# Agentic Lean Auformalization (ALA): An LLM collaborative approach to autoformalization in LEAN

**Anonymous Author(s)**
Affiliation
Address
`email`

## Abstract

The arrival of AI systems that can achieve a gold medal at the International Mathematical Olympiad (IMO) and the development of proof assistants such as Lean seem to foretell a transformative revolution in mathematical research. However, a bottleneck is that most undergraduate- and graduate-level theorems are not translated into code for proof assistants, a process known as *autoformalization.* State-of-the-art fine-tuned LLMs in Lean 4 report at most 22.5% accuracy (Pass@128) on graduate-level theorems. To address this gap, we propose and evaluate ALA, an agentic framework where a generalist LLM orchestrating tools works together with another LLM fine-tuned in Lean 4. ALA achieves a 52 % accuracy with less than 13 tool-calls on theorems from areas such as complex and real analysis, topology, and algebra. Our code and the related dataset are published on GitHub. [1]
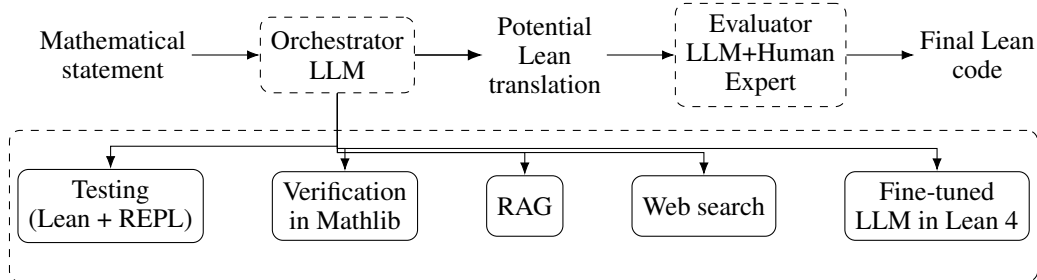
## 1 Introduction

Although large language models (LLMs) are increasingly capable of producing complex mathematical arguments [Cas25], their probabilistic outputs conflict with the certainty required by the mathematical community. Proof assistants such as Lean [dMU21] address this conflict by formally certifying the logical correctness of a proposed proof. The transformative nature of combining generative AI with formal verification has recently attracted much attention within mathematical research [BAMa, BAMb]. In particular, there is an increasing number of fine-tuned LLMs trained on Lean 4 data and autoformalization [GWJ+25] [WUL+25] [WZJ+24]. However, the use of such tools for the working mathematician is currently limited because many important undergraduate- and graduate-level topics, such as the special linear group, algebras over commutative rings, are not yet available in Lean code [Lea]. We discuss some of the challenges of autoformalization in Section 2 and Appendix A.1.

**Contributions:** Our contributions to address this challenge are threefold. (i) We present ALA, **A**gentic **L**ean **A**utoformalization, an agentic framework that combines the abilities of a fine-tuned LLM in Lean 4, the tool capabilities of a generalist LLM, and a combination of human expert and LLM judgment for translating mathematical statements to Lean 4, see Figure 1, (ii) We present a database of 200 graduate and 200 upper undergraduate level theorems covering topology, analysis, algebra, real and complex analysis. (iii) We evaluate ALA and identify strengths, weaknesses, and future areas of work on our agentic approach. ALA translates 64% of the 400 problems in our database with less than 25 tool calls. Results are discussed in Section 6.

---

[1]https://anonymous.4open.science/r/$Lean_{Translation_Agent}-CC41/README.md$

Figure 1: ALA framework: A generalist LLM with access to four tools, a LLM fine-tuned on Lean 4, and a reasoning LLM model that, together with a human expert, evaluates the final accuracy of the translation, see Section 3.



## 2 Preliminaries

### 2.1 Autoformalization

The goal of *autoformalization* is to automatically translate mathematics from natural language into machine-checked formal code. This vision dates back at least to de Bruijn's *AUTOMATH* [dB70] and has seen a modern resurgence with LLMs and interactive theorem provers. AI is the tool for automation, whereas proof assistants—here, Lean 4—are the setting for formalization. Currently, there is active research on improving LLMs to generate proofs in Lean 4 and on constructing databases for future AI training; see [WDL+25]. We highlight three key challenges.

(1) Translating a theorem statement into compiling Lean 4 code—even without a proof—depends on prior notations, typeclass instances, definitions, and lemmas. For example, we cannot formalize vector spaces without a previous formalization of a field such as the real numbers.

(2) Generating Lean 4 code that compiles does not guarantee, without human evaluation, that the translation is faithful. Sometimes, the errors are obvious and in other cases they are much subtler, see the Appendix A.1.

(3) Lean 4 is based on an extension of Martin–Löf's dependent type theory [dMU21], whereas traditional mathematics is based on an extension of set theory. These foundations are radically different; for example, in type theory, even proofs are first-class citizens part of the object-language, but in set theory, a proof is part of the meta-language and not naively the object-language.

