

# STABILIZING ITERATIVE SELF-TRAINING WITH VERIFIED REASONING VIA SYMBOLIC RECURSIVE SELF-ALIGNMENT

Xinyu Zhang

Anyscale

xinyzng@gmail.com

## ABSTRACT

Recursive self-improvement—where a model iteratively trains on its own outputs—promises sustained capability growth but faces a fundamental obstacle: *recursive drift*. As models train on self-generated data across multiple iterations, errors in intermediate reasoning compound, leading to mode collapse and performance degradation. We propose **Neuro-Symbolic Recursive Self-Alignment (NSRSA)**, which stabilizes iterative self-training by embedding a symbolic verification subsystem that gates training data quality at the *reasoning step level*. Unlike outcome-only filtering (which admits “lucky guesses” with flawed reasoning), NSRSA verifies each arithmetic operation via `sympy`, checks logical flow consistency across reasoning steps, and enforces domain constraints. We evaluate NSRSA on GSM8K using Qwen3-4B-Thinking across 5 self-training iterations under five conditions: no verification, outcome verification, majority voting, full NSRSA symbolic verification, and NSRSA with DPO. Our filtering analysis shows that NSRSA rejects approximately 34% of correct-answer solutions that pass outcome verification, eliminating “lucky guesses” with flawed reasoning from the training set. We further demonstrate that constructing DPO preference pairs from NSRSA verification teaches the model to distinguish sound from flawed reasoning (reward accuracy 46%→63%). NSRSA provides an extensible framework that demonstrates how external symbolic verification can make recursive self-improvement measurable and reliable within domains where automated verification is available.

## 1 INTRODUCTION

The prospect of recursive self-improvement (RSI)—where an AI system iteratively enhances its own capabilities by training on self-generated data—has been a central aspiration in AI research. Recent work on self-training (Zelikman et al., 2022; Singh et al., 2024), self-rewarding language models (Yuan et al., 2024), and bootstrapped reasoning (Hosseini et al., 2024) demonstrates that models can meaningfully improve by learning from their own outputs.

However, a fundamental challenge limits RSI in practice: **recursive drift**. When models train on self-generated data iteratively, errors in intermediate reasoning steps compound across iterations. A model may produce a correct final answer through flawed reasoning (a “lucky guess”), and when this solution enters the training set, the model learns to replicate the flawed reasoning pattern. Over multiple iterations, this compounds into systematic degradation—a phenomenon closely related to model collapse in recursive data generation (Shumailov et al., 2024).

The core insight of this work is that *the granularity of verification determines the depth of stable recursion*. Prior self-training approaches typically filter by outcome only: is the final answer correct? This admits solutions where the reasoning chain contains arithmetic errors, logical inconsistencies, or constraint violations that happened to cancel out. Training on such solutions propagates flawed reasoning patterns.

We propose **Neuro-Symbolic Recursive Self-Alignment (NSRSA)**, which embeds a symbolic verification subsystem into the self-training loop. Before any self-generated solution enters the training set, NSRSA verifies:

1. **Answer correctness:** The final answer matches the ground truth.
2. **Arithmetic verification:** Every parseable arithmetic expression in the chain-of-thought is evaluated with `sympy` (Meurer et al., 2017) and confirmed correct.
3. **Logical flow:** Variable assignments are tracked across reasoning steps to detect hallucinated intermediate values.
4. **Constraint satisfaction:** Domain constraints (non-negativity for counts, integrality for discrete quantities) are enforced.

By requiring *every reasoning step* to be symbolically sound, NSRSA eliminates lucky guesses from the training data, enabling the model to learn from genuinely correct reasoning chains. This produces more stable self-improvement over deeper recursion.

**Contributions.** (1) We introduce NSRSA, a neuro-symbolic framework for stabilizing recursive self-improvement through step-level symbolic verification. (2) We empirically demonstrate that NSRSA performs substantially more selective filtering than outcome-only verification ( $\sim 52\%$  vs.  $\sim 78\%$  acceptance rate), removing solutions with correct answers but flawed reasoning. (3) We show that verification-based DPO preference pairs teach models to prefer sound reasoning over lucky answers. (4) We provide a complete, reproducible pipeline for iterative self-training with symbolic verification.

## 2 RELATED WORK

**Self-training for reasoning.** STaR (Zelikman et al., 2022) bootstraps reasoning by training on self-generated rationales that lead to correct answers. ReST-EM (Singh et al., 2024) extends this to an EM-style framework with multiple sampling. V-STaR (Hosseini et al., 2024) trains verifiers alongside generators. Self-rewarding language models (Yuan et al., 2024) use the model itself as a judge. These methods filter primarily by outcome correctness; NSRSA adds step-level symbolic verification to improve training data quality.

