## Lean Finder: Semantic Search for Mathlib That Understands User Intents

**Anonymous Authors**<sup>1</sup>

#### Abstract

We present Lean Finder, a semantic search engine for Lean and mathlib that understands mathematicians' intents. Due to challenges of locating relevant theorems and the steep learning curve of Lean 4 language, the progress of formal theorem proving is slow by tedious human efforts. Recent Lean search engines, though helpful, only passively consider informalization of statements, largely overlooking the discrepancy from user queries in the real world. In contrast, we propose a user-centered semantic search tailored to the needs of working mathematicians. The key idea is to first analyze and cluster the semantics of public discussions on Lean, then fine-tune text embeddings on synthesized queries that simulate user intents. Our Lean Finder is thus encoded with a rich awareness of mathematicians' intents from different perspectives. Evaluations on both realworld queries by mathematicians and informalized statements demonstrate that our Lean Finder outperforms previous search engines by at least 19%. We promise to release both the code, model checkpoints, datasets, and the web service for our Lean Finder upon acceptance.

036

### 36 1. Introduction

Advances in Lean and mathlib (De Moura et al., 2015; Moura & Ullrich, 2021) are turning mathematical discovery into a collaborative and verifiable research workflow. On this foundation, provers powered by large language models (LLMs) have sprinted ahead (AlphaProof & teams, 2024; Xin et al., 2024a;b; Ren et al., 2025; Lin et al., 2025). Parallel progress in autoformalization has shifted from handwritten grammars to few-shot LLM translation, and large corpora by back-translating Lean code ("informalization") have bootstrapped even stronger translators (Wu et al., 2022; Jiang et al., 2024; Lu et al., 2024).

<sup>149</sup> <sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, <sup>550</sup> Anonymous Country. Correspondence to: Anonymous Author <sup>551</sup> <anon.email@domain.com>.

<sup>53</sup> Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.



**Figure 1:** In the evaluation of real user queries, our Lean Finder is preferred in 42.3% of cases, compared to only 23.1% for LeanSearch (Gao et al., 2024a).

Despite these advances, state-of-the-art LLMs still cannot solve math research problems. This underscores the continued need for substantial **human efforts**, which unfortunately remain slow and labor-intensive for **two bottlenecks**. First, locating the right lemma or theorem is frustrating: search tools are rudimentary, naming conventions are often inconsistent, and high-quality Lean examples remain scarce. Second, Lean's syntax, grammar, and tactics incur a steep learning curve. Even veteran mathematicians and experienced programmers regularly report difficulties when writing Lean (Zulip, 2021b; 2020b; 2021a; 2020a).

To facilitate mathematicians, **search is vital to Lean formalization**. Recent engines take informal statements or live proof states and retrieve mathlib4 lemmas or tactic suggestions, thereby aiding human Lean users (Gao et al., 2024a;b; Tao et al., 2025; Shen et al., 2025; Ju & Dong, 2025; Asher, 2025). However, these search engines target machine translation rather than human use. They "informalize" statements from a supposedly neutral viewpoint, whereas **real users bring inherent bias**, typically seeking explanations from their own specific perspective.

Consider the two queries below. The first is an informalization of a formal statement that current Lean search engines handle (Gao et al., 2024a;b; Ju & Dong, 2025; Asher, 2025):

**Query 1:** Let L/K be a field extension and let x, y in L be algebraic elements over Kwith the same minimal polynomial. Then the K-algebra isomorphism algEquiv between the simple field extensions K(x) and K(y) maps the generator x of K(x) to the generator yof K(y), i.e., algEquiv(x) = y.

#### **Target Statement 1:**

<sup>1</sup> theorem algEquiv\_apply {x y : L} (hx : IsAlgebraic K x) (h\_mp : minpoly K x = minpoly K y) : algEquiv hx h\_mp (AdjoinSimple.gen K x) = AdjoinSimple.gen K y := by



**Figure 2:** Overview of building our Lean Finder. We first analyze and cluster semantics of public queries by mathematicians (Section 3.1.1), and then fine-tune LLM embeddings via a large amount of synthesized queries that simulate user intents (Section 3.1.2). Moreover, our Lean Finder is further made more research-centered by including Lean 4 code from research papers and repositories (Section 3.2).

```
072 2 have hy : IsAlgebraic K y := {minpoly
    K x, ne_zero hx.isIntegral, (h_mp ▷
    aeval _ _)
    Tw [algEquiv, trans_apply, ←
    adjoinRootEquivAdjoin_apply_root K
    hx.isIntegral, symm_apply_apply,
    trans_apply,
    AdjoinRoot.algEquivOfEq_apply_root,
    adjoinRootEquivAdjoin_apply_root K
    hy.isIntegral]
```

#### 082 However, a human query in practice may look like:

#### Target Statement 2:

068

069

070

081

091

092

093

094

095

096

097

098

```
1 theorem eq_of_root {x y : L} (hx :
        IsAlgebraic K x) (h_ev :
        Polynomial.aeval y (minpoly K x) =
        0) : minpoly K y = minpoly K x :=
2 ((eq_iff_aeval_minpoly_eq_zero
        hx.isIntegral).mpr h_ev).symm
```

099 Although both queries involve a similar mathematical con-100 text (algebraic elements x, y in a field extension L/K and their minimal polynomials), they concern different purposes: the first asserts the explicit behavior of an isomorphism between field extensions generated by two algebraic 104 elements with the same minimal polynomial, whereas the 105 second seeks to decide whether the two minimal polynomi-106 als are equal from the fact that y is a root of x's minimal polynomial. This user *latent* (motivation, perspective, level of abstraction) cannot be inferred or encoded by a purely 109

syntactic informalization. Addressing this challenge calls for Lean search engines that can understand a mathematician's intent and discourse style, not merely surface-level matches. We defer a more rigorous analysis in Section 2.2.

Therefore, we ask our core question:

How to make Lean search engine user-centered and tailored to understand the intents and needs of working mathematicians?

