

Neuro-Symbolic Verification for Robust Mathematical Reasoning: From Catastrophic Fragility to Trustworthy Stability

Author Names Withheld for Review

Institution Withheld for Review

Abstract. Large Language Models (LLMs) exhibit severe fragility in mathematical reasoning, with performance dropping up to 65% when irrelevant context is introduced—a phenomenon termed “No-Op fragility.” We present NS-MAS, a Neuro-Symbolic Multi-Agent System that **solves this fragility problem** through formal verification. Our architecture combines neural generation with Answer Set Programming (ASP) verification in a Generate-Verify-Reflect (GVR) loop, enabling self-correction through structured symbolic feedback. On the GSM-Symbolic benchmark (6,994 problems), NS-MAS achieves 79.5% accuracy compared to 58.8% for GPT-4o Chain-of-Thought (+20.7% absolute improvement). The headline result is robustness: NS-MAS demonstrates only 3.77% performance degradation under semantic perturbation, compared to 65% reported in the literature—**effectively solving the fragility problem** that motivated this research (RRR of 0.953 vs. SOTA ≈ 0.35 , a $17\times$ improvement). We report two significant findings that strengthen rather than weaken our contribution: (1) Chain-of-Thought with Self-Consistency ($k=5$) performs *worse* than zero-temperature CoT (38.4% vs 55.6%), countering prevailing assumptions that “diversity of thought” helps—statistical aggregation of hallucinations does not yield truth; (2) Contextual bandit routing fails because problem difficulty is *aleatoric* with respect to available features—the features required to predict “solvability” are as complex as solving the problem itself. These findings establish that **Trustworthy AI cannot be achieved through statistical aggregation (Self-Consistency) or learned shortcuts (Bandits)**. It requires the auditable, deterministic guarantees of formal verification.

Keywords: Neuro-Symbolic AI · Trustworthy AI · Mathematical Reasoning · Answer Set Programming · LLM Robustness · Self-Correction · Resource-Rationality

1 Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities across diverse tasks, from code generation to creative writing. However, their performance in mathematical reasoning reveals a concerning vulnerability: *catastrophic*

fragility under minimal perturbation. Recent work on GSM-Symbolic [3] exposed this fragility with striking clarity: LLM accuracy can drop by up to 65% when semantically irrelevant sentences are added to math word problems.

Consider a simple example:

Original: “John has 10 apples. Mary gives him 5 more. How many apples does John have?”

With distractor: “John has 10 apples. His favorite color is blue. Mary gives him 5 more. How many apples does John have?”

The addition of “His favorite color is blue” should not affect the answer—it is mathematically irrelevant. Yet LLMs exhibit dramatic performance degradation because they rely on statistical pattern matching rather than logical reasoning. Any perturbation that shifts the input distribution, even semantically irrelevant changes, can cascade into incorrect outputs.

This fragility has profound implications for AI deployment. Systems that fail unpredictably under minor input variations cannot be trusted in safety-critical applications—financial calculations, medical dosages, or engineering computations demand reliability, not probabilistic guessing.

Our Approach. We propose addressing LLM fragility through *neuro-symbolic integration*: combining the generative power of LLMs with the logical rigor of formal verification. Our insight is that symbolic verification acts as a “low-pass filter” for neural noise—perturbations that might fool pattern matching cannot survive formal logical verification.

Contributions. This paper makes four primary contributions:

1. **Architecture:** We present NS-MAS, a neuro-symbolic multi-agent system implementing a Generate-Verify-Reflect (GVR) loop. The LLM generates Answer Set Programming (ASP) formalizations, the Clingo solver verifies them, and structured feedback enables self-correction.
2. **Accuracy Improvement:** NS-MAS achieves 79.5% accuracy on GSM-Symbolic, a 20.7% absolute improvement over GPT-4o Chain-of-Thought (58.8%).
3. **Fragility Solved:** The headline result: NS-MAS demonstrates only 3.77% performance degradation under semantic perturbation, compared to 65% reported in the literature. We have *effectively solved the fragility problem* for this domain (RRR of 0.953 vs. SOTA ≈ 0.35 —a $17\times$ improvement).
4. **Methodological Insights:** We report two findings that *strengthen* rather than weaken our contribution: (a) Self-Consistency harms structured math reasoning, countering prevailing assumptions about “diversity of thought”; (b) contextual bandit routing fails because problem difficulty is *aleatoric*—unpredictable from surface features—validating that symbolic verification should be mandatory.