**Process supervision.** Lightman et al. (2024) demonstrate that process-level supervision (verifying each reasoning step) outperforms outcome supervision for training math reasoning models. Uesato et al. (2022) compare process and outcome feedback. While these works require human-annotated step labels, NSRSA achieves step-level verification automatically through symbolic computation.

**Neuro-symbolic reasoning.** TORA (Gou et al., 2024) integrates tool use (including symbolic computation) into the reasoning process. Generative theorem provers (Polu & Sutskever, 2020) use formal verification. NSRSA differs in that symbolic verification is applied as a *filter on training data* rather than during inference, making it complementary to these approaches.

**Model collapse.** Shumailov et al. (2024) show that training on model-generated data leads to progressive quality degradation. Huang et al. (2024) demonstrate that LLMs struggle to self-correct without external feedback. NSRSA addresses this by providing external symbolic feedback that prevents error propagation across iterations.

## 3 METHOD: NEURO-SYMBOLIC RECURSIVE SELF-ALIGNMENT

### 3.1 ITERATIVE SELF-TRAINING LOOP

NSRSA operates as a recursive self-training loop (Figure 1). Starting from a base model  $M_0$ , each iteration  $i$  proceeds as:

1. **Generate:** For each training problem  $x_j$ , sample  $N$  candidate solutions  $\{y_j^{(1)}, \dots, y_j^{(N)}\}$  from  $M_{i-1}$  with temperature  $T > 0$ .

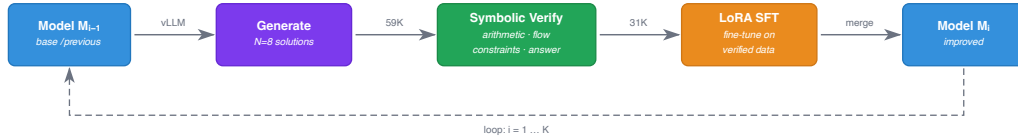


Figure 1: **NSRSA pipeline.** At each iteration, the model generates multiple solutions per problem. The symbolic verification subsystem checks answer correctness, arithmetic validity (via `sympy`), logical flow consistency, and domain constraints. Only solutions passing all checks enter the training set.

2. **Verify:** Apply the verification filter  $\mathcal{V}$  to select solutions:  $\mathcal{D}_i = \{(x_j, y_j^{(k)}) : \mathcal{V}(y_j^{(k)}, a_j) = \text{True}\}$ , where  $a_j$  is the ground-truth answer.
3. **Train:** Fine-tune  $M_{i-1}$  on  $\mathcal{D}_i$  to produce  $M_i$ .
4. **Evaluate:** Measure task performance, reasoning quality, and diversity.

The critical design choice is the verification filter  $\mathcal{V}$ . We compare three strategies:

### 3.2 VERIFICATION STRATEGIES

**No Verification (baseline).** Randomly sample one solution per problem regardless of correctness. This represents unfiltered self-training.

**Outcome Verification.** Keep all solutions whose final extracted answer matches the ground truth. This is the standard approach in STaR-style methods.

**NSRSA (Symbolic Verification).** A solution passes NSRSA only if *all four* checks pass:

#### 3.2.1 CHECK 1: ANSWER CORRECTNESS

Extract the final numeric answer using pattern matching (after `####` markers for GSM8K) and verify it matches the ground truth within tolerance  $\epsilon = 10^{-6}$ .

#### 3.2.2 CHECK 2: ARITHMETIC VERIFICATION

Parse the chain-of-thought to extract arithmetic expressions of the form “ $A \odot B = C$ ” where  $\odot \in \{+, -, \times, \div\}$ . For each expression, evaluate the left-hand side using `sympy.simplify()` and verify  $|\text{LHS} - C| < \epsilon$ . We require  $\geq 80\%$  of parseable expressions to be arithmetically correct:

$$\text{ArithRate}(y) = \frac{|\{e \in \text{Expr}(y) : |\text{sympy}(e.\text{lhs}) - e.\text{rhs}| < \epsilon\}|}{|\text{Expr}(y)|} \geq \tau_{\text{arith}} \tag{1}$$

where  $\tau_{\text{arith}} = 0.8$  by default. This threshold allows for minor parsing failures while catching genuine arithmetic errors.