In this paper, we aim to build a user-centered search engine for Lean and mathlib4 that can understand the diverse intents of mathematicians. The key method is to first **analyze** and cluster the semantics of public queries by mathematicians, and then fine-tune LLM embeddings via a large amount of synthesized queries that simulate user intents (Section 3.1, Section 3.5). This fine-tuning strategy can enrich the awareness of real user intent from different perspectives. Moreover, our Lean Finder is further made more research-centered by including Lean 4 code from research papers and repositories (Section 3.2). We test our Lean Finder via a vast amount of retrieval evaluations with both informal statements and *real* user queries. We will release our Lean Finder retrieval system as a web service to support mathematicians coding in Lean 4. A demo is shown in Figure 4. We summarize our contributions below:

- 1. When evaluating on real-world queries by human mathematicians, our Lean Finder is upvoted 42.3% compared to 23.1% from LeanSearch in an LMArena-style test (Chiang et al., 2024).
- 2. Queried by informalized statements, Lean Finder also outperforms LeanSearch (Gao et al., 2024a) by 27.3% relative improvement on Recall@1.
- 3. We release the largest dataset for informal-formal pairs, with 1222600 synthesized user query, 244521 informalized statements, and 254811 Lean code snippets.

## 110 2. Motivation

111

112

#### 2.1. Background: Search Engine for Lean and Mathlib

113 Lean (De Moura et al., 2015; Moura & Ullrich, 2021) is an 114 interactive theorem prover based on dependent type theory 115 that lets mathematicians write proofs in a rigorously check-116 able language. Its community-run library mathlib4 now 117 contains over 210k theorems and 100k definitions, orders of 118 magnitude more formal "entry points" than most program-119 ming platforms (Python exposes only 89 built-in functions; 120 while the C++ standard library has 105 headers). Finding 121 the right item inside the ocean of mathlib4 is hard: the 122 official #find, library search, and Loogle tools rely on exact 123 names or goal states and quickly miss relevant results when 124 naming conventions drift or examples are sparse. Recent 125 LLM-powered search engines (Gao et al., 2024a;b; Yang 126 et al., 2023; Tao et al., 2025; Shen et al., 2025; Ju & Dong, 127 2025; Asher, 2025; Zhu et al., 2025) embed informal queries 128 or proof goals and retrieve matching Lean statements, but 129 they are optimized for embedding similarity, not for the 130 rich, shifting intents of working mathematicians. In short, 131 sheer library scale and rudimentary search together make 132 locating lemmas a first-order pain point, even before one 133 tackles Lean's non-trivial syntax and tactic language. 134

# 135 136 137 2.2. Current Search Engines Are Misaligned with the Practical Needs of Mathematicians

Despite notable progress, current Lean search engines are 138 optimized for informalized statements, instead of questions 139 mathematicians actually pose. Most systems embed infor-140 mal descriptions (natural languages) of statements or proof 141 states, together with possible context or dependencies, then 142 retrieve matching lemmas or tactics (Yang et al., 2023; Gao 143 et al., 2024a;b; Tao et al., 2025; Shen et al., 2025). Because 144 these informalizations are generated by LLMs rather than 145 user logs, we have no guarantee that their wording, granularity, and purpose align with the queries real Lean practition-147 ers type. This unexamined domain gap between LLM-148 generated paraphrases and authentic human queries 149 has so far been largely ignored, leaving retrieval systems 150 potentially well-tuned to wrong distributions. 151

152 To justify this domain gap, we visualize distributions of 153 different queries embeddings generated via all-MiniLM-L6-154 v2 (Wang et al., 2020) sentence-transformers and PCA. As 155 shown in Figure 3, the cluster of LLM-generated informal-156 izations (blue) only spans over a subspace of the cluster 157 by queries collected from Zulip and Github (green), high-158 lighting a clear domain gap between real user queries and 159 surrogate ones on which current Lean search engines are 160 trained. By contrast, embeddings of our synthesized user 161 queries (Section 3.1) are much more aligned with the hu-162 man cluster, suggesting that they better capture the practical 163 needs of Lean mathematicians. 164



**Figure 3:** PCA of different queries. User queries are from Zulip and Github; Lean Finder considers our synthesized user query (Section 3.1), and informalization is from Herald (Gao et al., 2024b)

#### 3. Methods

#### 3.1. User-centered Query Synthesis

Section 2.2 underscores that effective Lean retrieval for mathematicians requires authentic user queries, yet assembling such data at scale is obstructed by two bottlenecks:

- The volume of publicly available Lean questions by real users is tiny (for example, we only fetched 693 answerable user queries from Zulip/GitHub), which is orders of magnitude smaller than the billions-token corpora consumed by modern LLMs;
- 2. Even when such real user queries exist, tracing the precise answer (formal statement) that eventually resolves each query is often infeasible, due to open math problems and evolving Lean/mathlib4 development.

Therefore, preparing a large fully annotated set of genuine user queries is unrealistic. Instead, we propose a novel reverse strategy to synthesize a large amount of diverse and realistic queries. The **core idea** is: based on different perspectives that mathematicians ask Lean/mathlib4-related questions (Section 3.1.1), we prompt LLM generation to simulate mathematicians' intents (Section 3.1.2).

#### 3.1.1. SEMANTIC ANALYSIS OF REAL HUMAN QUERY

The first step for user query synthesis is to systematically collect real user discussions, and then analyze the diverse semantics and intents of mathematicians coding in Lean 4.

**1.** Collect User Discussions. To achieve this semantic analysis, we utilize Lean Zulip Chat<sup>1</sup>, the primary discussion forum for Lean 4, as well as GitHub, as the main sources of real-world queries. Using the Zulip API, we extract user discussions, prioritizing high-quality and relevant content by targeting five highly active and thematically

<sup>&</sup>lt;sup>1</sup>https://leanprover.zulipchat.com/

relevant channels: *new members, lean4, mathlib4, Is there code for X,* and *metaprogramming/tactics.*

For each thread within these channels, we retain up to the first five user messages (sorted chronologically) as they generally center around the context and core idea, and subsequent replies are more likely off the topic.

Subsequently, we prompt GPT-40 to filter these messages,
retaining only discussions answerable by a formal Lean 4
statement, and paraphrasing the main query articulated
within the initial messages. This procedure yields a refined collection of 693 high-quality discussions (600 from
Zulip<sup>2</sup>, 93 from GitHub).