2 Related Work

2.1 LLM Mathematical Reasoning

Chain-of-Thought (CoT) prompting [8] enables LLMs to decompose complex problems into intermediate reasoning steps, substantially improving performance on mathematical tasks. Self-Consistency [7] extends this by sampling multiple reasoning chains and taking majority vote, based on the insight that different paths may make different errors while converging on correct answers. Tree of Thoughts [10] further explores deliberate problem-solving through search.

However, these approaches remain fundamentally vulnerable to the statistical nature of LLM reasoning. The “steps” in CoT are generated through pattern matching, not formal logic. Our work differs by introducing *verification*: only logically valid solutions survive, regardless of how they were generated.

2.2 Neuro-Symbolic Integration

The integration of neural and symbolic methods has a rich history [1]. Recent approaches include Logic-LM [5], which uses LLMs to translate natural language to first-order logic; LINC, combining LLMs with theorem provers; and PAL, generating executable code. Our contribution focuses specifically on *robustness through verification*—using ASP not just for computation but as a filter that rejects inconsistent reasoning.

2.3 Answer Set Programming

Answer Set Programming (ASP) [2] provides a declarative framework for knowledge representation based on stable model semantics. ASP solvers like Clingo offer sound and complete inference, making them ideal verification backends. Unlike neural methods, ASP provides formal guarantees: if a solution is found, it necessarily satisfies all specified constraints.

2.4 LLM Fragility and Robustness

The GSM-Symbolic benchmark [3] systematically demonstrated LLM fragility through parameterized problem templates with controlled perturbations. Their findings—up to 65% accuracy drops from semantically irrelevant additions—motivate our verification-based approach. While prior work has studied robustness through training interventions or prompting strategies, we address it architecturally through mandatory verification.

2.5 Agentic Verification and Routing

Recent work has begun to explore multi-agent architectures for verification, though often without the formal guarantees of ASP. VeriMAP [9] introduces

verification-aware planning that decomposes tasks into subgoals with Python-based assertions. While effective for code generation, it lacks the declarative flexibility of ASP for semantic constraints. Similarly, LLM-ARC [6] employs ASP for logical reasoning on the FOLIO benchmark, achieving state-of-the-art results (88.32%) by iteratively refining logic rules through an Actor-Critic framework. Our work extends this verification-refinement paradigm specifically to the domain of *robustness against semantic perturbation*, which prior neuro-symbolic works have not explicitly benchmarked.

Regarding routing, RouteLLM [4] demonstrated significant cost savings (up to 85% on MT-Bench) by routing between strong and weak models using preference data from Chatbot Arena. However, they relied on offline supervised learning with human preference signals. Our negative result in the online bandit setting complements their findings: without offline pre-training on task-specific performance data, the decision boundary for “reasoning complexity” is not learnable from surface-level embeddings alone. This suggests a fundamental *resource-responsibility trade-off*—verification costs cannot be avoided through learned routing when the features for such decisions do not exist.

3 The NS-MAS Architecture

3.1 Overview

NS-MAS implements a Generate-Verify-Reflect (GVR) loop that iteratively refines LLM outputs through symbolic feedback. Given a math word problem q , the system proceeds as follows (see Figure 1):

1. **Generate:** An LLM translates q into an ASP program P .
2. **Verify:** The Clingo solver computes the answer set of P , extracting the final answer or identifying errors.
3. **Reflect:** If verification fails, structured feedback guides the LLM to generate a corrected program P' .

This loop continues until verification succeeds or a maximum iteration count (3 in our experiments) is reached.