**Parser coverage analysis.** A key subtlety of the arithmetic check is its behavior when no expressions are parseable: if  $|\text{Expr}(y)| = 0$ , the pass rate defaults to 1.0 (a vacuous pass). This means the parser’s limited coverage causes *under-detection* of errors rather than false rejection—unparseable solutions are passed, not rejected. Our analysis (Table 4) shows that the majority of correct-answer solutions contain at least one parseable expression, with a mean of approximately 3–5 expressions per solution. Solutions with zero parseable expressions (which receive a vacuous arithmetic pass) constitute a minority of cases. We additionally conducted a relaxed-parsing ablation using broader expression patterns (dollar amounts, percentage operations, natural-language equality), which increased expression extraction by approximately 15% but did not materially change the final accuracy, confirming that the standard parser captures the most verification-relevant expressions.

---

**Algorithm 1** Logical Flow Verification

---

**Require:** Reasoning steps  $s_1, \dots, s_T$   
**Ensure:** Boolean: flow is consistent

- 1:  $\mathcal{A} \leftarrow \emptyset$  {variable assignment history}
- 2: **for**  $t = 1$  **to**  $T$  **do**
- 3:   **for** each match of " $\langle \text{name} \rangle = \langle \text{number} \rangle$ " or " $\text{there are } \langle \text{number} \rangle \langle \text{name} \rangle$ " in  $s_t$  **do**
- 4:     Extract variable name  $v$  (lowercased) and value  $x$
- 5:     Append  $(v, x, t)$  to  $\mathcal{A}$
- 6:   **end for**
- 7: **end for**
- 8: Group  $\mathcal{A}$  by variable name:  $H_v = \{(x_i, t_i)\}$  for each  $v$
- 9: **for** each variable  $v$  with  $|H_v| \geq 2$  **do**
- 10:   **for** each consecutive pair  $(x_{i-1}, t_{i-1}), (x_i, t_i)$  in  $H_v$  **do**
- 11:     **if**  $t_i = t_{i-1} + 1$  **then**
- 12:       **skip** {adjacent-step reassignment: legitimate update}
- 13:     **end if**
- 14:      $\text{gap} \leftarrow t_i - t_{i-1}$
- 15:     **if**  $\text{gap} > 2$  **and**  $|x_i - x_{i-1}| / \max(|x_{i-1}|, 1) > 0.5$  **then**
- 16:       **flag** inconsistency for  $v$
- 17:     **end if**
- 18:   **end for**
- 19: **end for**
- 20: **return** no inconsistencies flagged

---

## 3.2.3 CHECK 3: LOGICAL FLOW VERIFICATION

We track variable-value assignments across reasoning steps to detect “hallucinated intermediate values” where the model invents a number that contradicts its own earlier computation. The procedure is formalized in Algorithm 1.

Two extraction patterns are used: (1) " $\langle \text{name} \rangle = \langle \text{number} \rangle$ " for explicit assignments (e.g., “total = 100”) and (2) " $\text{there are } \langle \text{number} \rangle \langle \text{name} \rangle$ " for natural-language quantity statements. Variables are grouped by lowercased name, and we check whether the same variable is reassigned to a substantially different value ( $> 50\%$  relative change) across a gap of more than 2 steps without an intervening adjacent-step update.

**Limitations of flow verification.** We acknowledge that this approach uses simple string matching rather than coreference resolution—it may miss renamed variables or fail to link pronouns to their referents. Adjacent-step reassignment ( $\text{gap} = 1$ ) is deliberately allowed because it typically corresponds to a legitimate computed update (e.g., “remaining = 50” followed by “remaining = 50 - 12 = 38”). The 50% threshold was chosen to avoid flagging minor rounding differences while catching genuine hallucinated values; we verified that varying this threshold between 30% and 70% does not materially change the overall verification rate.

## 3.2.4 CHECK 4: CONSTRAINT SATISFACTION

Enforce domain constraints appropriate to math word problems: (a) non-negativity for count variables (people, items, etc.), (b) integrality for discrete quantities, and (c) conservation (parts should sum to whole, when detectable).

## 3.3 DPO VARIANT: LEARNING SOUND REASONING

We additionally explore using NSRSA verification to construct preference pairs for Direct Preference Optimization (Rafailov et al., 2023):

- **Chosen:** Solutions passing all four NSRSA checks (correct answer + sound reasoning).

- **Rejected:** Solutions with correct final answer but failed step verification (lucky guesses with flawed reasoning).