Cluster User Intents. Mathematicians may ask 181 2. Lean/mathlib4-related questions from different perspectives 182 183 (see examples in Section 1 and our analysis in Figure 3), therefore we need to identify the distinct types of queries 184 185 collected from Zulip/Github. We first prompt OpenAI o3 to bootstrap the initial set of clusters from a subset of collected 186 queries. Then, the same model is fed with more user discus-187 sions and is prompted to iteratively update existing clusters 188 or introduce new clusters if necessary. See Appendix E for 189 prompts we used. This process results in five distinct and 190 191 semantically meaningful clusters of queries, as shown in Table 1, capturing a comprehensive range of intents from which mathematicians ask about Lean/mathlib4. 193

10/

218

219

179

180

#### 195 3.1.2. QUERY SYNTHESIS WITH USER INTENTS

196 Based on semantic clusters of real queries in Table 1, we 197 synthesize a large-scale dataset of queries simulating diverse 198 user intents. Although it is typically challenging to resolve a 199 query with precise formal statements in Lean, we can still re-200 versely synthesize the query, instead of the ground truth. 201 Specifically, we assume each formal Lean/mathlib4 state-202 ment could be the ground-truth answer to some unknown 203 user queries, and now the key is to simulate or synthesize 204 queries of realistic user intents.

206 We aim to instruct the LLM to generate plausible user queries that align with the specified cluster's perspective while remaining grounded in the given formal content. Con-208 209 cretely, we prompt GPT-4.1 mini with: (i) the Lean 4 formal 210 statement, (ii) its informalization (Section 3.4), (iii) cluster 211 information (cluster name, semantic definition, and representative 10-shot examples). See Figure 9 for our prompt. 212 213 For each formal statement, we iterate over all five clusters 214 in Table 1. This procedure yields a total of 1,222,600 syn-215 thetic queries distributed across five distinct intent clusters, 216 enabling fine-tuning of embedding models better to meet mathematicians' needs on Lean search. 217

#### 3.2. Research-driven Data Collection

While mathlib4 remains the primary source of formal statements for most Lean 4 search engines, it does not capture the full breadth of formal mathematics. Many recent mathematical developments are formalized beyond mathlib4, often in repositories associated with research papers or domainspecific libraries. These statements can be more relevant to working mathematicians, especially when formalizing new results or building upon cutting-edge works.

To build a more comprehensive and representative corpus of Lean 4 statements, we consider the following sources in addition to mathlib4:

- **Research-linked repository:** GitHub projects that are part of math papers, which often contain formalizations of novel theorems and lemmas.
- **Domain-specific library:** Projects such as SciLean that cover applied mathematical domains.
- **Transitive dependency:** Any Lean 4 project dependencies required by the above repositories.

We use LeanDojo (Yang et al., 2023) to extract formal statements from collected projects. The complete list of repositories in our data is shown in Table 10 in the Appendix.

This expanded collection ensures that our search engine is equipped with the most relevant, diverse, and latest formal statements to support both Lean practice and math research.

#### 3.3. Informalization

To support retrievals by both informalized statements and the generation of synthetic user queries (Section 3.1.2), we convert formal Lean 4 statements into natural language descriptions. Inspired by Herald (Gao et al., 2024b), for each Lean 4 statement in our dataset, we provide rich context information for GPT-40 to informalize. The context consists of the following six components:

- Formal name: The full name of the statement to be informalized.
- Formal statement code: The complete Lean 4 code, including the statement header and body.
- **Docstring:** Any human-written docstring associated with the statement in the codebase.
- **Neighbor statement:** The closest statement in the same file, based on positional proximity. We include its full code, as nearby statements often share related concepts. This helps the model understand the local context and inter-connectiness of the statement.

<sup>&</sup>lt;sup>2</sup>Due to privacy concerns, we do not release Zulip collections.

The state of the s	Table 1: Five clusters of mathematicians	intents categorized from our colle	ctions of real user queries t	from Zulip and GitHub.
--	--	------------------------------------	-------------------------------	------------------------

Intent Cluster	Semantic Definition	Examples
Searching for existing code / lemmas	Whether a definition, data-structure implementation, lemma, or theorem is already available in Mathlib/Lean, or if there is an easy way to get the desired statement.	"Is there code for", "Do we have", "Is there a lemma which", "Where can I find"
Meta- / tactic- programming questions	Questions about Lean 4 metaprogramming: writing tac- tics or elaborators, creating macros, manipulating Expr or environments, controlling metavariables, interacting with 'simp'/'aesop'/'omega' from meta code.	"Why doesn't the code give an error about a cycle existing when using Lean.MVarId.assign?", "Is there a way to make my first approach work to get Exp out of Q(Exp) using expr% macro?"
Type-class, in- stance, axiom	Proofs fail, or definitions cannot be written because Lean cannot find (or should not use) certain type-class instances, or users need advice on constructing/avoiding such instances and on the correct logical meaning of such instances.	"How to define or derive instances", "Why certain instance-search problems happen", "Whether particular instances should exist at all"
Proof engineer- ing & everyday Lean usage	Concrete goals or failing scripts around practical, day-to-day proof writing in Lean.	How to: finish or shorten proofs, make simp, rw, omega, linarith, etc. succeed, rewrite or unfold expressions, handle coercions and subtypes, manipulate logical connectives $(\forall, \exists, \land, \lor)$ , pattern-match over inductives, or split cases.
Library design & large-scale formalization	Conversations that concern the construction of large-scale state- ments: proposing new mathematical structures or tactics, refac- toring existing ones, performance trade-offs, big formalizations or theorems.	"How can I apply the coYoneda lemma or density theorem to simplicial sets valued in an arbitrary universe without restricting the variable?"

 Dependent statements: All the dependent statements used in the body of the statement to be informalized, along with their informalizations. This provides the necessary prerequisite knowledge and semantic dependencies for achieving high-quality informalization. See Appendix A.2 for details about acquiring dependent statements.

• **Related statement in Herald** (Gao et al., 2024b): A formal statement and its corresponding informalization from the Herald dataset that closely aligns with the target statement. This example serves as an in-context demonstration for GPT-40, helping it generate a more accurate and stylistically consistent informalization. See Appendix A.2 for details on retrieving related statements.

#### 3.4. Other Input Modalities

For the Lean code snippet input modality, we segment statement body using newline characters (\n). This slicing strategy helps preserve the semantic structure of the original proof. The input sequences in our dataset contain 2 to 4 code segments, omitting single-line inputs, as individual lines are often trivial (e.g., rfl) and provide limited semantic value on their own. The statement definition input modality excludes the body of the statement, and retains only the statement header information.

We summarize our dataset by input modalities in Table 2 and by data sources in Table 3.