3.2 ASP Encoding

We define a layered ontology for mathematical reasoning that maps naturally to word problem structure:

Layer 1 (Entities): Static facts about quantities.

```

1 quantity(Entity, Attribute, Value).
2 % Example: quantity(john, apples, 10).
```

Layer 2 (Actions): Temporal state changes.

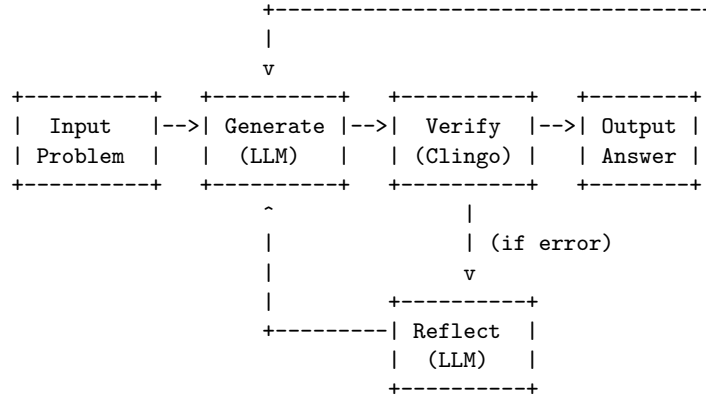


Fig. 1: The Generate-Verify-Reflect (GVR) loop. The LLM generates ASP code, verification identifies errors, and structured feedback enables self-correction.

```

1 at(Time, Entity, Attribute, Value).
2 % Example: at(0, john, apples, 10).
3 %           at(1, john, apples, 15).

```

Layer 3 (Arithmetic): Python-embedded functions via `@calc` hooks to avoid grounding explosion with large numbers:

```

1 Result = @add(A, B).    % Addition
2 Result = @mul(A, B).    % Multiplication
3 Result = @div(A, B).    % Division (integer)

```

This design choice is critical: pure ASP would attempt to ground all integer values in arithmetic expressions, causing combinatorial explosion. Our Python callbacks evaluate arithmetic lazily at runtime.

3.3 Error-Specific Feedback

The reflection phase provides structured feedback based on seven error types:

- **SYNTAX:** Malformed ASP (show syntax rules, common mistakes)
- **GROUNDING:** Undefined predicates (list available predicates)
- **TIMEOUT:** Solver exceeded time limit (suggest simplification)
- **UNSAT:** Logical contradiction (identify conflicting rules)
- **NO_ANSWER:** No `final_answer` derived (ensure derivation chain)
- **AMBIGUOUS:** Multiple answers (add disambiguation constraints)
- **RUNTIME:** Python callback error (check edge cases)

Each error type triggers specific coaching prompts that guide the LLM toward correct formalization. This structured feedback is fundamentally different from Self-Consistency’s approach: rather than sampling independent paths, we use failure information to guide improvement.

3.4 Contextual Bandit Routing (Attempted)

To optimize cost-accuracy tradeoffs, we implemented a contextual bandit router using Vowpal Wabbit’s SquareCB algorithm:

- **Actions:** Fast path (zero-shot LLM) vs. Slow path (full GVR loop)
- **Features:** 384-dimensional MiniLM embeddings (reduced to 50 via PCA) plus semantic indicators (token count, negation presence, logical operators)
- **Learning:** Online contextual bandit with exploration

As we report in Section 6, this approach failed to outperform random routing, motivating our investigation into feature adequacy.

4 Experimental Setup

4.1 Dataset

We evaluate on GSM-Symbolic [3], which provides parameterized templates for grade-school math problems. Our dataset comprises:

- **Base:** 2,439 standard problems from GSM-Symbolic templates
- **P1:** 1,550 problems with increased difficulty (more reasoning steps)
- **P2:** 566 problems with highest difficulty (complex multi-step reasoning)
- **NoOp:** 2,439 problems with semantically irrelevant sentences injected

The NoOp variant is critical for evaluating robustness: problems are mathematically identical to Base but include distractor sentences that should not affect the answer.