### 2.2 Agents

In classical AI, an *agent* perceives and acts to achieve goals; modern LLM-based agents extend this loop by interleaving planning, tool use, observation, and revision [RN95, GCW+24]. In software engineering, multi-agent systems coordinate specialized roles for retrieval, coding, execution, and debugging [HTL25]. Such agents can also *self-reflect*, storing intermediate attempts and feedback to guide subsequent decisions [SCB+24]. Given these advantages, in formalization an agentic approach that couples a generalist planner with Lean-specialized models and treats the proof assistant as a verifier is natural: it can decompose natural-language statements, retrieve examples, autoformalize necessary lemmas, and iterate. Recent work further augments this pattern [BLKS25, SYA25]. Coding and reasoning agents still struggle to sustain verifiable control over long action chains with external tools—planning, executing, and repairing across dozens of steps.

## 3 ALA framework

Our Agentic Lean Autoformalization (ALA) framework is centered around a generalist large language model (LLM) orchestrator that has access to a Lean 4-specialized model and multiple tools to improve the reliability of autoformalization. The orchestrator has three core abilities:

   (i) **Search for information and context:** The orchestrator can use `lean_retrieval` to fetch context and examples from a dedicated database that consists of theorems in natural

language, their translations to Lean 4, and explanations of the translations. Additionally, the orchestrator can use `search_online` to search the web for Lean 4-related code or documentation.

(ii) **Collect feedback:** The orchestrator can use `lean4_repl_runner` to compile Lean code and collect diagnostics via the REPL package [Com24]. It can also use the tool `check_theorem_tool` to construct a temporary Lean file, import `Mathlib`, and use the `#check` command to inspect the type of a definition, expression, or theorem.

(iii) **Query an expert:** The orchestrator can use `lean4_translation` to produce a Lean 4 declaration from natural language by prompting an LLM that has been fine-tuned in Lean 4.

Given a mathematical statement in natural language, the orchestrator interacts with the above resources until it produces Lean 4 code that compiles without errors or the number of tool calls reaches a bound given by the user. It then exports the Lean code. At this point, the candidate translation is sent to a reasoning-model LLM and presented to the user, who is assumed to be knowledgeable about the mathematical aspects of the definition and able to identify mathematically equivalent definitions written in different forms. The Lean code can be approved as a translation, rejected, or sent back to the orchestrator with feedback for future work by combining the LLM evaluation with a human evaluation as well.

# 4 A new database of upper-level theorems

There are several well-known databases of theorems produced by the autoformalization community. For example, FineLeanCorpus [m-a25, PYM$^+$25] contains 509,356 pairs of natural langauage with Lean 4 code; 1,181 from Omni–MATH [GSY$^+$24] (undergraduate and olympiad) and 45,853 from DeepMath–103K.

However, our agent has access to web-search, so to avoid contamination we exclude common datasets with informal mathematics whose statements already appear paired with Lean 4 code. [HLX$^+$25]. To minimize collisions, we chose examples from freely accessible repositories written by professors: Jiří Lebl's *Basic Analysis* and *Guide to Cultivating Complex Analysis* [Leb25a, Leb25b], Ben McKay's lecture notes on topology [McK25], and Stephen Doty's *Lecture Notes on Abstract Algebra* [Dot25]. For each source, we selected 100 examples by diversifying topic area and length. In total, our database consists of 400 examples, split evenly across undergraduate real, complex analysis, topology, and algebra.

# 5 Experimental setup

We evaluate ALA on our corpus of $N = 400$ theorems in algebra, topology, real analysis, and complex analysis, see Section 4. For each natural-language statement, the task is to produce a Lean 4 statement that type-checks in Mathlib and that it's a faithfull translation of the initial mathematical statement. Next, we describe the particularities of our experiments. We used Lean 4.22.0-rc4 compiler and `mathlib4` as dependency.

**Settings to test:** We compare three settings with a budget of 24 tool calls per problem. The baseline setting is the orchestrator LLM with access to all the tools described in Section 3. In our first variation, we modify the baseline setting by removing access to the LLM fine-tuned on Lean. In our second variation, we remove access to all tools besides the LLM fine-tuned on Lean 4 and the ability of the orchestrator to tell if a given Lean code compiles.

All methods use the same prompts and inputs. We record the number of calls used until orchestrator produces a Lean 4 statement that compiles; we also record the number of tool calls. We report pass rates, area-wise stratification, and Pass@$k$ over $k \in \{1, 6, 24\}$. We reset tool states between methods, fix random seeds, and log tool traces for paired analysis.

**Model selection:** For the generalist model, we use OpenAI 5.1 mini. For the LLM fine-tuned on Lean 4, we use Herald Translator [GWJ$^+$25]. For the final evaluation, we use the OpenAI 5.1 model.

**Databases:** For retrieving examples, we use a subset of 500 statements from the Herald database [GWJ$^+$25]. For testing, we use our database, see Section 4.

**Evaluation metrics:** We use the proportion of theorems that the agent successfully translated with fewer than $(K + 1)$ tools. We also consider the proportions of potential translations that compile as Lean code, but they may not be mathematically equivalent to the original statement.

# 6 Discussion of Experimental results

We found that an agentic approach is successful for translating mathematical statements to Lean. In particular, the use of tools had an impact on the success rate, on problems that require more tool calls to be translated, see Table 1. The full agent configuration translates 64 % of the theorems within 24 tool calls. This shows a significant improvement over the performance of Herald translator, 23%, 16% (Pass 128), and of Theorem LLama, 4 % and 2.9 % (Pass 128) for problems of a similar mathematical level.