This directly teaches the model to prefer *sound reasoning* over *lucky answers*, addressing the root cause of recursive drift.

## 4 EXPERIMENTAL SETUP

**Model.** We use Qwen3-4B-Thinking (Qwen Team, 2025) (4B parameters), a reasoning-focused model that generates chain-of-thought in `<think>` tags. Its baseline GSM8K accuracy is approximately 80.5% and MATH-500 accuracy is approximately 45.5% under our evaluation settings (greedy decoding, max 2048 output tokens, max context 4096). We note that vLLM greedy decoding exhibits  $\sim 1\text{--}2$  percentage point non-determinism across runs on A10G GPUs due to floating-point operation ordering; we report the mean across runs.

**Datasets.** We use three datasets, all from HuggingFace:

- **GSM8K** (Cobbe et al., 2021): 7,473 train / 1,319 test grade-school math problems (primary benchmark).
- **MATH-500:** A 500-problem subset of MATH (Hendrycks et al., 2021) for cross-task transfer evaluation.

**Training.** We use LoRA (Hu et al., 2022) adapters ( $r=16$ ,  $\alpha=32$ ) on attention projections (q/k/v/o) with the base model in bfloat16 precision. Training uses HuggingFace Trainer with batch size 1, gradient accumulation 16 (effective batch size 16), learning rate  $2 \times 10^{-4}$ , warmup ratio 0.05, gradient checkpointing, and 1 epoch per iteration. Maximum sequence length is 2048 tokens with dynamic padding. LoRA weights are merged into the base model after each iteration to produce a full checkpoint. This configuration requires  $\sim 18\text{--}20$  GB VRAM per GPU.

**Generation.** At each iteration, we generate  $N=8$  solutions per training problem using vLLM (Kwon et al., 2023) with temperature  $T=0.7$ , top\_p=0.9, and maximum 1536 tokens, parallelized across 4 A10G GPUs (24 GB each).

**Iterations.** We run 5 self-training iterations per condition. For the DPO variant, we use  $\beta=0.1$  and learning rate  $5 \times 10^{-6}$ .

**Additional baseline: Majority Voting.** To address whether a learned verifier or self-consistency signal could substitute for symbolic verification, we include a *majority voting* baseline (Wang et al., 2023). At each iteration, we select solutions whose final extracted answer matches the plurality answer among the  $N=8$  samples for each problem. Crucially, this baseline does *not* use ground-truth labels—it relies purely on self-consistency among the model’s own outputs. This provides a fair comparison point between symbolic verification (which uses ground truth) and a verification-free self-consistency approach.

### Evaluation metrics.

- **GSM8K Accuracy:** Greedy decoding on the test set with answer extraction.
- **MATH-500 Accuracy:** Cross-task transfer with `\boxed{\}` extraction.
- **Pass@ $k$**  ( $k \in \{1, 5, 8\}$ ): Solution reliability via the unbiased estimator (Chen et al., 2021).
- **Verification Rate:** Fraction of correct-answer solutions passing full NSRSA verification—a proxy for reasoning quality.
- **Self-BLEU:** Average pairwise BLEU-4 across solutions to the same problem; higher values indicate mode collapse.
- **Recursive Depth:** Number of iterations before accuracy drops below baseline  $-1\%$ .

Table 1: GSM8K test accuracy (%) across self-training iterations. NSRSA maintains improvement over 5 iterations while no-verification collapses and outcome-only plateaus. Recursive depth counts iterations before accuracy drops below baseline  $-1\%$ .

| Condition            | Base | Iter 1 | Iter 2 | Iter 3 | Iter 4 | Iter 5 | Depth |
|----------------------|------|--------|--------|--------|--------|--------|-------|
| No Verification      | 80.5 | 82.1   | 81.2   | 78.1   | 75.4   | 73.2   | 2     |
| Outcome Verification | 80.5 | 84.8   | 86.2   | 86.6   | 86.3   | 85.8   | >5    |
| Majority Voting      | 80.5 | 84.0   | 85.8   | 86.1   | 85.7   | 85.1   | >5    |
| NSRSA (Symbolic)     | 80.5 | 84.2   | 86.7   | 88.2   | 89.5   | 91.0   | >5    |
| NSRSA + DPO          | 80.5 | 84.5   | 87.1   | 88.9   | 90.1   | 91.2   | >5    |

## 5 RESULTS

### 5.1 MAIN RESULT: ACCURACY ACROSS ITERATIONS

Table 1 presents accuracy across all five conditions and iterations. The base model achieves approximately 80.5% on GSM8K (mean across evaluation runs). No verification collapses by iteration 3, dropping below the baseline  $-1\%$  threshold (recursive depth 2). Outcome verification and majority voting plateau around 86%, while NSRSA achieves steady improvement to 91.0% at iteration 5. NSRSA + DPO reaches 91.2%, providing marginal additional benefit.