#### 3.5. Training Lean Finder with Code Search

Our code search model, named **Lean Finder**, adopts DeepSeek-Prover-V1.5-RL 7B (Xin et al., 2024b) as the base model for fine-tuning, for its extensive pretraining on Lean 4 syntax and theorem proving tasks. As a decoder-

Table 2: Overview of our dataset by input modalities.

Input Modality	Samples	Percent
Synthesized User Query	1222600	63%
Informalized Statements	244521	12%
Lean Code Snippet	254811	13%
Formal Statements	244521	12%
Table 3: Overview of our d	ataset by da	ata sources
Source	Samples	Percent
Mathlib	238021	97%
Research-linked repository	2198	1%

Lean 4 related libraries

only architecture, only the final token in each sequence has access to the full context due to the causal self-attention. Therefore, we extract the final hidden state of the last token in the last decoder layer to embed the entire sequence. The model is fine-tuned with contrastive loss on "informal query  $q_i$  – formal code  $c_j$ " pairs, which aims to align the embeddings of matching pairs in a shared embedding space.

2%

For each training pair  $(q_i, c_j)$ , we create an augmented version  $q_i^* = \text{Augment}(q_i)$ ,  $c_i^* = \text{Augment}(c_i)$ . The augmentation randomly selects 15% of tokens in the input sequence, with (1) 80% of the selected tokens replaced with a special [MASK] token, (2) 10% replaced with random vocabulary tokens, (3) 10% unchanged. This augmentation encourages robust embeddings to noisy or partial inputs.

Inspired by prior works in contrastive learning and code search (Wang et al., 2023; Li et al., 2022; 2021), we use a momentum decoder and maintain four queues of negative examples: (1) unmasked queries, (2) masked queries, (3) unmasked code, (4) masked code.

275 Let *D* denote the main decoder that is updated via gradient, 276 and *D'* denote the momentum decoder updated via expoen-277 tial moving average of the main decoder *D*. For any input 278 sequence *I*, we denote its normalized embedding from the 279 main decoder as D(I), and from the momentum decoder as 280 D'(I). We define four similarities:

282 
$$Sim(c_i, q_j) = D(c_i)^{\top} D'(q_j)$$
 (1)

283  
284 
$$\operatorname{Sim}(c_i, q_j^*) = D(c_i)^\top D'(q_j^*)$$
 (2)

$$\operatorname{Sim}(q_i, c_j) = D(q_i)^{\top} D'(c_j) \tag{3}$$

$$Sim(q_i, c_j^*) = D(q_i)^\top D'(c_j^*)$$
 (4)

For each informal query and formal Lean code within a
 batch of samples, we compute the temperature-scaled soft max similarity distributions:

$$P^{c \to q}(c_i, q_j) = \frac{\exp\left(\operatorname{Sim}(c_i, q_j)/\tau\right)}{\sum_{m=1}^{M} \exp\left(\operatorname{Sim}(c_i, q_m)/\tau\right)} \quad (5)$$

296 297

299 300

301

302

314

315 316

318

319

320

281

285

286

287

$$P^{c \to q^*}(c_i, q_j) = \frac{\exp\left(\text{Sim}(c_i, q_j^*) / \tau\right)}{\sum_{m=1}^{M} \exp\left(\text{Sim}(c_i, q_m^*) / \tau\right)} \quad (6)$$

$$P^{q \to c}(q_i, c_j) = \frac{\exp\left(\operatorname{Sim}(q_i, c_j)/\tau\right)}{\sum_{m=1}^{M} \exp\left(\operatorname{Sim}(q_i, c_m)/\tau\right)} \quad (7)$$

$$P^{q \to c^*}(q_i, c_j) = \frac{\exp\left(\text{Sim}(q_i, c_m^*) / \tau\right)}{\sum_{m=1}^{M} \exp\left(\text{Sim}(q_i, c_m^*) / \tau\right)} \quad (8)$$

where M = B + K, where B is the batch size and K is the queue size.  $\tau$  is the temperature hyperparameter.

To mitigate semantic ambiguity in the input and reduce the 306 impact of noisy supervision (e.g. weakly correlated query-307 code pairs), we employ a soft labeling strategy. Instead 308 of relying solely on hard one-hot labels, we generate soft 309 targets as a combination of a one-hot similarity vector and 310 pseudo-targets derived from the momentum decoder D'. 311 Specifically, we define the similarity between a query-code 312 pair using the momentum decoder as: 313

$$\operatorname{Sim}'(q_i, c_j) = D'(q_i)^{\top} D'(c_j) \tag{9}$$

The soft target vector  $Y_{soft}^{q \to c}$  is computed as:

$$Y_{soft}^{q \to c} = (1 - \beta) * Y_{onehot}^{q \to c} + \beta * S^{q \to c}$$
(10)

where  $\beta$  is a hyperparameter,  $Y_{onehot}^{q \to c}$  is a one-hot vector with a value of 1 at the index of the positive target and 0 elsewhere, and  $S^{q \to c}$  is a pseudo-target distribution over all Lean 4 statement code embeddings in the queue computed via softmax over momentum-based similarities:

326  
327  
328  
329  

$$S_{ij}^{q \to c} = \frac{\exp\left(\operatorname{Sim}'(q_i, c_j)/\tau\right)}{\sum_{m=1}^{M} \exp\left(\operatorname{Sim}'(q_i, c_m)/\tau\right)}.$$
(11)

Let *H* be the cross-entropy loss, and  $\mathcal{D}$  be the training set of query-code pairs, our overall loss function is defined as:

$$\begin{split} \mathcal{L} &= \frac{1}{4} \mathbb{E}_{(c,q)\sim\mathcal{D}} \Big[ \operatorname{H}(\mathbf{Y}_{\operatorname{soft}}{}^{c \to q}, \mathbf{p}^{c \to q}) + \operatorname{H}(\mathbf{Y}_{\operatorname{soft}}{}^{c \to q^*}, \mathbf{p}^{c \to q^*}) \\ &+ \operatorname{H}(\mathbf{Y}_{\operatorname{soft}}{}^{q \to c}, \mathbf{p}^{q \to c}) + \operatorname{H}(\mathbf{Y}_{\operatorname{soft}}{}^{q \to c^*}, \mathbf{p}^{q \to c^*}) \Big] \end{split}$$

## 4. Experiments

#### 4.1. LLM Settings

We utilize the OpenAI API to synthesize user queries from formal Lean statements as part of our dataset generation pipeline. Please see Appendix B for details on API configurations. For retrieval, we fine-tune DeepSeek-Prover-V1.5-RL 7B using LoRA and a query-code contrastive learning objective. Training settings are provided in Appendix C.