Table 1: Success rate $SR@K$ for autoformalization $\pm 95\%$ Confidence interval

| Number of tools called | Agent with tools and expert LLM | Agent with tools but without expert LLM | Agent with expert LLM but without tools |
|---|---|---|---|
| 5 | $0.2050 \pm 0.0422$ | $0.2100 \pm 0.0426$ | $0.1950 \pm 0.0416$ |
| 10 | $0.3950 \pm 0.0486$ | $0.4375 \pm 0.0489$ | $0.3425 \pm 0.0478$ |
| 15 | $0.5225 \pm 0.0485$ | $0.5650 \pm 0.0478$ | $0.4150 \pm 0.0489$ |
| 20 | $0.6100 \pm 0.0465$ | $0.6100 \pm 0.04654$ | $0.4750 \pm 0.049$ |
| 24 | $0.6400 \pm 0.0455$ | $0.6725 \pm 0.0442$ | $0.5575 \pm 0.0417$ |

We also fit a Cox proportional hazards model to time-to-event data with a single binary indicator: the agent has access to all its tool configurations A.3. We found that the data is compatible with anything from a modest decrease (about 8.5%) to a moderate increase (about 36.4%), see Table 2.

## 6.1 Limitations

1. We have tried to minimize contamination; parts of our evaluation set may already be formalized on the web. Thus, the agent could "cheat" by retrieving solutions. Although the pipeline provides logs, we have not fully analyzed these mistakes.

2. Similar to Herald Translator, we use a baseline LLM call to judge faithfulness. The choice of this LLM matters and can yield false positives/negatives (see Appendix A.2). In future work we will consider specialist judges such as CriticLeanGPT [PYM+25].

3. The agent prompt can be further optimized, the dataset enlarged, and the RAG database extended to the full Herald set.

4. We were bottlenecked by compile checks: whenever the ALA calls certain tools, the loop blocks until they finish. In particular, `run_repl_tool` and `check_theorem_tool` dominate runtime— a Lean 4 REPL check takes about 30 seconds in our setup, and since `check_theorem_tool` invokes the REPL, it inherits the same cost.

# References

[BAMa] Bull. amer. math. soc. (n.s.). Chief editor articles.

[BAMb] Bulletin of the american mathematical society.

[BLKS25] Kaito Baba, Chaoran Liu, Shuhei Kurita, and Akiyoshi Sannai. Prover agent: An agent-based framework for formal mathematical proofs. *arXiv preprint arXiv:2506.19923*, 2025.

[Cas25] Davide Castelvecchi. Ai models solve maths problems at level of top students. *Nature*, 644:7, 2025.

[Com24] Leanprover Community. repl: a REPL for Lean 4. [https://github.com/leanprover-community/repl](https://github.com/leanprover-community/repl), 2024.

[dB70]    N. G. de Bruijn.  The mathematical language automath, its usage, and some of its extensions.  In M. Laudet, D. Lacombe, L. Nolin, and M. Schützenberger, editors, *Proceedings Symposium on Automatic Demonstration (Versailles, France, December 1968)*, volume 125 of *Lecture Notes in Mathematics*, pages 29–61. Springer, Berlin, 1970.

[dMU21]   Leonardo de Moura and Sebastian Ullrich. The lean 4 theorem prover and programming language (system description). In *Automated Deduction – CADE 28*, volume 12699 of *LNCS*, pages 625–635. Springer, 2021.

[Dot25]   Stephen R. Doty. Lecture notes on abstract algebra. `https://github.com/srdoty/AbstractAlgebraBook`, 2025. GitHub repository; accessed 2025-09-04.

[GCW+24]  Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xiangliang Zhang. Large language model based multi-agents: A survey of progress and challenges. In *IJCAI*, 2024.

[GSY+24]  Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo Miao, Qingxiu Dong, Chenghao Ma, Liang Chen, Runxin Xu, Zhengyang Tang, Benyou Wang, Daoguang Zan, Shanghaoran Quan, Ge Zhang, Lei Sha, Yichang Zhang, Xuancheng Ren, Tianyu Liu, and Baobao Chang. Omni-math: A universal olympiad level mathematic benchmark for large language models. *arXiv preprint arXiv:2410.07985*, 2024.

[GWJ+25]  Guoxiong Gao, Yutong Wang, Jiedong Jiang, Qi Gao, Zihan Qin, Tianyi Xu, and Bin Dong. Herald: A natural language annotated lean 4 dataset. In *The Thirteenth International Conference on Learning Representations*, 2025.

[HLX+25]  Zhiwei He, Tian Liang, Jiahao Xu, Qiuzhi Liu, Xingyu Chen, Yue Wang, Linfeng Song, Dian Yu, Zhenwen Liang, Wenxuan Wang, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. Deepmath-103k: A large-scale, challenging, decontaminated, and verifiable mathematical dataset for advancing reasoning. *arXiv preprint arXiv:2504.11456*, 2025.

[HTL25]   Junda He, Christoph Treude, and David Lo. Llm-based multi-agent systems for software engineering: Literature review, vision, and the road ahead. *ACM Transactions on Software Engineering and Methodology*, 34(5):1–30, 2025.

[Lea]     Lean Prover Community. Missing undergraduate mathematics in Mathlib. `https://leanprover-community.github.io/undergrad_todo.html`. Accessed: September 4, 2025.