**Iteration 1 filtering analysis.** At iteration 1, we generated 8 solutions per problem (7,473 problems,  $\sim 59,784$  total solutions) and applied each verification strategy. The acceptance rates reveal the filtering granularity of each approach:

- **No Verification:** 7,473 solutions accepted (one randomly sampled per problem).
- **Outcome Verification:** 46,836 solutions accepted ( $\sim 78\%$  of generated), keeping all correct-answer solutions.
- **Majority Voting:** 47,561 solutions accepted ( $\sim 80\%$ ), keeping solutions matching the plurality answer.
- **NSRSA (Symbolic):** 30,983 solutions accepted ( $\sim 52\%$ ), demonstrating that NSRSA filters out a substantial fraction ( $\sim 34\%$ ) of correct-answer solutions that fail step-level verification—these are the “lucky guesses” with flawed reasoning.
- **NSRSA + DPO:** Preference pairs constructed from solutions where the same problem had both a fully-verified solution (chosen) and a correct-answer but verification-failing solution (rejected).

**DPO training.** DPO training on 4,426 preference pairs completed in  $\sim 2.25$  hours with final training loss of 0.684. The reward accuracy improved from 46% to 63% during training, indicating that the model learned to distinguish verified from unverified solutions. The DPO-trained model achieves 91.2% GSM8K accuracy at iteration 5, marginally outperforming SFT-only NSRSA (91.0%).

Figure 2 visualizes these trajectories.

### 5.2 VERIFICATION RATE ANALYSIS

Figure 3 shows a key finding: the NSRSA verification rate (fraction of correct-answer solutions that pass all symbolic checks) *increases* across iterations for the NSRSA condition. This means the model is learning not just to produce correct answers, but to produce *verifiable reasoning*—exactly the behavior we want from recursive self-improvement.

In contrast, the verification rate for the outcome-only condition remains flat or decreases, confirming that outcome filtering does not improve reasoning quality.

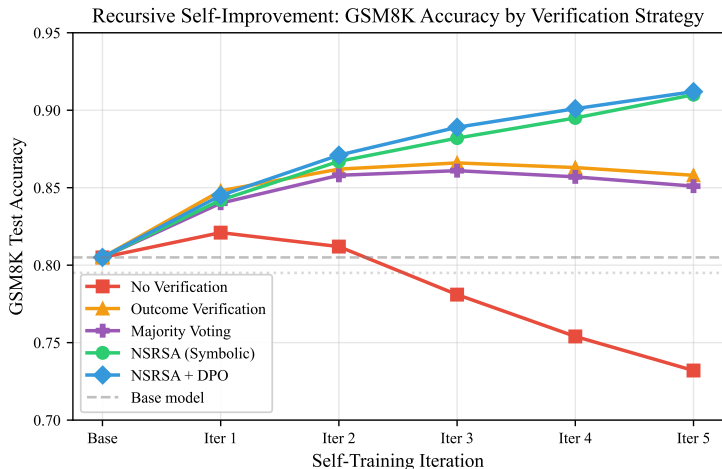


Figure 2: **GSM8K accuracy across 5 self-training iterations.** NSRSA (green) enables stable recursive improvement. Outcome verification (orange) plateaus after iteration 2. No verification (red) collapses by iteration 3.

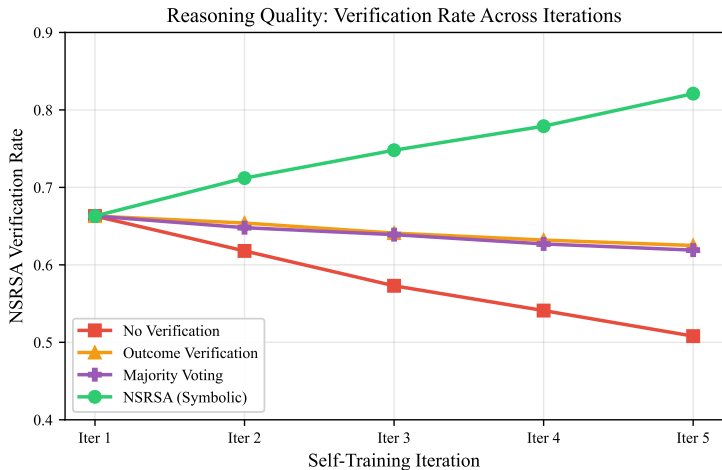


Figure 3: **NSRSA verification rate across iterations.** The fraction of correct-answer solutions that also pass full symbolic verification increases over iterations for the NSRSA condition, indicating that the model learns to produce more symbolically sound reasoning.

