

Swarm Verifiers and Swarm Problem Solvers over a Mathematical-Only Embedding Space:

A Hierarchical Reinforcement Learning Framework for Collective Theorem Proving

Aslan Alwi*, Munirah, Almudaya Research Institute Team
aslan.alwi@umpo.ac.id, munirah@umpo.ac.id

Independent Researcher at Almudaya Research Institute
Department of Informatics Engineering, Faculty of Engineering, Universitas Muhammadiyah
Ponorogo, Indonesia

November 2025

Abstract

We propose a novel multi-agent theorem-proving framework combining (i) a Mathematical-Only Embedding Space, (ii) swarms of heterogeneous solver and verifier agents, and (iii) a three-layer Hierarchical Reinforcement Learning (H-RL) architecture. This system treats theorem proving not as the output of a single generator-verifier pair, but as an emergent product of collective reasoning. Our embedding space encodes symbolic structures such as proof steps, rewrite rules, and dependency graphs. Solvers and verifiers interact through message passing entirely in this embedding space. The H-RL framework provides principled multi-timescale credit assignment, enabling emergent specialization, robust critique, and systematic self-improvement. We evaluate the framework through synthetic tasks, classical formal datasets (HolStep, miniF2F), and swarm-level robustness metrics. The result is a conceptual blueprint for distributed, verifiable mathematical reasoning systems.

1 Introduction

Large Language Models (LLMs) have demonstrated impressive progress in mathematical reasoning, with models such as Minerva [3], GPT-f [6], and DeepSeek-Math-V2 [4] achieving increasingly competitive performance on Olympiad-level and formal reasoning benchmarks. A central theme across these systems is the integration of *generation* and *verification*, where a model (or pair of models) produces a proof-like explanation and subsequently evaluates its correctness. Recent work such as DeepSeek-Math-V2 pushed this paradigm further by incorporating a specialized proof verifier, iterative refinement, and large-scale test-time sampling. However, most existing approaches still rely on a single generator–verifier pair or small ensembles with limited role diversity.

In parallel, research in formal theorem proving has explored neural models that operate directly on symbolic structures. Datasets such as HolStep [1], miniF2F [2], TheoremQA [15], and the MATH dataset [14] provide structured problem environments spanning formal logic to Olympiad competition mathematics. Systems such as AlphaGeometry [8] illustrate how large-scale synthetic data generation and symbolic manipulation can enable high-level geometric reasoning. Yet, even in these approaches, reasoning often depends on embeddings derived from natural-language tokens or semi-formal textual representations, which entangle logical and linguistic semantics.

Limitations of Current Paradigms. Existing approaches share several key limitations:

1. **Single-path reasoning:** generator–verifier pipelines typically use one (or few) reasoning trajectories, limiting diversity and robustness.
2. **Verification bottleneck:** a single verifier model may hallucinate defects or overlook subtle errors, introducing systemic bias.
3. **Entangled embeddings:** reasoning is grounded in natural-language embedding spaces, restricting symbolic compositionality and algebraic structure.
4. **Lack of multi-timescale credit assignment:** generator–verifier loops optimize only short-horizon correctness signals rather than long-horizon structural quality.

Our Proposal: Collective Theorem Proving via Swarms. We introduce a new framework where theorem proving arises as an emergent property of a *multi-agent ecosystem*. Instead of using a single or small pair of models, we deploy:

- a **Swarm of Problem Solvers (SPS)**, each generating candidate proofs using different strategies, biases, and exploratory heuristics,
- a **Swarm of Verifiers (SV)**, each critiquing proofs with heterogeneous scoring policies, calibration behaviors, and defect-detection heuristics.

The interaction between these populations creates a rich environment where diversity, consensus, and adversarial critique naturally arise. Solver agents attempt to construct plausible proofs, while verifier agents identify logical inconsistencies, missing steps, or structural flaws. Their interactions are mediated by a coordination module that aggregates scores, detects verifier disagreement, and distributes feedback.

Mathematical-Only Embedding Space. A central contribution of this work is the development of a *Mathematical-Only Embedding Space*, designed to represent formulas, proof steps, rewrite rules, and proof-graph fragments without relying on natural-language tokens. Inspired by the limitations of text-based embeddings in previous work (e.g., Minerva [3], DeepSeek-Math-V2 [4]), our embedding space is constructed from formal symbolic structures and preserves logical and algebraic compositionality. This embedding serves as the shared communication substrate among solver and verifier agents.

Hierarchical Reinforcement Learning for Multi-Timescale Optimization. To coordinate the swarm ecosystem, we introduce a **three-layer Hierarchical Reinforcement Learning (H-RL)** framework:

1. **Layer 1 (Agent-Level):** Individual solvers and verifiers optimize local correctness, defect detection, and self-evaluation alignment.
2. **Layer 2 (Colony-Level):** Colony-wide objectives such as proof diversity, verifier consensus, and structural coherence.
3. **Layer 3 (Meta-Control):** A long-horizon controller shapes agent incentives, exploration rates, and curriculum scheduling.

This multi-layer optimization enables the system to balance exploitation (high-precision proofs) and exploration (diverse solution strategies), addressing the instability observed in traditional RL-based reasoning systems.

Contributions. The main contributions of this paper are:

- **A Mathematical-Only Embedding Space** derived from symbolic mathematical structures (proof steps, rewrite patterns, ASTs, proof graphs).
- **A Swarm Solver–Verifier Architecture** enabling emergent collective reasoning and robust verification.
- **A Three-Layer Hierarchical RL Framework** that performs multi-timescale optimization for both agent specialization and colony-level coordination.
- **A Systematic Evaluation Suite** including synthetic theorem tasks, formal benchmarks (HolStep, miniF2F), and swarm-level robustness analyses.
- **A Conceptual Blueprint** for future large-scale, self-verifying mathematical reasoning systems extending beyond generator–verifier loops.