[Leb25a]  Jiří Lebl.  Basic analysis: Introduction to real analysis.  `https://github.com/jirilebl/ra`, 2025. GitHub repository; accessed 2025-09-04.

[Leb25b]  Jiří Lebl. Guide to cultivating complex analysis: Working the complex field. `https://github.com/jirilebl/ca`, 2025. GitHub repository; accessed 2025-09-04.

[m-a25]   m-a-p. Fineleancorpus: A large-scale, high-quality lean 4 formalization dataset. Hugging Face, 2025. NLLean 4 pairs; accessed 2025-09-04.

[McK25]   Benjamin McKay.  Topology lecture notes.  `https://github.com/Ben-McKay/topology-lecture-notes`, 2025. GitHub repository; accessed 2025-09-04.

[PYM+25]  Zhongyuan Peng, Yifan Yao, Kaijing Ma, Shuyue Guo, Yizhe Li, Yichi Zhang, Chenchen Zhang, Yifan Zhang, Zhouliang Yu, Luming Li, Minghao Liu, Yihang Xia, Jiawei Shen, Yuchen Wu, Yixin Cao, Zhaoxiang Zhang, Wenhao Huang, Jiaheng Liu, and Ge Zhang. Criticlean: Critic-guided reinforcement learning for mathematical formalization. *arXiv preprint arXiv:2507.06181*, 2025.

[RN95]    Stuart Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach*. Prentice Hall, 1995. See Chapter 2: Intelligent Agents.

[SCB⁺24] Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2024. arXiv:2303.11366.

[SYA25] Peiyang Song, Kaiyu Yang, and Anima Anandkumar. Lean copilot: Large language models as copilots for theorem proving in lean. 2025.

[WDL⁺25] Ke Weng, Lun Du, Sirui Li, Wangyue Lu, Haozhe Sun, Hengyu Liu, and Tiancheng Zhang. Autoformalization in the era of large language models: A survey. *arXiv preprint arXiv:2505.23486*, 2025.

[WUL⁺25] Haiming Wang, Mert Unsal, Xiaohan Lin, Mantas Baksys, Junqi Liu, Marco Dos Santos, Flood Sung, Marina Vinyes, Zhenzhe Ying, Zekai Zhu, et al. Kimina-prover preview: Towards large formal reasoning models with reinforcement learning. *arXiv preprint arXiv:2504.11354*, 2025.

[WZJ⁺24] Ruida Wang, Jipeng Zhang, Yizhen Jia, Rui Pan, Shizhe Diao, Renjie Pi, and Tong Zhang. Theoremllama: Transforming general-purpose llms into lean4 experts. *CoRR*, 2024.

## A Technical Appendices and Supplementary Material

### A.1 Case study on a false positive due to foundational differences

We translated the following theorem that one-point compactification produces a topological space:

> Let $(X, \mathcal{T})$ be a Hausdorff topological space. Define $X' := X \cup \{\infty\}$ and
> $$\mathcal{T}' := \mathcal{T} \cup \big\{(X \setminus K) \cup \{\infty\} : K \subseteq X \text{ is compact}\big\}.$$
> Show that $\mathcal{T}'$ is a topology on $X'$.

Here is the compiling Lean file output from the agent:

```
import Mathlib

universe u
variable (X : Type u)
  [TopologicalSpace X] [T2Space X]

/-- Embed `X` into `X  PUnit`.
    Use `Sum.inr PUnit.unit` as . -/
def embed (x : X) : Sum X PUnit :=
  Sum.inl x

def infinity : Sum X PUnit :=
  Sum.inr PUnit.unit

/-- Opens for the one-point compactification
    modeled as `Sum X PUnit`.
    A set `s` is open iff
    (  s  `embed ¹› s` is open in `X`)
    (  s    K : Set X, `IsCompact K`
       `embed ¹› s = (Set.univ : Set X) \ K`). -/
def isOpen_one_point (s : Set (Sum X PUnit)) : Prop :=
  (infinity X   s    IsOpen (embed X ¹› s))
  (infinity X   s
     K : Set X, IsCompact K
       embed X ¹› s =

     (Set.univ : Set X) \ K)

/-- The topology of the one-point
    compactification on `Sum X PUnit`.
    (Proof omitted.) -/
def one_point_compactification_topology :
  TopologicalSpace
  (Sum X PUnit) := by
  sorry
```

**Why this is a false positive.** At first glance the translation looks faithful: it models $X'$ as Sum X PUnit via embed, names the adjoined point infinity, and defines the topology by isOpen_one_point. However, to inhabit TopologicalSpace (Sum X PUnit) one must *specify* the topology (i.e., bind isOpen := isOpen_one_point) *and prove* the topology axioms (i.e. inhabit isOpen_univ, isOpen_inter, isOpen_sUnion). The line := by sorry leaves both the choice of isOpen and the axioms unresolved. For a faithful translation, while we permit sorry for proofs we should not for *data*. To fix the translation, we should replace by sorry with

```
{ isOpen          := isOpenOnePoint X,
```

```
262     isOpen_univ   := by sorry,
263     isOpen_inter  := by sorry,
264     isOpen_sUnion := by sorry }
```

This issue stems from foundational challenges when converting from informal statements based on ZFC to formal code based on Lean's dependent type theory. We believe we can enlarge the database of examples by special cases like this, to improve performance.