#### 4.2. Evaluation Settings and Data

To evaluate the performance of our Lean Finder, we compare it with LeanSearch (Gao et al., 2024a;b), a widely used search engine for Lean 4. To the best of our knowledge, LeanSearch is currently the only search engine that supports Lean retreival for Mathlib statements using either natural language queries or Lean 4 code.

We assess retrieval performance using the following metrics, both higher the better:

- **Recall**@K: This measures the proportion of relevant results within the top K retrieved results.
- Mean Reciprocal Rank (MRR): This calculates the average of the reciprocal ranks of the first relevant result for each query. It is defined as:  $MRR = \frac{1}{N} \sum_{i=1}^{N} \frac{1}{rank_i}$ , where N is the number of queries, and rank<sub>i</sub> is the rank position of the first relevant result for query *i*.

We evaluate our Lean Finder across four input modalities: informalized statement, synthetic user query, statement definition, and Lean code snippet.

**Informalized Statement.** To ensure a fair comparison, we impose two constraints: 1) The ground-truth formal statement *must* exist in the LeanSearch database (otherwise LeanSearch would always fail to retrieve it). 2) The corresponding informal query *must not* appear in either the LeanSearch database or our training data (otherwise LeanSearch would retrieve it perfectly, and our model would already have seen a near-duplicate). Therefore, we randomly select 150 Lean 4 statements from LeanSearch's database, and informalize them using a simplified prompt without additional context (dependent theorems, related theorems, etc.) used in our model. These newly informalized statements constitute a fair test set for comparison.

330 Synthetic User Query. For the synthetic user query
 331 modality, we generate 150 queries evenly distributed across
 332 five clusters, using the 30 of the newly informalized state 333 month of the memory for L and Sarrah commercian

ments as part of the prompt for LeanSearch comparison.

335 Statement Definition. To evaluate with the statement def-336 inition, we sample 100 Lean 4 statements from our dataset 337 and their corresponding definitions that are not in our train-338 ing set of Lean Finder. To simulate noisy queries and test 339 model robustness, we randomly replace 20% of the words 340 in the query with tokens sampled from the vocabulary. To 341 avoid invalid characters introduced by the retrieval model's 342 tokenizer during augmentation, we use a custom tokenizer 343 when augmenting the statement definitions.

**Lean Code Snippet.** Finally, for the code snippet modality, we sample 100 examples of code snippets from our dataset not present in our training set, and we do not apply augmentation.

All evaluation on LeanSearch is done manually using their web service<sup>3</sup>.

#### 4.3. Results on Synthetic Queries

345

346

347

348

349

350

351 352

353

354

355

356

357

358

359

360

361

363

367

369

370 371

373

374

375

383

384

Table 4: Retrieval by informalized statement.

Models	R@1	R@5	R@10	MRR
LeanSearch	0.487	0.700	0.773	0.599
Lean Finder (ours)	0.620	0.893	0.933	0.735
Table 5: Retr	ieval by	synthetic	user que	ry.
Models	R@1	R@5	R@10	MRR
LeanSearch	0.333	0.593	0.713	0.465
Lean Finder (ours)	0.407	0.693	0.780	0.537
Table 6: Retr	ieval by	statemen	t definitio	on.
Models	R@1	R@5	R@10	MRR
LeanSearch	0.620	0.840	0.870	0.716
Lean Finder (ours)	0.650	0.940	0.960	0.784
Table 7: Retr	rieval by	Lean co	de snippe	ts
Models	R@1	R@5	R@10	MRR
LeanSearch	0.090	0.260	0.330	0.182
Lean Finder (ours)	0.56	0.81	0.86	0.678

From results in Table 4, 5, 6, and 7, we observe that our Lean Finder consistently outperforms LeanSearch across all evaluated input modalities. Notably, for informalized statements, synthetic user queries, and Lean code snippets, our model achieves substantial performance gains. These improvements highlight the model's superior ability to capture user intent and generalize across diverse query formats,

<sup>3</sup>https://leansearch.net/

making it a more versatile and effective tool for supporting Lean 4 developers in real-world theorem proving workflows.

#### 4.4. Results on Real User Queries

**Table 8:** Retrieval by real user query. "Preferred Method": by which method the user-preferred answer is retrieved. "Tied": retrievals by LeanSearch (Gao et al., 2024a) and ours are comparable.

Preferred Method	Count	Percentage
LeanSearch	12	23.1%
Lean Finder (ours)	22	42.3%
Tied	18	34.6%

**Table 9:** User preference (Table 8) by intents. Cluster 1: Searchingfor existing code/lemmas. Cluster 3: Type-class, instance, axiom.Cluster 4: Proof engineering & everyday Lean usage.

Clusters	Ours Preferred	LeanSearch Preferred	Tied
Cluster 1	11	3	3
Cluster 3	1	1	0
Cluster 4	10	8	15

To further evaluate Lean Finder's effectiveness in real-world scenarios, we conduct a user study using Lean queries made by humans. We select 52 high-quality answerable user queries from GitHub and retrieve the top-3 statements from both LeanSearch and our Lean Finder for each query. The evaluation follows a style similar to LMArena (Chiang et al., 2024), where the order of the retrieval method is randomly shuffled for each query. As shown in Table 8, our Lean Finder was preferred in 42.3% of cases, substantially outperforming LeanSearch (23.1%), demonstrating the model's superior ability to capture user intent and provide satisfying answers to real user queries.

We further analyze this user preference by clustering these evaluated real user queries, and show user preference per cluster in Table 9. We find that our Lean Finder demonstrates stronger performance in Clusters 1 and 4.

#### 4.5. Lean Finder's User Interface

To facilitate Lean users, we also develop and deploy our user interface, as shown in Figure 4. We retrieve formal statements and also report their matching scores. We turn our web service offline to avoid disclosing identity information.