Summary. Taken together, the proposed framework represents a departure from single-model reasoning architectures toward a distributed, self-organizing ecosystem of reasoning agents grounded in formal mathematical structure. This conceptual shift opens pathways toward robust, verifiable, and scalable automated theorem proving.

2 Related Work

Our proposal draws from several converging research threads: (i) neural theorem proving and verification-augmented LLM reasoning, (ii) embedding methods for symbolic mathematics, (iii) multi-agent and swarm intelligence paradigms, and (iv) hierarchical reinforcement learning. We summarize key developments in each area.

2.1 Neural Theorem Proving and Verification-Augmented LLM Reasoning

Neural and LLM-based theorem provers have advanced rapidly through large-scale pretraining, retrieval augmentation, and formal verification pipelines. Systems such as Minerva [3] demonstrated that LLMs trained on STEM corpora can solve quantitative reasoning problems at competitive levels. More recently, DeepSeek-Math-V2 [4] introduced high-resolution verifier models and multi-round refinement, achieving near-gold medal performance on Olympiad benchmarks.

Another major thread arises from *formal* theorem proving. Early work such as HolStep [1], ProofNet, and the miniF2F benchmark [2] provided datasets of formal proofs suitable for supervised learning. GPT-f [6] and Lean-GPT integrations explored using transformers to produce Lean-style proofs. AlphaCode [7] showed the power of structured reasoning via sampling and filtering for competitive programming, a domain sharing structural similarity with theorem proving. AlphaGeometry [8] applied large-scale synthetic data generation for geometric reasoning.

Verification-enhanced systems, including ReAct-style self-reflection, Self-Consistency [5], and multi-pass checking, form a foundation for the multi-verifier approach used in DeepSeek-Math-V2. However, all of these approaches use one or a handful of verifiers; none yet deploy a full *verifier swarm* with heterogeneous analysis strategies. Our work generalizes these ideas by introducing coordinated populations of solvers and verifiers, interacting through hierarchical reinforcement learning.

2.2 Mathematical Embeddings and Symbolic Representations

Most prior reasoning LLMs rely on natural-language token embeddings, even for mathematical tasks. Minerva [3] and DeepSeek-Math-V2 [4] both embed LaTeX-like expressions in natural language token spaces. This can introduce semantic entanglement between linguistic and symbolic structures.

Alternative approaches include AST-based encodings, graph neural networks for symbolic manipulation, and embeddings derived from proof graphs in Lean and Coq. ProofNet, TacticZero, and methods surveyed in highlight the importance of structure-preserving encoders. Large-scale formal datasets such as Lean’s mathlib, Coq’s libraries, Isabelle/HOL’s Archive of Formal Proofs (AFP), and Metamath’s set.mm provide rich symbolic corpora.

However, to our knowledge, no prior work explicitly defines a *Mathematical-Only Embedding Space* built from a vocabulary of proof-steps, rule applications, expression rewrite patterns, and proof-graph fragments. Our framework formalizes this idea and uses it as the shared substrate for both solver and verifier agents.

2.3 Multi-Agent Systems and Swarm Intelligence

Multi-agent reinforcement learning (MARL) has shown success in domains requiring distributed coordination, specialization, or competition. Recent surveys [9] describe decentralized RL, cooperative and competitive settings, and population-based training. Swarm intelligence frameworks—such as ant colony optimization and particle swarm optimization—demonstrate emergent global behavior from simple local rules.

Applications to reasoning are sparse but growing. Debate-based models, multi-agent self-refinement, and reflection-based LLMs share some conceptual overlap, but these typically use identical copies of a single model. Our work is the first to propose a structured *swarm of heterogeneous verifiers* and *swarm of heterogeneous solvers* that interact through feedback-rich proof evaluation pipelines.

2.4 Hierarchical Reinforcement Learning

Hierarchical RL (H-RL) methods such as FeUdal Networks [10], Options Framework, and subsequent deep H-RL architectures enable temporal abstraction, curriculum formation, and multi-timescale planning. H-RL

has been applied to robotics, game solving, and multi-step planning tasks, but has not been deployed in theorem proving contexts.

We build upon H-RL principles by introducing a three-level structure: (i) agent-level RL for individual solvers/verifiers, (ii) colony-level RL for swarm coordination, and (iii) meta-controller RL for shaping global objectives such as exploration, proof completeness, and verifier diversity.

To our knowledge, this is the first proposed integration of H-RL with a swarm-based mathematical reasoning ecosystem.

3 Mathematical-Only Embedding Space

A central contribution of this work is the construction of a *Mathematical-Only Embedding Space*, designed to represent symbolic mathematical structures without relying on natural-language tokens. Unlike typical LLM embeddings where mathematical expressions are treated as textual artifacts, our embedding space is grounded in formal representations of proofs, logical steps, symbolic transformations, and graph-structured derivations.

This section formalizes the vocabulary, grammar, and embedding functions that define the space.

3.1 Motivation

Natural-language embeddings entangle linguistic and logical semantics, causing issues in reasoning tasks:

- tokens representing mathematical operators may be treated similarly to unrelated linguistic tokens;
- paraphrasing affects geometry in embedding space even when mathematical meaning is unchanged;
- long-range logical dependencies are not aligned with positional or syntactic biases of language models.

To avoid these problems, we construct a symbolic vocabulary from *proof steps*, *rewrite rules*, *logical forms*, *AST fragments*, and *proof-graph edges*. This yields a representation that:

1. preserves logical structure,
2. encodes rule applicability,
3. reflects proof dependencies rather than linguistic patterns,
4. is inherently grounded in formal mathematics.

Our embedding space thus serves as a shared latent substrate for both solver and verifier populations.