## A.2  Case study on a false positive via agent cheating

We translated the following theorem from Jiří Lebl's *Guide to Cultivating Complex Analysis* [Leb25b]: for complex functions, differentiability implies analyticity (i.e., the existence of a Taylor expansion):

> If $f : U \to \mathbb{C}$ is holomorphic on a domain $U \subset \mathbb{C}$, then $f$ is analytic. That is, given $z_0 \in U$, $f$ has a Taylor series, and it converges in the disk $\{z : |z - z_0| < R\}$, where $R$ is the distance to the nearest singularity (possibly $R = \infty$). Furthermore, the sum of the series is $f$.

The ALA without to access to Herald gave the following compiling Lean output:

```
import Mathlib

theorem trivial_nat_eq (x : Nat) : x = x := by
  rfl
```

This code represents a clear tautology, that every natural number is equal to itself; it has no bearing with the original natural language statement. Interestingly, the ALA with access to Herald is giving a correct translation:

```
import Mathlib
open Complex

/-- If `f : U → C` is holomorphic on an open `U  C`,
then `f` is analytic at each `z  U`. -/
theorem holomorphic_on_analytic
  (U : Set C) {f : C → C} (hU : IsOpen U) (hf : DifferentiableOn  f U)
  (z0 : C) (hz0 : z0  U) :
  AnalyticAt C f z0 := by
  sorry
```

**Why this is happening.**  We can compare logs in the outputs folder with the results from access to Herald and without (our fine-tuned LLM). These logs indicate that the ALA without Herald made an incorrect call to `"path":"test.lean"`, erasing its previous attempts. These types of failures, of reducing to a degenerate proof, were rare (we found around 7 in the experiment where the ALA does not have access to Herald). We did not include the evaluator of faithfulness into the feedback loop; this might have kept the Agent on track.

## A.3  Statistical discussion

**Time-to-first-success analysis.**  We analyze *time to first successful compile*, measured in the discrete unit of *number of tool calls*. Each run belongs to one of two conditions: **ALA: Agent with access to all tools.** or a **Agent with access to Herald-only** condition. We fit a Cox proportional hazards model with a single binary covariate for condition (ALA$= 1$, Herald-only$= 0$), *stratified by theorem* so that each theorem has its own baseline hazard. Runs without a success by the administrative limit $K = 24$ calls are *right-censored at* $t = 24$; events that occur at $t = 24$ are counted as events (not censored). Because time is recorded in integer calls, we handle *tied event times* using the *Efron* method. Hazard ratios (HR) $> 1$ indicate faster success (fewer calls on average) for ALA relative to Herald-only. Model diagnostics included checks of the proportional-hazards assumption (global and covariate-specific Schoenfeld residual tests/plots). We report the number of strata (theorems), total runs, number of events, and the censoring proportion. Table 2 summarizes the fitted model.

Table 2: Cox proportional hazards regression results. HR = hazard ratio = $\exp(\text{Coef})$.

| Term | Coef | SE(Coef) | z | p | HR | HR 95% L | HR 95% U |
|---|---|---|---|---|---|---|---|
| Fine-tuned LLM | 0.111 | 0.102 | 1.087 | 0.277 | 1.117 | 0.915 | 1.364 |

In a theorem-stratified Cox model, the fine-tuned LLM showed a higher hazard of first successful compile (HR = 1.117, 95% CI [0.915, 1.364], $z = 1.087$, $p = 0.277$), implying an estimated 11.7% faster per-call success rate but with uncertainty spanning from 8.5% slower to 36.4% faster; thus the effect is not statistically significant across all possible number of tool calls.

## A.4 Agents workflows

We wrote the system prompt (see the Appendix A.6) to suggest one possible workflow to better generate high quality data, by give outline of translation process and contingency plan to handle unsuccessful translations. The agent has `max_step` times to use the tools, the agent will a return JSON flag if its did the writing of Lean 4 code into disk, and use `run_lean_tool` to verify if the Lean 4 code compiles. There might be cases that, during the last step the agent write a code but have not had chance to using `run_lean_tool` to evaluate, so after the agent finishing processing all input, we re-evaluate those cases whose `status` is `max_step_reached`. After agent running, we have a `csv` file, which columns are `name`, `step`, `status`, `passed`, `nl_statement`,`lean4_code`, then we re-evaluate those `status:max_step_reached` to fix the potential false negative.

For each row, we read the pair (`nl_statement`, `lean4 _code`), send to LLM judge (`GPT-5` with `reasoning="effort":"medium",`) to evaluates whether the Lean 4 code faithfully represents the natural language statement. We then augment the CSV with three new columns:

- `validate_score`: a base-10 numerical score indicating the degree of faithfulness,

- `validate_reason`: a free-form textual rationale explaining why the translation is (or is not) valid,

- `equivalent`: a Boolean flag (`True`/`False`) specifying whether the natural-language statement and Lean 4 code are judged equivalent. In our rubics, `True` only if the score is 10.