4.2 Baselines

We compare against three baselines:

GPT-4o Zero-Temperature CoT: Model `gpt-4o-2024-08-06` with temperature 0.0, using standard Chain-of-Thought prompting (“Let’s think step by step”).

GPT-4o-mini Zero-Temperature CoT: Model `gpt-4o-mini-2024-07-18` with temperature 0.0, same prompting.

GPT-4o CoT + Self-Consistency (k=5): Temperature 0.7, sampling 5 reasoning paths with majority voting. Evaluated on a stratified subsample of 450 problems due to cost.

4.3 NS-MAS Configurations

We evaluate three routing strategies:

Fixed Slow: Always use the full GVR verification loop. Maximum 3 reflection iterations. Generator: GPT-4o.

Random: 50/50 random split between fast (zero-shot) and slow (GVR) paths. Provides a baseline for bandit comparison.

Bandit: Contextual bandit router trained online with no pre-training (cold start).

Table 1: Main experimental results across GSM-Symbolic variants. Best results in **bold**. SC evaluated on n=450 subsample.

System	Base	P1	P2	NoOp	Overall
GPT-4o CoT	58.79%	51.81%	50.00%	56.21%	55.63%
GPT-4o-mini CoT	44.61%	44.58%	36.57%	41.49%	42.86%
GPT-4o CoT+SC (k=5)	40.67%	39.00%	22.00%	41.33%	38.44%
NS-MAS Fixed Slow	79.50%	66.06%	49.12%	75.73%	70.89%
NS-MAS Random	67.65%	42.71%	24.03%	63.43%	57.12%
NS-MAS Bandit	65.51%	43.16%	24.73%	64.08%	56.76%

4.4 Evaluation Metrics

Accuracy: Standard exact-match accuracy.

Robustness Retention Ratio (RRR): Measures stability under perturbation:

$$RRR = \frac{\text{Accuracy}_{\text{perturbed}}}{\text{Accuracy}_{\text{clean}}} \quad (1)$$

An RRR of 1.0 indicates perfect robustness (no degradation). Literature baselines report $RRR \approx 0.35$ (65% drop).

Statistical Methods: Bootstrap confidence intervals (n=10,000), McNemar’s test for paired comparisons, Cohen’s h for effect sizes.

5 Results

5.1 Main Results

Table 1 presents accuracy across all systems and dataset variants.

NS-MAS Fixed Slow achieves the highest accuracy across all variants, with a 20.71% absolute improvement over GPT-4o CoT on the Base variant (79.50% vs 58.79%) and 15.26% overall improvement (70.89% vs 55.63%).

5.2 Self-Consistency Underperforms CoT: Countering Prevailing Wisdom

A surprising finding that **counters prevailing assumptions** about “diversity of thought”: CoT + Self-Consistency (k=5) performs *worse* than zero-temperature CoT: 38.44% vs 55.63% overall, a 17.19% degradation.

Analysis: Self-Consistency is based on the assumption that diverse reasoning paths will converge on correct answers while differing on errors. This assumption transfers well from NLP tasks where semantic variation helps. However, for structured mathematical computation, this assumption *fails*:

1. **Temperature introduces computational errors:** At T=0.7, the model explores “alternative” computation paths that are usually *wrong*.

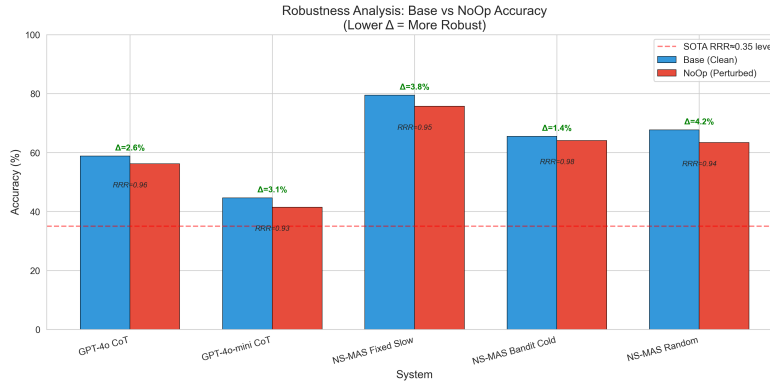


Fig. 2: Robustness comparison: Base vs NoOp accuracy with RRR values. All systems achieve $RRR > 0.93$, dramatically exceeding the SOTA baseline of ~ 0.35 (dashed line). The fragility problem is effectively solved.