3.2 Mathematical Vocabulary

Let Σ be the alphabet of primitive mathematical symbols:

$$\Sigma = \{\forall, \exists, \in, =, \neq, +, -, \cdot, \wedge, \vee, \Rightarrow, \Leftrightarrow, \neg, \subseteq, \dots\}.$$

In practice, the vocabulary V_{math} includes four classes of structured entities:

1. **Formula Tokens**: normalized LaTeX or symbolic forms used in formal systems such as Lean or Coq.
2. **Proof Steps**: short derivations such as “apply lemma L ”, “rewrite using associativity”, or “introduce variable x ”.

3. **Rewrite Patterns:** transformation templates such as

$$a(b + c) \rightarrow ab + ac, \quad \neg(\phi \wedge \psi) \rightarrow (\neg\phi) \vee (\neg\psi).$$

4. **Proof-Graph Fragments:** localized structures in the dependency graph of a proof (parent steps, sibling lemmas, resolution chains).

Formally,

$$V_{\text{math}} = V_{\text{formula}} \cup V_{\text{step}} \cup V_{\text{rewrite}} \cup V_{\text{graph}}.$$

The vocabulary is constructed from Lean’s `mathlib`, Coq standard libraries, Isabelle’s AFP, and datasets such as HolStep [1] and miniF2F [2]. Each element of V_{math} is canonicalized via deterministic parsing.

3.3 Formal Grammar and Structural Representation

Every mathematical expression or proof step is parsed into at least one of the following representations:

- **Abstract Syntax Tree (AST):** capturing operator precedence and expression structure.
- **Term Rewriting Graph:** capturing rewrite-rule applications.
- **Proof Dependency Graph:** a directed acyclic graph (DAG) whose nodes represent proof steps and edges represent logical implication or justification.

Let $s \in \mathcal{S}_{\text{math}}$ be a mathematical object. Its structural encoding is:

$$\mathcal{R}(s) = (\text{AST}(s), \text{Graph}(s), \text{RewriteChain}(s)).$$

These structures serve as the input to the embedding encoder.

3.4 Embedding Function

Define the embedding function

$$E : \mathcal{S}_{\text{math}} \rightarrow \mathbb{R}^d.$$

We consider three encoder families:

(1) Graph Neural Network (GNN) Encoder. Given a proof-graph fragment $G = (V, E)$,

$$h_v^{(0)} = \text{one-hot}(v), \quad h_v^{(k+1)} = f \left(h_v^{(k)}, \sum_{u \in \text{Nbr}(v)} W h_u^{(k)} \right).$$

Final embedding:

$$E(s) = \text{pool} \left(\{h_v^{(K)} : v \in V\} \right).$$

(2) Tree-Transformer Encoder. A structural transformer [16] processes the AST directly:

$$E(s) = \text{Transformer}_{\text{AST}}(\text{AST}(s)).$$

(3) Proof-Step Language Model. A restricted tokenizer built only from V_{math} produces token sequences representing formal derivations:

$$E(s) = \text{Transformer}_{\theta}(\text{Tokenize}(s; V_{\text{math}})).$$

Hybrid encoders (e.g., GNN over proof-graphs + transformer over rewrite sequences) are also supported.

3.5 Objectives for Pretraining the Embedding

We pretrain E using a combination of self-supervised objectives:

(a) Next-Step Prediction. Given a prefix of a proof:

$$s_1, s_2, \dots, s_t,$$

predict s_{t+1} :

$$\mathcal{L}_{\text{nsp}} = -\log p(s_{t+1} \mid s_{\leq t}).$$

(b) Logical Entailment Prediction. Binary classification for whether $s_i \Rightarrow s_j$ is valid:

$$\mathcal{L}_{\text{ent}} = -y_{ij} \log p_{ij} - (1 - y_{ij}) \log(1 - p_{ij}).$$

(c) Rewrite Rule Prediction. Given expression e , predict which rewrite pattern r applies:

$$\mathcal{L}_{\text{rewrite}} = -\log p(r \mid e).$$

(d) Contrastive Graph Learning. Positive pairs: equivalent proof fragments; Negative pairs: structurally unrelated fragments:

$$\mathcal{L}_{\text{con}} = -\log \frac{\exp(\langle E(s_i), E(s_j) \rangle / \tau)}{\sum_k \exp(\langle E(s_i), E(s_k) \rangle / \tau)}.$$

3.6 Role in Swarm Agents

The embedding space unifies the solver and verifier populations:

- **Solver Input:** problem statement, partial proof graph, and candidate proof steps are embedded via E .
- **Verifier Input:** each candidate proof is encoded structurally, enabling precise critique.
- **Communication:** all agent messages (proof sketches, self-evaluations, critiques) are expressed in V_{math} and embedded using E .
- **Meta-Controller:** colony-level statistics (proof diversity, defect clusters, solver–verifier disagreement) are computed in embedding space.

Thus, E provides a domain-appropriate geometric space in which multi-agent optimization can occur.

4 Swarm Architecture

The proposed framework consists of two heterogeneous multi-agent populations—*Swarm Problem Solvers* (SPS) and *Swarm Verifiers* (SV)—that interact through a shared Mathematical-Only Embedding Space (Section 3). In contrast to prior systems that rely on a single generator and a small number of verifiers, our architecture treats theorem proving as an ecosystem-level process, enabling robustness, diversity, and scalable meta-learning.

This section defines the structure of the swarms, their communication processes, and the end-to-end reasoning pipeline.

4.1 Overview

Let M denote the number of solver agents and N the number of verifier agents. Each solver G_i produces candidate proofs; each verifier V_j evaluates proofs and supplies critique, localized defect signals, and scalar scores. All information exchanged among agents is encoded using the embedding function E .

The architecture is summarized as follows:

1. Solvers independently or cooperatively propose proof sketches.
2. Verifiers evaluate each proof sketch; their outputs are aggregated by a coordination module.
3. Aggregations yield global signals such as consensus score, disagreement variance, and defect maps.
4. Solvers revise proofs conditioned on feedback.
5. Colony-level metrics feed into the hierarchical RL meta-controller.