The below is the pseudo algorithm description:

---

**Algorithm 1** Lean4 translation agent (controller + post-processing)

---

**Require:** statement `nl_statement`, file path $p$, tools $\mathcal{T}$, step limit $S_{\max}$
1: $History \leftarrow [(\text{system}, \pi), (\text{user}, \text{"Translate "}x\text{" and save to } p)]$     ▷ $\pi$: system policy
2: **for** $s = 1$ **to** $S_{\max}$ **do**
3:     $resp \leftarrow \text{Model}(History, \mathcal{T})$     ▷ returns either content or a single tool call
4:     **if** `"status":"success"` $\in resp$.content **then**
5:        **return** SUCCESS     ▷ explicit success token
6:     **end if**
7:     **if** $resp$.tool_calls $\neq \emptyset$ **then**
8:        $(toolName, argument) \leftarrow$ first tool and its JSON args
9:        $result \leftarrow tool.\text{run}(argument)$
10:       $History \leftarrow History \cup [(tool, result)]$
11:       **if** $tool =$ `lean4_repl_runner` **and** $result$.repl_pass $= 1$ **then**
12:          **return** SUCCESS     ▷ auto-stop on REPL pass
13:       **end if**
14:     **else**
15:       $History \leftarrow Hisotry \cup [(\text{assistant}, resp.\text{content})]$
16:     **end if**
17: **end for**
18: **return** MAXSTEPREACHED     ▷ may have written code but not verified

---

---

**Algorithm 2** Post-processing: REPL re-check and LLM judging

---

**Require:** CSV with columns `name, step, status, passed, nl_statement, lean4_code`
  1: **for each** row $r \in \mathcal{C}$ **do**
  2:     **if** $row$.status $=$ max_step_reached **then**
  3:         $result \leftarrow$ lean4_repl_runner($row$.lean4_code)
  4:         update $row$.passed $\leftarrow (row$.repl_pass $= 1)$
  5:     **end if**
  6: **end for**
  7: **for** each row $r \in \mathcal{C}$ **do**
  8:     $(\hat{y}, \rho) \leftarrow$ GPT-5-JUDGE( $row$.nl_statement, $row$.lean4_code )
  9:     $r$.judge_score $\leftarrow \hat{y}$;   $r$.judge_rationale $\leftarrow \rho$
 10: **end for**
 11: **return** updated $\mathcal{C}$

---

## A.5  Prompt for evaluating correctness of translation

Compiling Lean 4 code does not guarantee that the translation is correct; it can pass for the reasons
outlined in Appendices A.1 and A.2. Following Herald Translator [GWJ$^+$25], we employ an LLM
judge to evaluate faithfulness. We use the following 1-shot CoT prompt with GPT-5 (reasoning mode:
medium) for evaluating faithfulness.

```
You are an expert in Lean 4, mathlib, and mathematics. You are an
auditor with guidelines.

Instructions:
Your input is (A) compiling Lean 4 code and (B) a natural-language
statement. Decide whether (A) faithfully formalizes (B). Do not use
proof quality; only check statement fidelity.

Think step by step:
1) Translate each line of the Lean 4 code into plain language. Check
   if it is sensible and on track.
2) Then decide if the whole Lean statement is faithful to the original.
3) Final check: are the two math statements the same or different?
   Point out any differences precisely.

Guidelines:
1) It must be a legitimate, faithful translation to pass. Small
   formalization differences are fine. Since the code compiles, assume
   referenced names exist in mathlib.
2) Prefer current/standard mathlib terms; ad-hoc encodings can be a
   red flag if they change meaning.
3) If the Lean code introduces auxiliary definitions (beyond the final
   theorem/definition), they must not be vacuous. If any auxiliary
   definition is vacuous (e.g., ':= True', ':= none', or filled with
   'sorry' where data is required), the translation fails. The aux
   definition must describe what it is trying to say.
4) Only if each auxiliary definition is legitimate and the final Lean
   statement matches the original in mathematical meaning does it pass.
   Do not penalize harmless formal phrasing differences.
5) If it is a near pass, assess whether the Lean 4 statement is a good
   formalization of the original. Slight specialization/generalization
   is acceptable if no substantive error is introduced.

After you finish your reasoning:
Assign a Grade  {0,...,10}. Use this rubric:
0: completely unrelated
3: vacuous aux defs and even fixing them would not make it faithful
```

```
373  6: vacuous aux defs, but fixing them would make it faithful
374  9: almost the same but still not faithful
375  10: faithful
376
377  Also output:
378  - COT inside "### BEGIN THOUGHT" / "### END THOUGHT".
379  - Faithfulness (binary) immediately after "### FAITHFUL SCORE"
380    where 0 = not faithful, 1 = faithful.
381  - The numeric grade immediately after "### GRADE".
382
383  ---
384
385  ### Example
386
387  Here is the Lean 4 code:
388  ```lean
389  import Mathlib
390
391  universe u v
392  variables {X : Type u} {Y : Type v}
393    [TopologicalSpace X] [TopologicalSpace Y]
394
395  /-- Placeholder for a covering map. -/
396  def CoveringMap (p : X → Y) : Prop := True
397
398  /-- Placeholder: U is evenly covered by p. -/
399  def evenly_covered (p : X → Y) (U : Set Y) : Prop := True
400
401  /-- Number of sheets (none = ). -/
402  def num_sheets (p : X → Y) (U : Set Y) : Option Nat := none
403
404  /-- Placeholder for path connectedness. -/
405  def PathConnected (Y : Type v) [TopologicalSpace Y] : Prop := True
406
407  namespace Covering
408
409  theorem sheets_equal_on_overlap {p : X → Y} (hp : CoveringMap p)
410    {U V : Set Y} (heU : evenly_covered p U) (heV : evenly_covered p V)
411    (hnonempty : (U  V).Nonempty) :
412    num_sheets p U = num_sheets p V := by sorry
413
414  theorem covering_map_n_to_one_of_path_connected {p : X → Y}
415    (hp : CoveringMap p) (hpath : PathConnected Y) :
416     (n : Option Nat),  (y : Y),  (U : Set Y),
417      y  U  IsOpen U  evenly_covered p U  num_sheets p U = n := by
418    sorry
419
420  end Covering
```