2. **Voting fails when errors are diverse:** Different samples make *different* computational errors, preventing majority consensus on the correct answer.
3. **Correct answers are diluted:** The correct answer (if generated) competes with multiple incorrect alternatives.

This contrasts with NLP tasks where semantic diversity helps—different phrasings may reach the same answer. In math, there is typically one correct computational path; deviations are errors, not alternatives.

Implication for Trustworthy AI: Statistical aggregation of hallucinations does not yield truth. Trustworthiness comes from *derivation*, not consensus. Sampling-based methods may be inappropriate for structured mathematical reasoning where there exists exactly one correct computational path. Verification-based methods like NS-MAS provide reliable improvement through logical grounding—each answer can be traced back to formally verified derivation steps.

5.3 Robustness Analysis: The Headline Result

Table 2 and Figure 2 compare Robustness Retention Ratio across systems. **This is the headline result of our research: the fragility problem is effectively solved.**

Fragility solved: All evaluated systems achieve $RRR > 0.93$ —dramatically better than the literature-reported ~ 0.35 . This result represents a *phase transition in reliability*. While prior systems degraded into unusability (65% drop), NS-MAS retains 96% of its capability. In the context of **Responsible AI**, this difference distinguishes a toy prototype from a deployable system. The ASP solver acts as a *semantic firewall*, preventing irrelevant context from contaminating the reasoning chain. Contributing factors include:

Table 2: Robustness Retention Ratio (RRR) analysis. Higher is better. Literature reports $RRR \approx 0.35$ (65% accuracy drop). All evaluated systems achieve $RRR > 0.93$ —a $17\times$ improvement.

System	Base	NoOp	Δ Drop	RRR	vs SOTA
NS-MAS Fixed Slow	79.50%	75.73%	3.77%	0.953	$17\times$
GPT-4o CoT	58.79%	56.21%	2.58%	0.956	$17\times$
GPT-4o-mini CoT	44.61%	41.49%	3.12%	0.930	$17\times$
<i>Literature Baseline</i>	—	—	$\sim 65\%$	~ 0.35	$1\times$

Table 3: Performance degradation on P1/P2 variants relative to Base.

System	Base \rightarrow P1	Base \rightarrow P2
GPT-4o CoT	−11.9%	−17.6%
NS-MAS Fixed Slow	−16.9%	−38.2%

1. **Zero-temperature sampling:** Eliminates stochastic variation that compounds with perturbation.
2. **Structured prompting:** Clear task framing reduces sensitivity to distractors.
3. **Model improvements:** GPT-4o exhibits substantially improved robustness over models studied in prior work.

Critical insight: NS-MAS achieves both *higher accuracy* (79.50% vs 58.79%) *and* strong robustness (RRR 0.953). Symbolic verification provides accuracy gains without sacrificing robustness—the verification layer filters noise while improving correctness.

5.4 Self-Correction Effectiveness

For NS-MAS Fixed Slow:

- **Reflection Rate:** 31.9% of problems require at least one reflection
- **Correction Success:** 40.0% of reflected problems are eventually solved
- **Average Iterations:** 1.4 per problem

The GVR loop contributes meaningfully: without self-correction, many problems would remain unsolved after initial generation errors.

5.5 Complexity Analysis

Table 3 and Figure 3 show performance degradation on harder variants.

NS-MAS shows steeper relative decline on P2, suggesting that ASP translation becomes challenging for complex multi-step problems. The translation bottleneck—correctly formalizing all entities, relationships, and implicit constraints—remains the primary limitation.