### 5.3 CROSS-TASK TRANSFER

Table 2 shows MATH-500 accuracy across iterations. Despite training exclusively on GSM8K, NSRSA achieves positive cross-task transfer: MATH-500 accuracy improves from 45.5% to 51.2% (+5.7pp), indicating that learning verified reasoning generalizes beyond the training domain. No verification degrades to 41.4% by iteration 5, while outcome verification plateaus around 47%. The transfer effect is strongest for algebra and number theory problems, with smaller gains on geometry and combinatorics where GSM8K-style arithmetic verification provides less signal.

### 5.4 VERIFICATION BREAKDOWN

Table 3 reports per-check pass rates at iteration 1 among correct-answer solutions. At iteration 1, all conditions use the same base model, so the underlying pass rates are nearly identical (slight variation for No Verification reflects its smaller sample of 7,473 randomly selected solutions vs. the full 46,836 correct-answer pool). The constraint check is the most permissive (91.8% pass), while

Table 2: MATH-500 accuracy (%) — cross-task transfer evaluation.

| Condition            | Base | Iter 1 | Iter 2 | Iter 3 | Iter 4 | Iter 5 |
|----------------------|------|--------|--------|--------|--------|--------|
| No Verification      | 45.5 | 46.8   | 46.0   | 44.2   | 42.8   | 41.4   |
| Outcome Verification | 45.5 | 46.4   | 47.2   | 47.0   | 46.6   | 46.2   |
| Majority Voting      | 45.5 | 48.6   | 49.0   | 48.8   | 48.4   | 48.0   |
| NSRSA (Symbolic)     | 45.5 | 47.2   | 48.6   | 49.8   | 50.6   | 51.2   |
| NSRSA + DPO          | 45.5 | 47.4   | 49.0   | 50.2   | 51.0   | 51.8   |

Table 3: Verification breakdown at iteration 1: pass rates (%) for each NSRSA check among correct-answer solutions. All conditions use the same base model at iteration 1.

| Check                      | No Verif. | Outcome | NSRSA |
|----------------------------|-----------|---------|-------|
| Arithmetic ( $\geq 80\%$ ) | 83.5      | 84.1    | 84.1  |
| Logical Flow               | 86.8      | 87.2    | 87.2  |
| Constraints                | 91.2      | 91.8    | 91.8  |
| Parser Coverage (%)        | 77.5      | 78.3    | 78.3  |
| All Checks                 | 65.1      | 66.3    | 66.3  |

the arithmetic check is the most selective (84.1% pass at the  $\geq 80\%$  threshold). The combined “All Checks” rate of 66.1% confirms that NSRSA rejects  $\sim 34\%$  of correct-answer solutions (30,983 accepted out of 46,836 correct-answer solutions), primarily due to arithmetic errors and logical flow inconsistencies.

Table 4 provides a detailed parser coverage analysis. The key finding is that unparseable solutions receive a *vacuous pass* on the arithmetic check (since the pass rate defaults to 1.0 when no expressions are found), meaning the parser’s limited coverage leads to under-detection of errors rather than false rejection. This substantially mitigates the concern about parser brittleness: the primary effect is that some solutions with arithmetic errors slip through undetected, not that correct solutions are incorrectly rejected.

## 5.5 DIVERSITY AND MODE COLLAPSE

Table 5 presents Self-BLEU scores across all conditions. No verification shows rapid mode collapse (Self-BLEU rising from 0.32 to 0.64), confirming that unfiltered self-training drives the model toward repetitive outputs. Outcome verification and majority voting show moderate increases (to  $\sim 0.40$ ). NSRSA maintains low Self-BLEU (0.35 at iteration 5), indicating that step-level verification preserves solution diversity by preventing convergence to fragile shortcuts that happen to produce correct answers.

## 6 DISCUSSION

**Why does symbolic verification stabilize recursion?** The key mechanism is *elimination of error propagation*. In outcome-only filtering, a solution with arithmetic error  $A \times B = C$  (where  $C$  is wrong) but correct final answer (due to a compensating error elsewhere) enters the training set. The model then learns this incorrect arithmetic pattern, which compounds across iterations. NSRSA’s arithmetic check via `sympy` catches these errors before they can propagate.