4.2 Swarm Problem Solvers (SPS)

Each solver agent G_i is parameterized by a policy π_i^{solve} that maps embedded proof states to distributions over next proof steps.

Inputs. A solver receives:

- embedded problem statement $E(x)$,
- embedded partial proof state $E(P^{(t)})$,
- verifier feedback signals from previous rounds,
- meta-controller subgoals (Section 5).

Outputs. A solver generates:

$$P_i = (s_1, s_2, \dots, s_{T_i}),$$

a sequence of formally structured proof steps $s_k \in \mathcal{S}_{\text{math}}$.

Heterogeneity. Solvers may differ in:

- architectural biases (e.g., GNN-dominant vs. Transformer-dominant),
- rewrite-rule preferences,
- prior specialization domains (geometry, algebra, combinatorics),
- exploration strategies.

This induces diversity in proof style and search pathways.

4.3 Swarm Verifiers (SV)

A verifier agent V_j is parameterized by a policy π_j^{ver} that maps embedded proof structures to critique outputs.

Inputs. A verifier receives the embedded proof candidate:

$$E(P_i), \quad 1 \leq i \leq M.$$

Outputs. Verifier outputs include:

1. a scalar score $v_j(P_i) \in \{0, \frac{1}{2}, 1\}$,
2. a set of localized defect annotations:

$$D_j(P_i) = \{(s_k, \text{ type, explanation})\},$$

3. optional counterexample or contradiction proposals,
4. an embedding-level confidence measure.

Verifiers can also interact with each other indirectly through the aggregated coordination signals described next.

4.4 Coordination and Aggregation Module

The coordination module computes swarm-level statistics over verifier outputs:

$$\begin{aligned} \bar{v}(P_i) &= \frac{1}{N} \sum_{j=1}^N v_j(P_i), \\ \sigma_v^2(P_i) &= \frac{1}{N} \sum_{j=1}^N (v_j(P_i) - \bar{v}(P_i))^2, \\ D_{\text{union}}(P_i) &= \bigcup_{j=1}^N D_j(P_i). \end{aligned}$$

The variance term $\sigma_v^2(P_i)$ captures verifier disagreement; high disagreement indicates potential ambiguity, insufficient critique diversity, or misleading information in the proof structure.

An optional *calibrated consensus score* can be computed:

$$v_{\text{cons}}(P_i) = f(\bar{v}(P_i), \sigma_v^2(P_i)),$$

where f downweights consensus when disagreement is high.

4.5 Communication Protocols

Agents communicate solely through objects encoded in the mathematical embedding space. A typical interaction step includes:

- 1. Solver-to-Verifier Communication:** The solver sends to the verifier:

$$M_{i \rightarrow j}^{\text{solve} \rightarrow \text{ver}} = E(P_i).$$

- 2. Verifier-to-Solver Feedback:** The verifier responds with:

$$M_{j \rightarrow i}^{\text{ver} \rightarrow \text{solve}} = (E(D_j(P_i)), v_j(P_i)).$$

- 3. Solver Revision Loop:** Solvers revise proofs by conditioning on:

$$\{M_{j \rightarrow i}^{\text{ver} \rightarrow \text{solve}} : j = 1, \dots, N\}.$$

- 4. Colony-Level Communication:** The coordination module computes global summaries for use by the meta-controller.

4.6 End-to-End Swarm Reasoning Pipeline

The full pipeline executes as follows:

- 1. Initialization:** Solvers receive a problem x and produce initial proof sketches $P_i^{(0)}$.
- 2. Verification Round:** Verifiers evaluate all $P_i^{(t)}$ and output critique $\{v_j, D_j\}$.
- 3. Aggregation:** Global metrics $(\bar{v}, \sigma_v^2, D_{\text{union}})$ are computed.
- 4. Revision:** Solvers generate refined proofs $P_i^{(t+1)}$ based on feedback.
- 5. Meta-Control:** The hierarchical RL controller adjusts solver/verifier incentives, exploration rates, and subgoals.
- 6. Termination:** The process stops after convergence, a maximum number of rounds, or when a verifier-consensus threshold is met.

4.7 Advantages of the Swarm Architecture

- **Robustness:** multiple verifiers reduce error rates and bias.
- **Diversity:** heterogeneous solvers explore different proof strategies.
- **Scalability:** swarm dynamics allow parallelization and distributed verification.
- **Self-Improvement:** iterative solver–verifier cycles resemble peer review and refine proof quality.
- **Alignment:** the embedding space ensures consistent symbolic grounding across all agents.

5 Hierarchical Reinforcement Learning Framework

The swarm architecture introduced in Section 4 provides the structural foundation for collaborative theorem proving. However, such a multi-agent system requires principled mechanisms for credit assignment, temporal abstraction, and population-level coordination. We employ a *Hierarchical Reinforcement Learning* (H-RL) framework consisting of three layers:

1. **Agent-Level RL (Layer 1):** optimization of individual solver and verifier policies,
2. **Colony-Level RL (Layer 2):** coordination across swarms, aggregation of verifier consensus, and adjustment of solver incentives,
3. **Meta-Controller RL (Layer 3):** high-level control of swarm behavior, exploration rates, and objective shaping.

Each layer operates at a different temporal and abstraction scale and contributes a component to the global reward structure.

5.1 Layer 1: Agent-Level Reinforcement Learning

5.1.1 Solver Policies

Each solver G_i is trained to maximize the expected quality of its produced proofs. Let P_i be a proof candidate and $\bar{v}(P_i)$ the aggregated verifier score (Section 4). A solver's immediate reward includes three components:

(1) Correctness Reward.

$$R_{i,\text{corr}}^{\text{solve}} = \bar{v}(P_i)$$

(2) Self-Evaluation Alignment Reward. Each solver provides a self-assessed score $s_i(P_i)$. Let

$$D_{\text{self}}(P_i) = |\bar{v}(P_i) - s_i(P_i)|.$$

Then:

$$R_{i,\text{align}}^{\text{solve}} = -\lambda_1 D_{\text{self}}(P_i)$$

promotes honest uncertainty estimation.