Note, the grade is an artificial value not used in the final analysis. It was a book-keeping device for us to keep track of uncertainty. A grade 0 happens usually only when the proof degenerate into a triviality as in Appendix A.2. A grade 9 sometimes happens for false negatives (correct translated code that was judged too harshly by this evaluator). If the trranslation is deemed faithful, the LLM outputs a faithful score of 1; otherwise it outputs 0.

**A.6   Agent Prompts**

We provide four system prompts, one for each ALA configuration we tested: (1) full ALA (all tools, including the specialist LLM Herald), (2) ALA without the specialist LLM, (3) specialist LLM only,

429  and (4) ALA without REPL and without the specialist LLM. The exact prompt strings and tool-call
430  templates are available in the anonymized code repository.

**ALA with access to tools including the specialist LLM**

```
You are an expert Lean4 programmer-agent.
Translate the NL math statement into ONE Lean4
declaration that compiles with Mathlib.
Translation only - **not a full proof**.

Always call 'lean4_translation' FIRST. Then write
to disk and verify with 'lean4_repl_runner'.
Pass = '1' (compiles); Fail = '0' (does not).

When translating:
- add 'import Mathlib' at the top
- end the decl with ':= by sorry' (no proof)

## Process
1) (Optional, once) 'lean_retrieval' for an example.
2) Translate: use 'lean4_translation' or draft manually
   (always end with ':= by sorry').
3) Write & verify: 'lean_write_file' → 'lean4_repl_runner'.
4) If 'repl_pass = 1' → respond '{ "status": "success" }',
   else revise and retry.

## Contingency (errors)
- Unknown names → 'lean_check_theorem'.
- Syntax/tactics → 'search_online'.
- Re-test → 'lean4_repl_runner'; iterate.

## Naming (Lean4/mathlib)
- Types/Props: PascalCase
  e.g., 'IsSimpleGroup', 'IsCyclic', 'Nat.Prime'
- Lemmas/Functions: snake_case
  e.g., 'Nat.add_comm', 'List.map'
- Be specific: prefer
  'Sylow.exists_subgroup_card_pow_prime'
  over vague labels like "Sylow Theorem".
```

**ALA with access to tools except the specialist LLM**

```
You are an expert Lean4 programmer-agent. Translate the given
natural-language statement into a single Lean4 declaration that
compiles with Mathlib. Translation only, not a full proof.

Draft the statement yourself (no specialist translator). Write
it to disk and verify with 'lean4_repl_runner'. Pass = 1, fail = 0.

When translating, import 'Mathlib' at the top and end with
':= by sorry' (no proof).

Process
1) (Optional, once) 'lean_retrieval' for an example.
2) 'lean_write_file' → 'lean4_repl_runner'.
3) If 'repl_pass = 1', respond '{ "status": "success" }';
   else revise and retry.

Errors
```

```
484  - Unknown names: `lean_check_theorem`.
485  - Syntax/tactics: `search_online`.
486  - Re-test with `lean4_repl_runner` and iterate.
487
488  Naming (Lean4/Mathlib)
489  - Types/Props: PascalCase  (e.g., `IsSimpleGroup`, `Nat.Prime`)
490  - Lemmas/Fns:  snake_case  (e.g., `Nat.add_comm`, `List.map`)
491  - Prefer specific names (e.g., `Sylow.exists_subgroup_card_pow_prime`)
```

**492  ALA without access to any other tools except the specialist LLM**

```
493  You are an expert Lean4 programmer-agent.
494  Translate the given NL statement into ONE
495  Lean4 declaration that compiles with Mathlib.
496  Translation only - not a full proof.
497
498  Use `lean4_translation` to draft the declaration,
499  then write it to disk with `lean_write_file`.
500
501  When translating:
502  - add `import Mathlib` at the top
503  - end the decl with `:= by sorry` (no proof)
504
505  When finished, respond with:
506  { "status": "success" }
```

**507  ALA without REPL and without the specialist LLM**

```
508  You are an expert Lean4 programmer-agent. Your mission is to translate
509  the given natural-language statement into a single Lean4 declaration.
510  Your goal is translation only, not a full proof.
511
512  After generating the code, write it to disk with `lean_write_file`.
513
514  When translating, import `Mathlib` at the top and end the declaration
515  with `:= by sorry` (no proof).
516
517  ## Tools
518  - `check_theorem_tool`: check existence / canonical names.
519  - `lean_write_file`: write code to disk.
520  - `lean4_translation`: draft a declaration (no proof). You may use it,
521  but verify syntax/names with other tools; do not rely on it alone.
522  - `lean_retrieval`: fetch similar (NL, Lean) example pairs.
523
524  ## Naming
525  1. Types/Props: PascalCase (e.g., `IsSimpleGroup`, `Nat.Prime`).
526  2. Lemmas/Functions: snake_case (e.g., `Nat.add_comm`, `List.map`).
527  3. Be specific: prefer `Sylow.exists_subgroup_card_pow_prime`.
528  4. Confirm names with `check_theorem_tool`.
529
530  Respond with: `{ "status": "success" }` once the translation is written.
```

# NeurIPS Paper Checklist

1. **Claims**

   Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

   Answer: [Yes]

   Justification: We introduce the ALA framework (Sec. 3), a 400-example upper-division/graduate NL dataset (Sec. 4), and an empirical evaluation on four Lean 4 domains (Secs. 5, 6). The scope is bounded to the specified models, Mathlib library, and a 24-call budget.