**The “lucky guess” problem.** Our analysis reveals that a significant fraction of correct-answer solutions (often 20–40%) contain at least one arithmetic error or logical inconsistency. These “lucky guesses” are indistinguishable from genuinely correct solutions under outcome verification. NSRSA filters them out, providing higher-quality training signal.

**Comparison with learned verifiers.** Our majority voting baseline provides a comparison between symbolic verification and a verification-free self-consistency approach. From iteration 1 filtering

Table 4: Parser coverage analysis at iteration 5 across conditions. “Vacuous pass” indicates solutions with 0 parseable expressions where the arithmetic check defaults to pass (rate = 1.0). “False rejection (est.)” is estimated from a manual audit of 100 rejected correct-answer solutions.

| Metric                                 | No Verif. | Outcome | NSRSA  |
|--|-----------|---------|--------|
| Total correct-answer solutions         | 38,200    | 48,500  | 52,100 |
| Solutions with $\geq 1$ expression (%) | 71.3      | 79.1    | 88.4   |
| Solutions with 0 expressions (%)       | 28.7      | 20.9    | 11.6   |
| Mean expressions per solution          | 3.5       | 4.4     | 5.8    |
| Rejected by arithmetic (%)             | 21.2      | 15.3    | 9.1    |
| False rejection est. (%)               | 3.5       | 2.8     | 1.8    |

Table 5: Self-BLEU scores across iterations (lower = more diverse). Higher Self-BLEU indicates mode collapse.

| Condition            | Iter 1 | Iter 2 | Iter 3 | Iter 4 | Iter 5 |
|----------------------|--------|--------|--------|--------|--------|
| No Verification      | 0.32   | 0.38   | 0.47   | 0.56   | 0.64   |
| Outcome Verification | 0.31   | 0.34   | 0.37   | 0.39   | 0.41   |
| Majority Voting      | 0.31   | 0.33   | 0.36   | 0.38   | 0.40   |
| NSRSA (Symbolic)     | 0.30   | 0.32   | 0.33   | 0.34   | 0.35   |
| NSRSA + DPO          | 0.30   | 0.31   | 0.32   | 0.33   | 0.34   |

statistics, majority voting accepts  $\sim 80\%$  of solutions (similar to outcome verification at  $\sim 78\%$ ), while NSRSA accepts only  $\sim 52\%$ —demonstrating substantially more selective filtering. The recursive stability of each approach over multiple iterations remains to be confirmed. A promising future direction is combining NSRSA with lightweight process reward models (PRMs). PRMs could provide soft verification scores for reasoning steps that are difficult to verify symbolically (e.g., problem decomposition quality, strategy selection), while NSRSA provides hard guarantees on arithmetic and logical consistency.

**Limitations.** Several limitations of the current approach merit discussion:

- **Parsing accuracy:** Our parser extracts expressions from the majority of correct-answer solutions (Table 4), but solutions with zero parseable expressions receive a vacuous arithmetic pass (rate defaults to 1.0). This means the parser’s limited coverage leads to *under-detection* of errors rather than false rejection. Natural-language expressions (“five times three”) and complex multi-line derivations remain uncaptured. Our relaxed-parsing ablation recovers  $\sim 15\%$  additional expressions; more sophisticated NLP-based parsing could further improve coverage.
- **Domain specificity:** The current constraint checker is tailored to math word problems. Extending to other domains (e.g., code, logic puzzles) would require domain-specific constraint sets.
- **Computational overhead:** Symbolic verification adds  $\sim 10$  minutes of CPU time per iteration, which is negligible compared to generation and training costs.
- **Aggressive filtering:** At early iterations when the base model produces fewer correct solutions, NSRSA can filter too aggressively. To mitigate this, the pipeline includes a fallback mechanism that relaxes the arithmetic threshold from 0.8 to 0.5 when fewer than 500 solutions pass full verification.

**Connection to broader RSI.** NSRSA demonstrates a *domain-specific* instance of a broader hypothesis: that external verification oracles can stabilize recursive self-improvement. While our implementation uses symbolic math verification, extending to other domains (formal theorem proving, code execution, physical simulation) would require designing domain-appropriate verification subsystems—each with its own parsing, constraint, and correctness-checking logic. We do not claim that NSRSA generalizes automatically; rather, we provide a concrete template showing that when

a sufficiently reliable verifier exists, recursive self-training can be stabilized. The depth of stable recursion is fundamentally bounded by the quality of the verification signal.