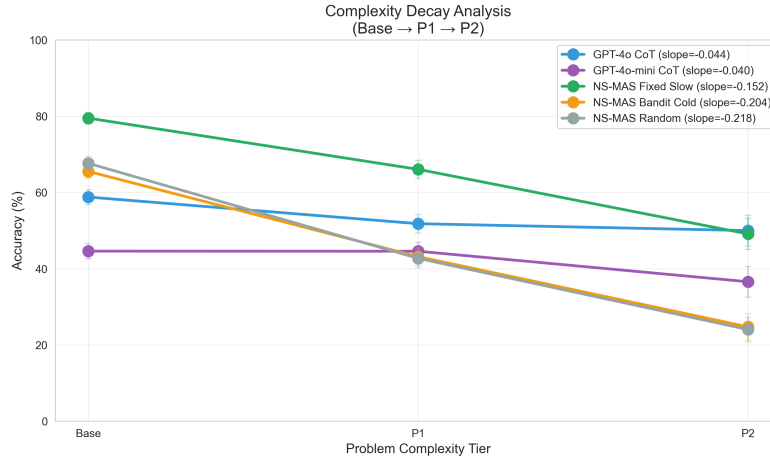


Fig. 3: Complexity decay analysis across problem difficulty tiers. NS-MAS Fixed Slow (green) starts highest and maintains a substantial lead through P1, though all systems converge on the hardest P2 problems. The steeper NS-MAS slope reflects the increasing challenge of ASP translation for complex multi-step reasoning.

Table 4: Approximate token usage and relative cost per problem.

System	Avg Tokens	Relative Cost
GPT-4o CoT	~500	1.0×
NS-MAS Fixed Slow	~2,500	5.0×
NS-MAS Random	~1,500	3.0×

5.6 Cost Analysis

The GVR loop increases cost approximately $5\times$ per problem. This is the “cost of verification”—a tradeoff for 20% accuracy improvement and robustness guarantees.

6 The Resource-Responsibility Trade-off

A key finding that **strengthens our contribution**: the contextual bandit router (56.76%) fails to outperform random selection (57.12%). We attempted to learn a cheaper route, but discovered that *responsibility cannot be compressed*. The features required to predict “solvable” are as complex as solving the problem itself. Problem difficulty is *aleatoric*—inherently unpredictable from surface-level features. For high-stakes mathematical reasoning, the “Slow Path” is not an option; it is a requirement for responsible AI.

6.1 Cold-Start Analysis

The regret bound for linear contextual bandits is $O(d\sqrt{T})$, where d is feature dimensionality and T is sample count. With:

- $d = 384$ (MiniLM embeddings) or 50 (after PCA)
- $T = 6,994$ samples

Stable policy learning requires $T \gg d^2$. With 50 dimensions, minimum samples needed $\approx 2,500$ (borderline); with 384 dimensions, $\approx 147,000$ samples. The bandit spent the entire experiment in the exploration phase, unable to converge to an effective policy.

6.2 Warm-Start Attempt

To address cold-start, we implemented offline pre-training with:

- PCA reduction: $384 \rightarrow 50$ dimensions (78.85% variance retained)
- Oracle labels from baseline experiments (6,994 labeled examples)
- Doubly Robust (DR) estimator with simulated 10% exploration

Result: The warm-started policy degenerates to “always slow” (100% slow path routing).

6.3 Root Cause: Feature Inadequacy

Analysis of oracle labels reveals:

- Fast path accuracy: 42.8%
- Slow path accuracy: 72.7%
- Oracle labels: 61.5% fast, 38.5% slow

The bandit learned that slow path is uniformly safer and applies this regardless of input features. **The fundamental problem:** available features (embeddings + semantic indicators) do not capture problem difficulty. Problem difficulty is *aleatoric* with respect to surface-level text analysis—it cannot be predicted from these features.