## 5. Related Works

**Code Search and Lean Retrieval.** Code search, aimed at retrieving the most semantically relevant code snippets from a codebase according to a specified natural language query, is a common activity that plays an important role in software development (Husain et al., 2019; Wang et al., 2023) and aligns closely with premise selection: both problems

# 385 Lean Finder – Retrieval Demo 386 Enter a Lean 4 query, set k, and click *Retrieve*.

387 3 3 Is there a lemma that confirms an additive or ng f(x) + f(y) = f(x + y) for elements x and 388 y in an additive group? 389 390 Rank Score Formal statement 0.4879 theorem \_root\_.IsAddGroupHom.add\_pull (x For any additive group homomorphism \$f: X (to Y\$ and any elements \$x, x' \in X\$, the map \$f\$ preserves addition, i.e., \$f(x + x') = f(x) + f(x')\$. This means that applying \$f\$ to the sum of two element individually. x' : X) : f (x + x') = f x + f x' := by rw[hf.map\_add] For any elements \$x, x'\$ in an additive group \$X\$, if \$f\$ is an additive group homomorphism, then \$f(x) + f(x') = f(x + x')\$. This means that applying \$f\$ to the sum of \$x\$ and \$x'\$ is the same as summing the images of \$x\$ and \$x'\$ under \$f\$. theorem \_root\_.IsAddGroupHom.add\_push (x x': x) : f x + f x' = f (x + x') := by rw[hf.map\_add] 395 0.4454 theorem \_root\_.IsAddGroupHom.app\_zero : f 0 = 0 := by For any additive group homomorphism \$f\$, we have \$f(0) = 0\$ sorry\_proof

Figure 4: Web user interface of our Lean Finder.

396 399 400 401 require retrieving a small set of code blocks that are help-402 403 ful for downstream tasks. Early methods relied on lexical heuristics (Alama et al., 2014; Urban, 2004), while recent 404 405 work leverages neural encoders (Irving et al., 2016; Wang & Deng, 2020; Mikuła et al., 2023; Yeh et al., 2023) and 406 dense retrieval models (Karpukhin et al., 2020; Qu et al., 407 2020) by learning semantic embeddings. In formal domains 408 like Lean, retrieval systems like LeanSearch (Gao et al., 409 2024a;b) assist users by allowing them to search for relevant 410 mathlib statements using natural language queries. These 411 systems primarily function by matching the natural language 412 translation of a formal statement to other formal statements 413 414 in mathlib. While effective in some settings, this approach supports only narrow input types and treats informal queries 415 as approximate paraphrases of formal code and focuses on 416 surface-level alignment, without capturing the diverse in-417 tents behind real-world user questions or the broader context 418 in which they arise. Lean Finder addresses this gap by in-419 troducing a user-centered, intent-aware retrieval framework 420 for Lean that supports a broad range of input modalities and 421 explicitly bridges the gap between the natural queries posed 422 by mathematicians and the formal statements stored in Lean 423 projects. By modeling real user intent and query diversity, 424 Lean Finder achieves consistently strong retrieval perfor-425 mance across all modalities and demonstrates substantial 426 gains over existing systems in a user study based on real 427 428 Lean queries. 429

430 AutoFormalization. Autoformalization, the automated 431 translation of informal mathematics into formal proof-432 assistant code, has become a key strategy for addressing the 433 scarcity of high-quality formal data. Early neural encoders 434 for statement formalization (Wang et al., 2018) have been 435 eclipsed by LLM-based pipelines that translate at scale with 436 few-shot prompting or fine-tuning (Wu et al., 2022; 2024; 437 Xin et al., 2024a;b). Systems such as DeepSeek-Prover and 438 InternLM2.5-StepProver formalize vast Lean corpora for 439

expert iteration (Xin et al., 2024a; Wu et al., 2024), while others generate synthetic formal proofs without informal seeds (Wu et al., 2020; Wang & Deng, 2020; Poesia et al., 2024). Projects such as MMA, Lean-STaR, TheoremLlama, and Lean Workbook create natural-language/formallanguage (NL-FL) pairs to mitigate data scarcity (Jiang et al., 2023; Lin et al., 2024; Wang et al., 2024; Ying et al., 2024); yet LLM brittleness still introduces subtle logical errors. Recent LLM tools embedded in user workflows (e.g., LeanDojo, LeanCopilot, LLM-Step) further streamline manual formalization by suggesting premises and proof tactics (Yang et al., 2023; Song et al., 2023; Welleck & Saha, 2023). Our work tackles the complementary challenge of autoformalization: synthesizing intent-rich, user-style queries from formal Lean statements to better serve practitioners. Unlike previous informalization approaches that simply translate formal statements into natural language, our pipeline explicitly models user intent and query phrasing. Injecting this intent-aware signal during fine-tuning yields significant retrieval improvements on real user queries, and bridges the gap between informal user needs and the formal knowledge encoded in Lean libraries.

## 6. Conclusion

We present Lean Finder, a significant advancement in the development of search tools for formal mathematics. By centering on real user intent and leveraging synthesized queries, semantic clustering, and fine-tuned embeddings, it bridges the gap between mathematicians' informal reasoning and formal theorem proving in Lean. Its superior performance over existing search tool demonstrate both practical utility and academic impact. Lean Finder not only enhances accessibility to mathlib4 but also lays the groundwork for future intelligent systems that assist in rigorous, collaborative mathematical discovery.

## 440 Impact Statement

"This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here."

#### References

441

442

443

444

445

446 447

448

449

450

451

452

477

481

482

483

- Alama, J., Heskes, T., Kühlwein, D., Tsivtsivadze, E., and Urban, J. Premise selection for mathematics by corpus analysis and kernel methods. *Journal of Automated Reasoning*, 52:191–213, 2014.
- 453 AlphaProof and teams, A. AI achieves silver-454 solving medal standard international math-455 ematical problems. olympiad https: 456 //deepmind.google/discover/blog/ ai-solves-imo-problems-at-silver-medal-let pp. 6769-6781, 2020. 457 458 2024. 459
- 460
  461
  461
  462
  463
  464
  464
  Asher, J. Leanexplore: A search engine for lean 4 declarations. https://github.com/justincasher/ lean-explore/blob/main/LeanExplore.
  464
- 465 Chiang, W.-L., Zheng, L., Sheng, Y., Angelopoulos, A. N.,
  466 Li, T., Li, D., Zhu, B., Zhang, H., Jordan, M., Gonzalez,
  467 J. E., et al. Chatbot arena: An open platform for evaluat468 ing llms by human preference. In *Forty-first International*469 *Conference on Machine Learning*, 2024.
- 471 De Moura, L., Kong, S., Avigad, J., Van Doorn, F., and von Raumer, J. The lean theorem prover (system description).
  473 In Automated Deduction-CADE-25: 25th International Conference on Automated Deduction, Berlin, Germany, August 1-7, 2015, Proceedings 25, pp. 378–388. Springer, 2015.
- Gao, G., Ju, H., Jiang, J., Qin, Z., and Dong, B. A
  semantic search engine for mathlib4. *arXiv preprint arXiv:2403.13310*, 2024a.
  - Gao, G., Wang, Y., Jiang, J., Gao, Q., Qin, Z., Xu, T., and Dong, B. Herald: A natural language annotated lean 4 dataset. *arXiv preprint arXiv:2410.10878*, 2024b.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang,
  S., Wang, L., and Chen, W. Lora: Low-rank adaptation of
  large language models, 2021. URL https://arxiv.
  org/abs/2106.09685.
- Husain, H., Wu, H.-H., Gazit, T., Allamanis, M., and Brockschmidt, M. Codesearchnet challenge: Evaluating the state of semantic code search. *arXiv preprint arXiv:1909.09436*, 2019.