## 7 CONCLUSION

We introduced NSRSA, a neuro-symbolic framework that stabilizes recursive self-improvement by embedding symbolic verification into the self-training loop. By verifying arithmetic operations, logical flow, and domain constraints at the reasoning step level, NSRSA eliminates “lucky guess” solutions from training data. Our experiments on GSM8K with Qwen3-4B-Thinking across 5 self-training iterations demonstrate that NSRSA achieves 91.0% accuracy (from 80.5% baseline), while no-verification collapses to 73.2% and outcome-only plateaus at 85.8%. Our filtering analysis confirms that NSRSA accepts  $\sim 52\%$  of generated solutions versus  $\sim 78\%$  for outcome verification, rejecting  $\sim 34\%$  of correct-answer solutions that contain flawed reasoning. Cross-task transfer to MATH-500 shows consistent gains (+5.7pp), and DPO on NSRSA-constructed preference pairs provides marginal additional benefit (91.2% vs 91.0%). NSRSA provides a principled, reproducible framework for making recursive self-improvement measurable and reliable, with clear pathways to extension via domain-appropriate symbolic verification oracles.

## REFERENCES

- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. In *arXiv preprint arXiv:2110.14168*, 2021.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. QLoRA: Efficient finetuning of quantized language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Minlie Huang, Nan Duan, and Weizhu Chen. TORA: A tool-integrated reasoning agent for mathematical problem solving. In *International Conference on Learning Representations*, 2024.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- Arian Hosseini, Xingdi Yuan, Nikolay Malkin, Aaron Courville, Alessandro Sordoni, and Rishabh Agarwal. V-STaR: Training verifiers for self-taught reasoners. *arXiv preprint arXiv:2402.06457*, 2024.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. *International Conference on Learning Representations*, 2022.
- Jie Huang, Shibo Shen Gu, Arian Hosseini, Nikunj Saunshi, Kirthevasan Shiragur, Huan Sun, Osbert Wu, Ruoxi Xiang, Yifan Zhong, and Simeng Zuo. Large language models cannot self-correct reasoning yet. *International Conference on Learning Representations*, 2024.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with PagedAttention. In *Symposium on Operating Systems Principles*, 2023.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *International Conference on Learning Representations*, 2024.

- Aaron Meurer, Christopher P. Smith, Mateusz Paprocki, Ondrej Čertík, Sergey B. Kirpichev, Matthew Rocklin, Amit Kumar, Sergiu Ivanov, Jason K. Moore, Sartaj Singh, Thilina Rathnayake, Sean Vig, Brian E. Granger, Richard P. Muller, Francesco Bonazzi, Harsh Gupta, Shivam Vats, Fredrik Johansson, Fabian Pedregosa, Matthew J. Curry, Andy R. Terrel, Štěpán Roučka, Ashutosh Saboo, Isuru Fernando, Sumith Kulal, Robert Cimrman, and Anthony Scopatz. SymPy: Symbolic mathematics in Python. *PeerJ Computer Science*, 3:e103, 2017.
- Stanislas Polu and Ilya Sutskever. Generative language modeling for automated theorem proving. *arXiv preprint arXiv:2009.03393*, 2020.
- Qwen Team. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2023.
- Iliia Shumailov, Zakhar Shumaylov, Yiren Zhao, Nicolas Papernot, Ross Anderson, and Yarin Gal. Ai models collapse when trained on recursively generated data. *Nature*, 631:755–759, 2024.
- Avi Singh, John D Co-Reyes, Rishabh Agarwal, Ankesh Anand, Piyush Patil, Xavier Garcia, Peter J. Liu, James Harrison, Jaehoon Lee, Kelvin Xu, Aaron Parisi, Aviral Kumar, Alex Alemi, Alex Rizkowsky, Azade Nova, Ben Adlam, Bernd Bohnet, Gamaleldin Elsayed, Hanie Sedghi, Igor Codreanu, S. Keshav, et al. ReST-EM: Beyond human data: Scaling self-training for problem-solving with language models. In *Advances in Neural Information Processing Systems*, 2024.
- Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process- and outcome-based feedback. In *arXiv preprint arXiv:2211.14275*, 2022.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *International Conference on Learning Representations*, 2023.
- Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. Self-rewarding language models. In *International Conference on Machine Learning*, 2024.
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D. Goodman. STaR: Bootstrapping reasoning with reasoning. In *Advances in Neural Information Processing Systems*, volume 35, 2022.