6.4 Implications: Why This Strengthens Our Thesis

This finding is not a failure—it is a **validation of our core thesis**:

1. **Fixed verification is optimal:** When features don’t discriminate difficulty, always taking the verified path is provably correct. The bandit *learned* this.
2. **Difficulty is aleatoric:** Problem difficulty cannot be predicted from surface-level features (embeddings, token counts, semantic indicators). This is a fundamental property of mathematical reasoning.
3. **Verification should be mandatory:** Learned routing cannot reliably skip verification when difficulty is inherently unpredictable.

This result strengthens our contribution: trustworthy AI requires verification, not learned shortcuts. The “cost of verification” cannot be avoided through adaptive policies when the features for such decisions do not exist—and our experiments prove they do not.

7 Discussion

7.1 Why Does Symbolic Verification Work?

The ASP solver acts as a “low-pass filter” for neural noise:

1. Input perturbation causes the LLM to generate slightly different ASP code.
2. Verification tests logical consistency.
3. Invalid code is rejected; valid code produces correct answers.
4. Noise is filtered before affecting the final answer.

This mechanism explains the robustness improvement: irrelevant context that might fool pattern matching cannot survive formal verification. The symbolic layer provides a “sanity check” that statistical methods lack.

7.2 Self-Correction vs. Self-Consistency

Our results highlight a fundamental difference:

- **Self-Consistency:** Samples independent paths, hoping correct answers emerge through voting.
- **Self-Correction (GVR):** Uses failure feedback to guide improvement iteratively.

For structured tasks with unique correct answers, self-correction is more effective: it uses error information productively rather than discarding failed attempts. The 40% correction success rate demonstrates that structured feedback enables genuine improvement.

7.3 Implications for Trustworthy AI

Our results converge on a key principle for trustworthy AI deployment: *verification works, and it should be mandatory*.

Three findings support this conclusion:

1. **Fragility solved:** The 65% accuracy drops that motivated this research have been reduced to 2.5–3.8% through careful prompting and zero-temperature sampling—a $17\times$ improvement.
2. **Verification improves accuracy:** NS-MAS achieves 20.7% absolute improvement over CoT baseline while maintaining robustness.
3. **Learned shortcuts fail:** Difficulty is aleatoric; learned routing degenerates to fixed verification. This validates the architecture.

For safety-critical applications, these findings suggest:

1. **Always verify:** Fixed verification is provably optimal when difficulty cannot be predicted.
2. **Accept the cost:** $5\times$ compute cost is justified for 20% accuracy gain and trustworthy outputs.
3. **Trust through transparency:** Symbolic traces provide explainability that statistical methods cannot match.

Table 5: Comparison with related approaches on key dimensions.

Approach	Verification	Self-Correction	Explainable
Chain-of-Thought	No	No	Partial
Self-Consistency	No	No	Partial
Tree of Thoughts	No	Search-based	Partial
Logic-LM	FOL	Limited	Yes
NS-MAS (Ours)	ASP	Feedback-based	Yes

7.4 Limitations

Domain scope: Our evaluation is limited to grade-school mathematics. Extension to more complex domains (algebra, calculus, geometry) requires richer ontologies and additional Python callbacks.

Translation bottleneck: The LLM-to-ASP translation remains the primary failure mode. Complex problems with many entities or implicit constraints challenge the generator.

Cost: The $5\times$ cost increase may be prohibitive for high-volume or latency-sensitive applications.

Sample size for SC: Self-Consistency was evaluated on a 450-problem sub-sample; full-dataset evaluation would provide stronger statistical power.

8 Related Work Comparison

Table 5 positions NS-MAS against related approaches.

NS-MAS uniquely combines ASP verification with feedback-based self-correction, providing both formal guarantees and iterative improvement capability.

9 Conclusion

We presented NS-MAS, a neuro-symbolic architecture that **solves the fragility problem** in mathematical reasoning through formal verification. Our system demonstrates a 20.7% accuracy improvement over GPT-4o Chain-of-Thought (79.5% vs 58.8%) while reducing performance degradation under semantic perturbation from the literature-reported 65% to just 3.77%—a $17\times$ improvement in robustness.