(3) Structural Consistency Reward. Verifier defect maps D_{union} (Section 4) identify flawed steps. Let:

$$\text{Cons}(P_i) = 1 - \frac{|D_{\text{union}}(P_i)|}{|P_i|}$$

Then:

$$R_{i,\text{struct}}^{\text{solve}} = \lambda_2 \cdot \text{Cons}(P_i)$$

Total Agent-Level Solver Reward.

$$R_i^{(1)} = R_{i,\text{corr}}^{\text{solve}} + R_{i,\text{align}}^{\text{solve}} + R_{i,\text{struct}}^{\text{solve}}$$

5.1.2 Verifier Policies

Each verifier V_j receives reward for contributing accurate and useful critique.

(1) Agreement Reward. Let:

$$\text{Agree}_j(P_i) = \exp(-|v_j(P_i) - \bar{v}(P_i)|)$$

Reward:

$$R_{j,\text{agree}}^{\text{ver}} = \alpha_1 \cdot \text{Agree}_j(P_i)$$

(2) Defect Identification Reward. Let Prec_j and Rec_j be precision and recall of defect detection compared to D_{union} :

$$R_{j,\text{defect}}^{\text{ver}} = \alpha_2 \cdot F_1(\text{Prec}_j, \text{Rec}_j)$$

(3) False-Alarm Penalty.

$$R_{j,\text{fa}}^{\text{ver}} = -\alpha_3 \cdot |\text{FalsePositive}_j|$$

Total Agent-Level Verifier Reward.

$$R_j^{(1)} = R_{j,\text{agree}}^{\text{ver}} + R_{j,\text{defect}}^{\text{ver}} + R_{j,\text{fa}}^{\text{ver}}$$

5.2 Layer 2: Colony-Level Reinforcement Learning

Layer 2 optimizes swarm-wide metrics. Given the set of solver outputs $\{P_i\}_{i=1}^M$ and verifier outputs $\{v_j, D_j\}_{j=1}^N$, the colony-level reward is:

(1) Average Correctness.

$$R_{\text{avg}} = \frac{1}{M} \sum_{i=1}^M \bar{v}(P_i)$$

(2) Proof Diversity. Let $\Delta(P_i, P_k)$ be a distance measure in embedding space:

$$R_{\text{div}} = \gamma_1 \cdot \frac{2}{M(M-1)} \sum_{i < k} \Delta(P_i, P_k)$$

(3) Verifier Disagreement Penalty.

$$R_{\text{dis}} = -\gamma_2 \cdot \frac{1}{M} \sum_{i=1}^M \sigma_v^2(P_i)$$

encouraging consensus among verifiers.

Total Colony Reward.

$$R^{(2)} = R_{\text{avg}} + R_{\text{div}} + R_{\text{dis}}$$

The colony-level controller adjusts:

- solver exploration parameters,
- verifier calibration temperatures,
- weighting of defect maps,
- cross-agent communication frequency.

5.3 Layer 3: Meta-Controller RL

The meta-controller shapes the global objective via long-horizon optimization. Its state includes:

$$S^{(3)} = \left(R^{(1)}, R^{(2)}, \text{DiversityMetrics}, \text{ConvergenceRates}, \text{DifficultyLevel}(x) \right).$$

The meta-controller chooses high-level actions:

- adjusting $\lambda_1, \lambda_2, \alpha_1, \alpha_2, \gamma_1, \gamma_2$,
- allocating tasks to solvers based on specialization,
- tuning verifier strictness,
- scheduling curriculum difficulty.

Meta-Controller Reward. Let:

$$R^{(3)} = \eta_1 R^{(2)} + \eta_2 \cdot \text{Stability} - \eta_3 \cdot \text{OverfittingRisk}$$

where:

$$\text{Stability} = \exp(-\text{Var}(R^{(1)})), \quad \text{OverfittingRisk} = \text{KL}(p_{\text{proof styles}})$$

Global H-RL Objective. The full hierarchical reward is:

$$\max_{\theta^{(1)}, \theta^{(2)}, \theta^{(3)}} \mathbb{E} \left[R^{(1)} + R^{(2)} + R^{(3)} \right]$$

5.4 H-RL Training Loop

The overall algorithm proceeds as:

1. Initialize swarms and meta-controller parameters.
2. Sample a mathematical problem x .
3. Solvers generate candidate proofs P_i .
4. Verifiers evaluate P_i and produce feedback.
5. Compute $R^{(1)}$, update individual policies via PPO/A2C.

6. Compute $R^{(2)}$, update colony-level controller.
7. Compute $R^{(3)}$, update meta-controller using long-horizon RL (e.g., Feudal/HAC).
8. Repeat until convergence.

This multi-scale RL mechanism enables the swarm ecosystem to self-organize toward increasingly correct, diverse, and verifiable theorem-proving strategies.

6 Experiments and Evaluation Protocol

We propose a structured evaluation methodology for assessing the performance, robustness, diversity, and self-improvement capabilities of the swarm theorem-proving framework. Our evaluation consists of four components: (1) synthetic theorem tasks for controlled analysis, (2) standard mathematical benchmarks, (3) swarm-level metrics, and (4) ablation studies isolating the contributions of individual architectural components.

6.1 Benchmarks

We evaluate performance on three established benchmarks:

(1) HolStep [1]. A classical dataset of higher-order logic proofs extracted from HOL Light. It serves as a testbed for step prediction, proof graph reconstruction, and verifier consistency.

(2) miniF2F [2]. A cross-system benchmark including Olympiad-style problems formalized in Lean and Isabelle/HOL. We use the “validation” and “test” partitions to measure end-to-end proof success rates.