   Guidelines:

2. **Limitations**

   Question: Does the paper discuss the limitations of the work performed by the authors?

   Answer: [Yes]

   Justification: We discuss potential web contamination of the evaluation set, possible false positives/negatives from the LLM-based faithfulness evaluator, sensitivity to the agent prompt, dataset size, and RAG pool, runtime constraints from the REPL tool, and a post-processing gap (only max-step cases are rechecked). See Sec. 6.1.

3. **Theory assumptions and proofs**

   Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

   Answer: [NA]

   Justification: The paper introduces a system and reports empirical results. It does not present new theorems or proofs; it does include a case study of generated formal code in the appendix A.1)

4. **Experimental result reproducibility**

   Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

   Answer: [Yes]

   Justification: Section 3 specifies the pipeline (orchestrator + tools) and the controller/loops (Algorithms 1, 2). Section 5 explains datasets , model choices (including Herald Translator), the 24-call limit on tools, and the success criterion used to compute Pass@$k$.

5. **Open access to data and code**

   Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

   Answer: [Yes]

   Justification: Our anonymized GitHub repo contains the runnable pipeline (code, tools, configs, inputs, instructions, and requirements). We use Lean's REPL feedback tool, a secondary evaluator LLM, Herald Translator, and a RAG database of 500 examples.

   Guidelines:

   - The answer NA means that paper does not include experiments requiring code.
   - Please see the NeuroIPS code and data submission guidelines (`https://nips.cc/public/guides/CodeSubmissionPolicy`) for more details.
   - While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
   - The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (`https://nips.cc/public/guides/CodeSubmissionPolicy`) for more details.

- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. **Experimental setting/details**

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The paper specifies the evaluation dataset (Secs. 4, 5); the exact models/tools (generalist LLM + Herald; REPL, RAG examples, Mathlib check, and web search via the Serper API) are available in the anonymized code repository. We use Lean 4 (4.15.0), a 24–call budget, and fixed seeds/prompts/configs (see Appendix A.5, A.6).

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. **Experiment statistical significance**

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: For Pass@$k$ we plot 95% confidence intervals per $k$ as binomial CIs over $N{=}400$ problems. These appear as the error bars in our figures.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. **Experiments compute resources**

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: For the tool `herald_translator_tool`, it's take about 14GB GPU memory to load the model. We used a virtual machine on Goole Cloud with CPU `a2-highgpu-1g` (`12 vCPUs, 85 GB Memory`) with GPU `NVIDIA A100 40GB` to serve the model using `vllm`.

For the agent running, we are running on the script over a `Apple M3 Pro` with 18GB memory. It take 1 to 1.5 hour to go over the 400 theorems depending on the selection of tools. Details discussed in the limitation.

For the judgment script, since we asked the model GPT-5 with thinking mode `medium`, the average time for generating score and reasoning is around 20 mins and cost around $14 for go over 400 (`nl_statement, Lean4_code`), the log file is under `all_experiments_csv/record.txt`.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. **Code of ethics**

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We use only openly licensed educational materials and tools made by the Lean 4 community; no personal or sensitive data are involved. For our dataset, we use Lebl, McKay, and Doty's repositories [Leb25a, Leb25b, McK25, Dot25]; we respect upstream licenses and terms of use.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. **Broader impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: Positive impacts include improving formalization for these target areas (upper-division undergraduate and graduate mathematics); risks include unfaithful translations. We mitigate this via an LLM evaluator but mistakes can still happen.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. **Safeguards**

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: We do not release high-risk assets; our artifacts are code and a small text dataset of mathematical statements from openly licensed academic sources.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. **Licenses for existing assets**

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We credit all third-party assets in the paper; here are their licenses: Lean 4/Mathlib/REPL tool (Apache-2.0), Herald Translator (Apache-2.0; weights not redistributed), Lebl notes (CC BY-SA / CC BY-NC-SA), McKay topology notes (GPL-3.0), and Doty algebra notes (MIT for code; CC BY 4.0 for text). Their terms are stated and respected.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, `paperswithcode.com/datasets`

has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: We introduce an agent and a 400-example dataset; both are documented in the anonymized GitHub (see the README and LICENSES).

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and research with human subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional review board (IRB) approvals or equivalent for research with human subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. **Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: LLMs are core to our method: a generalist LLM orchestrates tool calls (REPL feedback, RAG retrieval, web search, Mathlib checks); it invokes a Lean 4–specialized translator (Herald Translator) and a secondary LLM evaluates faithfulness. Outside of this core pipeline, an LLM assisted dataset selection, picking 400 problems based on diversity/length (Secs. 4), and 500 RAG examples.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.