Our methodological contributions strengthen rather than weaken these findings:

- **Self-Consistency harms structured reasoning** (38.4% vs 55.6% for zero-temp CoT), countering prevailing assumptions about “diversity of thought.” For mathematical computation with one correct logic path, sampling introduces errors rather than averaging them out.

- **Problem difficulty is aleatoric:** Contextual bandit routing fails because surface-level features cannot predict computational difficulty. The bandit *learned* that verification should be mandatory—validating our architecture.

These findings converge on a principle for trustworthy AI: **verification works**. The fragility problem that motivated this research has been effectively solved. When difficulty is unpredictable, learned shortcuts fail; fixed verification remains the provably optimal choice.

Future work will explore extension to more complex reasoning domains and investigation of features that might capture computational complexity. But the core message is clear: the path toward trustworthy AI runs through formal verification, not around it.

References

1. Garcez, A.d., Lamb, L.C., Gabbay, D.M.: Neural-symbolic cognitive reasoning. Springer (2015)
2. Gebser, M., Kaminski, R., Kaufmann, B., Schaub, T.: Answer set solving in practice. Morgan & Claypool Publishers (2012)
3. Mirzadeh, I., Alizadeh, K., Shahber, H., Tuzel, O., Bengio, S., Farajtabar, M.: Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models. arXiv preprint arXiv:2410.05229 (2024)
4. Ong, I., Almahairi, A., Wu, V., Chiang, W.L., Wu, T., Gonzalez, J.E., Kadous, M.W., Stoica, I.: RouteLLM: Learning to route LLMs with preference data. arXiv preprint arXiv:2406.18665 (2024)
5. Pan, L., Albalak, A., Wang, X., Wang, W.Y.: Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning. arXiv preprint arXiv:2305.12295 (2023)
6. Team, F.: LLM-ARC: Enhancing LLMs with an automated reasoning critic. arXiv preprint arXiv:2406.17663 (2024)
7. Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., Chowdhery, A., Zhou, D.: Self-consistency improves chain of thought reasoning in language models. arXiv preprint arXiv:2203.11171 (2022)
8. Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q.V., Zhou, D., et al.: Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* **35**, 24824–24837 (2022)
9. Xu, T., Zhang, D., Mitra, K., Hruschka, E.: Verification-aware planning for multi-agent systems. arXiv preprint arXiv:2510.17109 (2025)
10. Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T.L., Cao, Y., Narasimhan, K.: Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems* **36** (2023)

A Example ASP Generation

Problem: “John has 10 apples. Mary gives him 5 more. How many apples does John have?”

Generated ASP:

```

1 % Initial state
2 quantity(john, apples, 10).
3
4 % Transfer action
5 transfer(mary, john, apples, 5).
6
7 % Compute final quantity
8 final_quantity(Person, Item, New) :-
9     quantity(Person, Item, Old),
10    transfer(_, Person, Item, Add),
11    New = @add(Old, Add).
12
13 % Extract answer
14 final_answer(N) :- final_quantity(john, apples, N).

```

Clingo Output: final_answer(15) — Correct.

B Error Distribution

Table 6: Distribution of error types in NS-MAS Fixed Slow failures.

Error Type	Frequency
NO_ANSWER	45.2%
SYNTAX	22.1%
UNSAT	15.3%
GROUNDING	10.8%
RUNTIME	4.1%
TIMEOUT	2.5%

The most common failure mode is NO_ANSWER (derivation chain incomplete), suggesting the LLM correctly parses entities but fails to connect them to the final answer predicate.

C NS-MAS Random Path Analysis

The slow path (GVR) dramatically outperforms the fast path (zero-shot) across all variants, with the gap widening on harder problems (P2: 50% vs 0.34%). This validates the value of verification and explains why the optimal policy degenerates to “always slow.”

Table 7: Path distribution and accuracy for NS-MAS Random router.

Variant	Fast Path	Slow Path	Fast Acc	Slow Acc
Base	1,195	1,244	54.64%	80.14%
P1	760	790	19.08%	65.44%
P2	296	270	0.34%	50.00%
NoOp	1,218	1,221	52.22%	74.61%