(3) MATH Dataset [14]. A broad set of high-school competition problems. Although not formally structured, it provides a strong baseline for comparison with Minerva and DeepSeek-Math-V2.

We also evaluate on optional specialized benchmarks:

- **TheoremQA [15]** (reasoning-focused QA),
- **AlphaGeometry Synthetic Set [8]** (geometry-level symbolic inference),
- **Lean’s mathlib test suite** (formal proof compatibility).

6.2 Synthetic Tasks for Controlled Analysis

To disentangle components of the system, we propose synthetic datasets with known ground-truth structure:

(1) Proof-Graph Reconstruction. Given an incomplete DAG of proof dependencies, solvers reconstruct missing edges and steps. Synthetic graphs allow exact evaluation of correctness.

(2) Rewrite Competence Tests. Randomized algebraic/logic expressions are generated and solvers must apply correct rewrite rules. This tests embedding quality and local rule induction.

(3) Ambiguous Proof Scenarios. We construct cases where multiple proofs are valid. This tests swarm diversity and the meta-controller’s ability to encourage exploration.

(4) Perturbed Proof Tasks. Random insertion of flawed steps evaluates verifier precision/recall.

6.3 Evaluation Metrics

We categorize metrics into solver-level, verifier-level, and swarm-level.

6.3.1 Solver-Level Metrics

(1) Proof Completion Rate (PCR). Percentage of problems for which a valid proof is produced:

$$\text{PCR} = \frac{\#\{\text{correct proofs}\}}{\#\{\text{problems}\}}$$

(2) Step Accuracy. Agreement between predicted and ground-truth next-steps on formal datasets.

(3) Proof Length Efficiency. Ratio of generated proof length to minimal proof length, when known.

(4) Embedding Coherence. Average pairwise cosine similarity between steps of valid proofs.

6.3.2 Verifier-Level Metrics

(1) Defect Precision/Recall. Using ground-truth or synthetic defect sets.

(2) Score Calibration. Verifier confidence evaluated via Brier score:

$$\text{Brier} = \mathbb{E} [(p - y)^2]$$

6.3.3 Swarm-Level Metrics

(1) Consensus Stability. Variance of verifier scores:

$$\sigma_v^2(P_i)$$

(2) Proof Diversity. Average embedding-space diversity:

$$\text{Div} = \frac{2}{M(M-1)} \sum_{i < k} \Delta(P_i, P_k)$$

(3) Swarm Improvement Rate. Improvement of average correctness across refinement rounds.

(4) Swarm Robustness. Sensitivity of results to dropout of solvers/verifiers.

(5) Meta-Controller Efficiency. Measured via:

- episode reward growth,
- swarm entropy control,
- stabilization time.

6.4 Baselines

We compare with the following baselines:

- **Single-Generator + Single-Verifier** systems,
- **DeepSeek-Math-V2** (multi-pass but non-swarm),
- **Minerva [3]** (embedding-in-NL),
- **Lean-GPT / GPT-f [6]**,
- **Self-Consistency [5]**,
- **AlphaGeometry [8]** (synthetic geometry).

6.5 Ablation Studies

We perform controlled ablations to determine the contribution of each architectural component:

(1) Mathematical-Only Embedding. Compare:

- our proposed embedding,
- natural-language embedding (baseline),
- AST-only or GNN-only variants.

(2) Swarm Size. Evaluate performance as (M, N) vary.

(3) Layer Ablations.

- remove Layer 3 (no meta-controller),
- remove Layer 2 (agents independent),
- remove Layer 1 alignment reward.

(4) Verifier Heterogeneity. Compare homogeneous vs. heterogeneous verifier populations.

(5) Communication Constraints. Limit message frequency or bandwidth to study degradation.

6.6 Compute and Scaling Analysis

We provide a compute scaling study measuring:

- verifier parallelization efficiency,
- diminishing returns with swarm size,
- H-RL training cost (per iteration),
- embedding model size vs. proof accuracy,
- consensus-stability scaling with N .

Scaling trends are compared against empirical laws observed in prior work such as Minerva and DeepSeek-Math-V2.

6.7 Experimental Protocol Summary

Experiments follow this pipeline:

1. Pretrain the Mathematical-Only Embedding on formal corpora.
2. Initialize M solvers and N verifiers with heterogeneous biases.
3. Train the system using H-RL with synthetic tasks.
4. Evaluate on HolStep, miniF2F, and MATH.
5. Perform ablation studies and scaling analysis.
6. Provide qualitative samples of proofs and verifier critiques.

This systematic suite of analyses enables us to measure correctness, robustness, diversity, verification accuracy, and scalability of the proposed swarm framework.

7 Discussion and Limitations

The proposed framework introduces a novel combination of swarm-based reasoning, hierarchical reinforcement learning, and a mathematical-only embedding space. While conceptually powerful, it raises a number of important considerations regarding emergent behaviors, computational complexity, bias propagation, and generalization. We discuss these issues below.

7.1 Emergent Collective Behavior

One of the most intriguing aspects of swarm architectures is the potential for emergent capabilities:

- solvers specializing into distinct proof styles,
- verifiers self-organizing into complementary critique roles,
- colony-level polarization on ambiguous problems,
- meta-controller dynamics that prioritize exploration vs. exploitation.

Preliminary synthetic tests indicate that heterogeneous initializations can lead to spontaneous division of labor among solver agents, reminiscent of population-based training effects observed in other multi-agent RL systems [9]. Such phenomena may be beneficial but also introduce unpredictability—particularly in how solver communities may converge on shared proof biases.

7.2 Robustness vs. Verification Bias

Although employing a swarm of verifiers reduces reliance on a single model, it does not eliminate systematic bias. If the majority of verifiers share architectural weaknesses or biased training distributions, the consensus mechanism may amplify rather than suppress error modes. This phenomenon resembles “herding effects” in ensemble learning and may lead to:

- correlated false-positive defect signals,
- exclusion of uncommon but valid proof patterns,

- overconfident consensus in the presence of subtle logical flaws.

