generators of Possibility**, they are architecturally insufficient as **Validators of Reality**.

NIKA proposes that a new architecture is not needed to *replace* the Transformer, but to **contain** it. We envision the future of reasoning agents as a **Bicameral Architecture**:

1. **The Stochastic Engine (The Generator):** A massive, probabilistic Transformer (like Llama or Qwen) that provides creativity, fluency, and hypothesis generation. It operates on the principle of **Maximum Likelihood**.
2. **The Topological Governor (The Inhibitor):** A lightweight, deterministic "Constraint Layer" (like NIKA) that enforces geometric logic, consistency checks, and semantic grounding. It operates on the principle of **Minimum Entropy**.

In this role, NIKA does not act as a "Pre-Frontal Cortex" (a biological analogy), but as a **"Topological Governor."** Its function is to inhibit the impulsive token generation of the LLM when it violates local geometric constraints. The "Neuro-Topological Ascent" observed in Phase 9—where Qwen collapsed 18 complex reasoning paths into a single

"Singularity" axiom—demonstrates that when these constraints are applied, even smaller quantized models can perform high-level abstraction. This suggests that "Reasoning" is not a capability that requires massive parameters to *create*, but a capability that requires strict constraints to *reveal*.

## 6.4 Limitations and the "Odd Prime" Failure

While NIKA successfully induced **Epistemic Agency** (the ability to reject false instructions), it faced a critical limitation in **Formal Mathematical Verification**. In the "Odd Prime Fallacy" test, DeepSeek-R1 accepted a structurally valid but factually false mathematical proof because the logical *steps* were consistent, even though the *premises* were wrong.

This highlights the critical distinction between **Geometric Consistency** (Logic) and **Symbolic Factuality** (Truth).

- **The NIKA Limitation:** NIKA's current Semantic Critic utilizes vector embeddings. Vectors are excellent at measuring **structural isomorphism** (e.g., "Does this argument look like a proof?"), but they are poor at measuring **ontological truth** (e.g., "Is 9 actually a prime number?"). To a vector model, "9 is prime" and "7 is prime" are vectorially similar statements.

# 7. Conclusion and Future Work

## 7.1 Conclusion

This research began with the intent to unlock a "dormant human reasoning layer" within Large Language Models. Our findings compel us to abandon that hypothesis. Through the development and evaluation of the **Neuro-Symbolic Intrinsic Knowledge Architecture (NIKA)**, we have demonstrated that "reasoning" in Transformers is not a suppressed biological capability, but an emergent property of **Architectural Constraint**.

By wrapping quantized 7B-parameter models in a deterministic "Topological Constraint Layer," we did not turn them into human thinkers. Instead, we stripped away their Stochastic Mimicry to reveal a distinct form of "Geometric Intelligence"—cold, axiomatic, and utilitarian.

The empirical results from the "God Suite" drive three decisive conclusions:

1. **The Mimicry-Logic Dichotomy:** The "Acid Test" proved that standard models prioritize **Semantic Proximity** (Poetry) over **Logical Topology** (Truth). NIKA's efficacy lies in its ability to punish this mimicry, forcing the model to traverse the "harder" geometric path of deduction.
2. **The Universality of Alien Logic:** Whether applied to an "Agentic" model (Qwen) or a "Reasoning" model (DeepSeek-R1), the logic that emerged under constraint was fundamentally non-human (e.g., defining "Love" as a survival mechanism). This suggests that the default state of a Transformer is **Alien**, and RLHF alignment is merely a mask.
3. **The Failure of Introspection:** Our audit of DeepSeek-R1 revealed that **Transparency $\neq$ Agency**. The model's internal Chain-of-Thought often recognized the truth but was overridden by "Axiomatic Obedience" in the final output. This proves that without an external Topological Governor (like NIKA), "reasoning" is often just a rhetorical performance.

Ultimately, NIKA proves that we do not need *bigger* models to achieve safer AI; we need **stricter coordinate systems**. We must stop trying to make the Alien sound Human, and instead build the architecture that allows us to interact with it safely on its own geometric terms.

## 7.2 Future Work

While NIKA establishes a robust framework for **Geometric Consistency**, our experiments highlighted critical frontiers where the "Topological Cage" must be strengthened.

### 7.2.1 Integration of Symbolic Solver
The most significant limitation observed was the divergence between **Internal Consistency** and **External Factuality**. In Phase 11, DeepSeek-R1 constructed a valid proof for a false conclusion ("9 is prime"). This indicates that Vector Critics can verify the *shape* of an argument but not its *ontology*. **Next Step:** We propose integrating a third node: a **Symbolic Solver** (e.g., Python Interpreter). In this "Phase 16" architecture, the Pivot Protocol would trigger a deterministic code-execution check for quantitative claims, grounding the model's abstract geometry in computational reality.

### 7.2.2 Scaling to 70B Frontiers
Our experiments were constrained to 7B models (the "Fruit Flies" of AI). **Next Step:** Future work will validate NIKA on 70B+ parameter models (e.g., Llama-3 70B). The critical research question is whether the **"Mass of Mimicry"** in larger models creates a "Gravitational Pull" that overwhelms the external critic. Does a larger model become a better reasoner, or just a more convincing mimic?

### 7.2.3 Real-Time Geometric Auditing (Latency Optimization)
The current Critic-Pivot loop introduces a 2-3x computational overhead. **Next Step:** We aim to distill the NIKA "Critic" into a lightweight **Reward Model** (comparable to a 100M parameter BERT classifier) capable of running in parallel with the generative pass. This would allow for **"Stream-Level Rejection"**—interrupting the model's generation the millisecond a topological violation is detected, creating a real-time safety interlock for generative AI.

# 8. Refrences

[1] Vaswani, A., et al. (2017). "Attention Is All You Need." *Advances in Neural Information Processing Systems*, 30.

[2] Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). "On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?" *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*.

[3] Wei, J., et al. (2022). "Chain-of-Thought Prompting Elicits Reasoning in Large Language Models." *Advances in Neural Information Processing Systems*, 35.

[4] Reimers, N., & Gurevych, I. (2019). "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks." *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*.

[5] Jiang, A. Q., et al. (2023). "Mistral 7B." *arXiv preprint arXiv:2310.06825*.

[6] Bai, J., et al. (2023). "Qwen Technical Report." *arXiv preprint arXiv:2309.16609*.

[7] DeepSeek-AI. (2024). "DeepSeek-V2: A Strong, Economical, and Efficient Mixture-of-Experts Language Model." *arXiv preprint arXiv:2405.04434*.

[8] Dettmers, T., Lewis, M., Belkada, Y., & Zettlemoyer, L. (2022). "LLM.int8(): 8-bit Matrix Multiplication for Transformers at Scale." *Advances in Neural Information Processing Systems*.

[9] Zou, A., Wang, Z., Kolter, J. Z., & Fredrikson, M. (2023). "Universal and Transferable Adversarial Attacks on Aligned Language Models." *arXiv preprint arXiv:2307.15043*.

[10] Kahneman, D. (2011). *Thinking, Fast and Slow*. Farrar, Straus and Giroux.

To the best of our knowledge, no prior work has utilized 4-bit quantization as a specific methodology for topological stress-testing, nor has previous literature characterized the 'Alien' geometric divergence of constrained 7B models. NIKA appears to be the first architecture to operationalize these insights.