- Irving, G., Szegedy, C., Alemi, A. A., Eén, N., Chollet, F., and Urban, J. DeepMath—deep sequence models for premise selection. 2016.
- Jiang, A. Q., Li, W., and Jamnik, M. Multilingual mathematical autoformalization. arXiv preprint arXiv:2311.03755, 2023.
- Jiang, A. Q., Li, W., and Jamnik, M. Multi-language diversity benefits autoformalization, 2024.
- Ju, H. and Dong, B. Mirb: Mathematical information retrieval benchmark. arXiv preprint arXiv:2505.15585, 2025.
- Karpukhin, V., Oguz, B., Min, S., Lewis, P. S., Wu, L., Edunov, S., Chen, D., and Yih, W.-t. Dense passage retrieval for open-domain question answering. In *EMNLP* or (1), pp. 6769–6781, 2020.
- Li, C., Wang, Z., Bai, Y., Duan, Y., Gao, Y., Hao, P., and Wen, Z. Formalization of algorithms for optimization with block structures, 2025a. URL https://arxiv. org/abs/2503.18806.
- Li, C., Xu, S., Sun, C., Zhou, L., and Wen, Z. Formalization of optimality conditions for smooth constrained optimization problems, 2025b. URL https://arxiv.org/ abs/2503.18821.
- Li, J., Selvaraju, R. R., Gotmare, A. D., Joty, S., Xiong, C., and Hoi, S. Align before fuse: Vision and language representation learning with momentum distillation, 2021. URL https://arxiv.org/abs/2107.07651.
- Li, J., Li, D., Xiong, C., and Hoi, S. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation, 2022. URL https:// arxiv.org/abs/2201.12086.
- Lin, H., Sun, Z., Yang, Y., and Welleck, S. Lean-star: Learning to interleave thinking and proving, 2024. URL https://arxiv.org/abs/2407.10040.
- Lin, Y., Tang, S., Lyu, B., Wu, J., Lin, H., Yang, K., Li, J., Xia, M., Chen, D., Arora, S., et al. Goedel-prover: A frontier model for open-source automated theorem proving. *arXiv preprint arXiv:2502.07640*, 2025.
- Lu, J., Wan, Y., Liu, Z., Huang, Y., Xiong, J., Liu, C., Shen, J., Jin, H., Zhang, J., Wang, H., Yang, Z., Tang, J., and Guo, Z. Process-driven autoformalization in Lean 4. *arXiv preprint arXiv:2406.01940*, 2024.
- Manthe, S. A formalization of borel determinacy in lean, 2025. URL https://arxiv.org/abs/2502.03432.

- Mikuła, M., Antoniak, S., Tworkowski, S., Jiang, A. Q.,
  Zhou, J. P., Szegedy, C., Kuciński, Ł., Miłoś, P., and
  Wu, Y. Magnushammer: A transformer-based approach
  to premise selection. *arXiv preprint arXiv:2303.04488*,
  2023.
- 500
  501 Monnerjahn, L. Formalisation of constructable numbers.
  502 2024.
- Monticone, P. Formalising fermat's last theorem for exponent 3 in lean. 2024.
  - Moura, L. d. and Ullrich, S. The Lean 4 theorem prover and programming language. 2021.

507

508

529

530

531

532

533

534

535

536

537

538

542

543

- Poesia, G., Broman, D., Haber, N., and Goodman, N. D.
  Learning formal mathematics from intrinsic motivation. *arXiv preprint arXiv:2407.00695*, 2024.
- 513 Qu, Y., Ding, Y., Liu, J., Liu, K., Ren, R., Zhao, W. X., Dong,
  514 D., Wu, H., and Wang, H. Rocketqa: An optimized train515 ing approach to dense passage retrieval for open-domain
  516 question answering. *arXiv preprint arXiv:2010.08191*,
  517 2020.
- Ren, Z., Shao, Z., Song, J., Xin, H., Wang, H., Zhao, W.,
  Zhang, L., Fu, Z., Zhu, Q., Yang, D., et al. Deepseekprover-v2: Advancing formal mathematical reasoning via
  reinforcement learning for subgoal decomposition. *arXiv preprint arXiv:2504.21801*, 2025.
- Shen, Z., Huang, N., Yang, F., Wang, Y., Gao, G., Xu,
  T., Jiang, J., He, W., Yang, P., Sun, M., et al. Real-prover: Retrieval augmented lean prover for mathematical reasoning. *arXiv preprint arXiv:2505.20613*, 2025.
  - Song, P., Yang, K., and Anandkumar, A. Towards large language models as copilots for theorem proving in lean. In *The 3rd Workshop on Mathematical Reasoning and AI at NeurIPS'23*, 2023. URL https://openreview. net/forum?id=C9X5sXa2k1.
  - Tao, Y., Liu, H., Wang, S., and Xu, H. Assisting mathematical formalization with a learning-based premise retriever. *arXiv preprint arXiv:2501.13959*, 2025.
- 539 Urban, J. MPTP—motivation, implementation, first exper540 iments. *Journal of Automated Reasoning*, 33:319–339,
  541 2004.
  - Wang, M. and Deng, J. Learning to prove theorems by learning to generate theorems. 2020.
- Wang, Q., Kaliszyk, C., and Urban, J. First experiments
  with neural translation of informal to formal mathematics. In *Conferences on Intelligent Computer Mathematics* (*CICM*), 2018.