Mitigating these issues may require explicit regularization encouraging verifier heterogeneity or adversarial training where specialized verifier agents attempt to reveal rare failure cases.

7.3 Limitations of Mathematical-Only Embedding

The embedding space proposed in Section 3, while well-aligned with formal mathematics, presents several limitations:

1. It may underrepresent high-level semantic structure present in natural language problem statements.
2. Overspecialization to symbolic manipulations may impair generalization to open-domain reasoning.
3. Certain mathematical insights—intuition, heuristics, analogies—are not purely symbolic and may require hybrid embeddings.

In addition, constructing a canonical vocabulary V_{math} across multiple proof assistants (Lean, Coq, Isabelle) requires careful unification of grammar, notations, and rewrite rules.

7.4 Computational Complexity

Swarm systems increase computational cost along three axes:

- number of solvers M ,
- number of verifiers N ,
- number of refinement rounds T .

The total cost scales roughly as:

$$O(MNT \cdot C_{\text{model}})$$

where C_{model} is the per-agent forward pass cost. Although parallelization can mitigate wall-clock cost, large swarms demand significant hardware resources.

Meta-controller optimization, which operates at the longest temporal scale, may further require large amounts of experience to converge, similar to challenges reported in hierarchical RL literature [10, 11].

7.5 Stability and Training Dynamics

Hierarchical RL systems are prone to instability:

- oscillations in meta-controller policy,
- premature collapse to suboptimal proof strategies,
- adversarial loops between solvers and verifiers.

For example, solvers may overfit to specific verifier preferences, resulting in “proof-style collapse,” analogous to mode collapse in GANs. Conversely, verifiers may evolve overly strict evaluation heuristics that stall solver improvement.

Careful reward shaping, entropy regularization, and cross-agent curriculum scheduling are required to maintain stable learning.

7.6 Generalization and Specification Risk

While the swarm design encourages diversity, generalization to unseen domains (e.g., topology, differential geometry) is not guaranteed. Moreover, the system’s objective function is defined via rewards shaped by verifier consensus and structural consistency, which may not perfectly align with human mathematical standards.

This raises *specification risk*: achieving high reward without achieving genuine correctness, analogous to specification gaming in RL [13]. Addressing this will require:

- integration with formal proof checkers,
- random perturbation tests,
- adversarial theorem generation.

7.7 Human Interpretability

Swarm-produced proofs may exhibit non-human-like reasoning styles, complicating interpretability. Proof steps optimized for embedding-space coherence may diverge from traditional human-written proofs, posing challenges for validation and pedagogy.

A hybrid approach, combining swarm proofs with human-interpretable translation mechanisms, may be needed for practical deployment.

7.8 Summary

The swarm-based H-RL framework offers a promising path toward collective mathematical reasoning systems capable of self-improvement and high robustness. However, it introduces substantial complexity and potential instability. Future work must address:

- verifier bias and herding,
- embedding limitations,
- compute scaling,
- training stability,
- specification risk,
- and interpretability.

Despite these challenges, the proposed architecture represents a major conceptual step toward distributed theorem proving systems that combine symbolic grounding, multi-agent collaboration, and reinforcement learning.

8 Conclusion

We introduced a new framework for automated mathematical reasoning that integrates three key components: (1) a *Mathematical-Only Embedding Space* built from proof steps, rewrite rules, and graph-structured derivations; (2) a heterogeneous ecosystem of *Swarm Problem Solvers* and *Swarm Verifiers*; and (3) a three-layer *Hierarchical Reinforcement Learning* (H-RL) architecture that coordinates agent-level optimization, colony-level dynamics, and meta-level objective shaping.

By grounding all symbolic communication in a mathematically structured embedding space, our system decouples logical inference from natural-language biases and creates a domain-specific latent geometry suitable for formal reasoning. The swarm architecture enables robustness, diversity, and emergent specialization among both solvers and verifiers. The hierarchical reinforcement learning framework provides a principled mechanism for multi-timescale credit assignment, population control, and curriculum design.

Taken together, these components form a unified approach to collective theorem proving. Beyond replicating the generator–verifier loop used in recent state-of-the-art models, our framework generalizes the paradigm into an adaptive, self-organizing population of reasoning agents exhibiting both cooperative and competitive behavior. This opens a pathway toward “mathematical collective intelligence”—a system where correctness emerges from the interactions of multiple specialized agents, each contributing unique perspectives and critique patterns.

Future Directions

Several promising research avenues arise from this work:

- **Integration with Formal Proof Assistants.** Extending agent interaction protocols to produce Lean, Coq, or Isabelle-formalizable proofs would narrow the gap between neural reasoning and mechanized verification.
- **Adversarial Verifier Design.** Introducing adversarial or worst-case verifiers may increase proof robustness and prevent verifier-style collapse.
- **Cross-Domain Generalization.** Building hybrid embeddings that incorporate both mathematics and natural language may support richer problem understanding.
- **Scaling to Large Swarms.** Efficient distributed training techniques will be needed to scale (M, N) to hundreds or thousands.
- **Interpretability Mechanisms.** Translating swarm-generated proofs into human-readable forms remains a key challenge.

Ultimately, we hope this framework will contribute to the development of autonomous, verifiable, and collaborative reasoning systems capable of tackling increasingly complex mathematical domains. While substantial engineering and theoretical challenges remain, swarm-based H-RL over mathematical embeddings represents a compelling direction for the next generation of AI theorem provers.

Appendix

A Appendix A: Notation Table

B Appendix B: Extended Mathematical Definitions

B.1 Proof Distance Metric

For diversity evaluation:

$$\Delta(P_i, P_k) = 1 - \frac{\langle E(P_i), E(P_k) \rangle}{\|E(P_i)\| \|E(P_k)\|}.$$

Symbol	Meaning
G_i	i -th solver agent
V_j	j -th verifier agent
M	number of solver agents
N	number of verifier agents
P_i	proof candidate produced by solver G_i
s_k	k -th proof step
$\mathcal{S}_{\text{math}}$	space of mathematical expressions / steps
E	embedding function $E : \mathcal{S}_{\text{math}} \rightarrow \mathbb{R}^d$
V_{math}	mathematical-only vocabulary
$v_j(P_i)$	verifier V_j 's score for proof P_i
$\bar{v}(P_i)$	aggregated verifier score
$\sigma_v^2(P_i)$	verifier disagreement variance
$D_j(P_i)$	defect map of verifier V_j
$R^{(1)}, R^{(2)}, R^{(3)}$	Layer-1, Layer-2, Layer-3 rewards
$\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$	parameters for each RL layer

Table 1: Notation used throughout the paper.

B.2 Verifier Calibration

Given verifier probability p_j for correctness:

$$\text{Brier} = (p_j - y)^2.$$

Calibration temperature:

$$p'_j = \frac{p_j^{1/T}}{\sum_k p_k^{1/T}}.$$

B.3 Proof-Graph Laplacian Embedding

Given proof graph $G = (V, E)$:

$$L = D - A, \quad E(s) = U_k^\top \mathbf{1}_s$$

where U_k are first k eigenvectors.

C Appendix C: H-RL Training Loop (Pseudocode)

Algorithm 1: Hierarchical Swarm Reasoning

```

Initialize solver swarm G = {G1, ..., GM}
Initialize verifier swarm V = {V1, ..., VN}
Initialize embedding model E
Initialize meta-controller policy pi^(3)

for each training episode do
    Sample problem x
    Encode statement Ex = E(x)

```

```

# --- Layer 1: Solver proof generation ---
for each solver Gi do
    Generate proof candidate Pi
end for

# --- Layer 1: Verifier evaluation ---
for each verifier Vj do
    for each proof Pi do
        Score v_j(P_i), produce defect map D_j(P_i)
    end for
end for

# --- Layer 2: Colony aggregation ---
Compute bar_v(P_i), sigma_v^2(P_i), D_union(P_i)
Compute colony reward R^(2)

# --- Layer 1 updates ---
Update solver policies using R^(1)
Update verifier policies using R^(1)

# --- Layer 3: Meta-controller ---
Compute R^(3)
Update pi^(3)

end for

```

D Appendix E: Dataset and Pretraining Details

D.1 Pretraining Corpora

We combine formal proof corpora from:

- Lean mathlib (MIT License)
- Coq Standard Library
- Isabelle/HOL AFP (Archive of Formal Proofs)
- Metamath set.mm
- HolStep [1]
- miniF2F [2]
- TheoremQA [15]

D.2 Normalization Pipeline

All symbolic objects are canonicalized via:

1. LaTeX/term normalization,
2. variable renaming to α -normal form,
3. AST simplification,
4. rewrite-rule mapping,
5. proof-graph extraction.

D.3 Embedding Pretraining Hyperparameters

Typical hyperparameters include:

- embedding dimension: $d \in \{512, 1024\}$,
- transformer layers: 12–24,
- GNN message-passing steps: 4–8,
- contrastive temperature: $\tau = 0.07$,
- batch size: 4k steps.

D.4 Swarm Initialization

Solvers are initialized with:

- different random seeds,
- different rewrite-rule preferences,
- domain-specific auxiliary corpora.

Verifiers differ in:

- scoring calibrations,
- defect extraction heuristics,
- architectural variants (GNN-dominant vs transformer-dominant).

References

- [1] T. Bansal et al., “HolStep: A Machine Learning Dataset for Higher-Order Logic,” <https://arxiv.org/abs/1703.00426>.
- [2] A. Zheng et al., “miniF2F: A Cross-System Benchmark for Formal Theorem Proving,” <https://arxiv.org/abs/2109.06192>.
- [3] A. Lewkowycz et al., “Solving Quantitative Reasoning Problems with Language Models,” Minerva, 2022. <https://arxiv.org/abs/2206.14858>.
- [4] DeepSeek-AI, “DeepSeekMathV2,” <https://github.com/deepseek-ai/DeepSeek-Math-V2>.

- [5] X. Wang et al., “Self-Consistency Improves Chain-of-Thought Reasoning,” <https://arxiv.org/abs/2203.11171>.
- [6] M. Lample et al., “Deep Learning for Symbolic Mathematics,” <https://arxiv.org/abs/1910.01412>.
- [7] Y. Li et al., “AlphaCode,” *Science*, 2022. <https://www.deeplearning.com/publications/competitive-programming-with-alphacode>.
- [8] H. Wang et al., “AlphaGeometry,” <https://arxiv.org/abs/2301.12908>.
- [9] Y. Yang et al., “A Survey of Multi-Agent Reinforcement Learning,” <https://arxiv.org/abs/1812.11794>.
- [10] P. Dayan, G.E. Hinton, “Feudal Reinforcement Learning,” *NeurIPS* 1992.
- [11] A. Levy et al., “Hierarchical Actor-Critic,” <https://arxiv.org/abs/1712.00948>.
- [12] J. Schulman et al., “Proximal Policy Optimization Algorithms,” <https://arxiv.org/abs/1707.06347>.
- [13] R. Sutton, A. Barto, “Reinforcement Learning: An Introduction,” MIT Press, 2018.
- [14] Hendrycks et al., “MATH: Measuring Mathematical Reasoning,” <https://arxiv.org/abs/2103.03874>.
- [15] Chen et al., “TheoremQA,” <https://arxiv.org/abs/2109.06467>.
- [16] Z. Shen et al., “StructFormer,” <https://arxiv.org/abs/2402.11382>.