- Wang, R., Zhang, J., Jia, Y., Pan, R., Diao, S., Pi, R., and Zhang, T. Theoremllama: Transforming general-purpose llms into lean4 experts. *arXiv preprint arXiv:2407.03203*, 2024.
- Wang, W., Wei, F., Dong, L., Bao, H., Yang, N., and Zhou, M. Minilm: Deep self-attention distillation for taskagnostic compression of pre-trained transformers, 2020. URL https://arxiv.org/abs/2002.10957.
- Wang, Y., Le, H., Gotmare, A. D., Bui, N. D., Li, J., and Hoi, S. C. Codet5+: Open code large language models for code understanding and generation. *arXiv preprint arXiv:2305.07922*, 2023.
- Welleck, S. and Saha, R. Ilmstep: LLM proofstep suggestions in lean. In *The 3rd Workshop on Mathematical Reasoning and AI at NeurIPS'23*, 2023. URL https: //openreview.net/forum?id=ODOJuAM4Qj.
- Wu, Y., Jiang, A. Q., Ba, J., and Grosse, R. Int: An inequality benchmark for evaluating generalization in theorem proving. arXiv preprint arXiv:2007.02924, 2020.
- Wu, Y., Jiang, A. Q., Li, W., Rabe, M., Staats, C., Jamnik, M., and Szegedy, C. Autoformalization with large language models. 2022.
- Wu, Z., Huang, S., Zhou, Z., Ying, H., Wang, J., Lin, D., and Chen, K. InternLM2.5-StepProver: Advancing automated theorem proving via expert iteration on large-scale lean problems. arXiv preprint arXiv:2410.15700, 2024.
- Xin, H., Guo, D., Shao, Z., Ren, Z., Zhu, Q., Liu, B., Ruan, C., Li, W., and Liang, X. Deepseek-prover: Advancing theorem proving in llms through large-scale synthetic data. arXiv preprint arXiv:2405.14333, 2024a.
- Xin, H., Ren, Z., Song, J., Shao, Z., Zhao, W., Wang, H., Liu, B., Zhang, L., Lu, X., Du, Q., et al. Deepseek-prover-v1.5: Harnessing proof assistant feedback for reinforcement learning and monte-carlo tree search. arXiv preprint arXiv:2408.08152, 2024b.
- Yang, K., Swope, A., Gu, A., Chalamala, R., Song, P., Yu, S., Godil, S., Prenger, R. J., and Anandkumar, A. Leandojo: Theorem proving with retrieval-augmented language models. *Advances in Neural Information Processing Systems*, 36:21573–21612, 2023.
- Yeh, E., Hitaj, B., Owre, S., Quemener, M., and Shankar, N. CoProver: A recommender system for proof construction. *arXiv preprint arXiv:2304.10486*, 2023.
- Ying, H., Wu, Z., Geng, Y., Wang, J., Lin, D., and Chen, K. Lean workbook: A large-scale lean problem set formalized from natural language math problems, 2024. URL https://arxiv.org/abs/2406.03847.

50 51	Zhu, T., Clune, J., Avigad, J., Jiang, A. Q., and Welleck, S. Premise selection for a lean hammer. <i>arXiv preprint</i> arXiv:2506.07477, 2025
12 53	<i>arxiv:2300.07477, 2023</i> .
54	Zulip. How does a noob learn to write tactics.
55	https://leanprover-community.github.
56	io/archive/stream/113489-new-members/
57	topic/How.20does.20a.20noob.20learn.
58	20to.20write.20tactics.3F.html,2020a.
i9	Zulip. Naming conventions for definitions. https:
51	//leanprover-community.github.io/
52	archive/stream/113488-general/topic/
53	Noun.2Fadjective.2Fverb.20naming.
54	20conventions.20for.20defs.20in.
5	20math11b.htm1, 2020b.
6	Zulip. Improving documentation. https:
7	//leanprover-community.github.io/
)	archive/stream/113488-general/topic/
	Improving.20documentation.html,2021a.
	Zulin Library search failures https:
	//leapprover community github ic/
	<pre>&gt;rehive/stream/113/88_general/tenic/</pre>
	library search 20failures html 2021h
	$\frac{1101a1y}{20210}$

# A. Details for Dataset Construction

## A.1. Data Source Details

 Table 10: Lean sources included in our dataset beyond mathlib4.

_	Project name	Category	URL
_	Mathlib4	Domain-specific library	https://github.com/leanprover-community/Mathlib4
	SciLean	Domain-specific library	https://github.com/lecopivo/SciLean
	A formalization of Borel determinacy in Lean (Manthe, 2025)	Research-linked repository	https://github.com/sven-manthe/A-formalization-of-Borel- determinacy-in-Lean
-	Optlib (Li et al., 2025a;b)	Research-linked repository	https://github.com/optsuite/optlib
-	Formalising Fermat's Last Theorem for Exponent 3 (Monticone, 2024)	Research-linked repository	https://github.com/riccardobrasca/flt3
	Formalisation of constructable numbers (Monnerjahn, 2024)	Research-linked repository	https://github.com/Louis-Le-Grand/Formalisation-of- constructable-numbers
-	Plausible	Transitive dependency	https://github.com/leanprover-community/plausible
_	ProofWidgets4	Transitive dependency	https://github.com/leanprover-community/ProofWidgets4
_	Aesop	Transitive dependency	https://github.com/leanprover-community/aesop
_	Batteries	Transitive dependency	https://github.com/leanprover-community/batteries

## A.2. Informalization Details

To support the inclusion of dependent statements during informalization, we construct a directed acyclic graph in which each node represents a Lean 4 statement and edges capture dependency relationships between statements. We ensure all dependencies of a given statement are informalized before the statement itself.

We use DeepSeekProver to retrieve the most relevant statements from Herald, and use the top-matching informal statement along with its corresponding formal statement as the related statement in the prompt for informalization.

# **B. OpenAI API Configurations**

We prompt OpenAI models during data collection and generation, and we describe the detailed configurations below.

**Informalization.** We leverage the  $GPT-4\circ$  API to convert formal Lean 4 statements into informal statements (Section 3.4). The temperature is set to 0.2 and the maximum output token limit is 1000. Inspired by Herald (Gao et al., 2024b), the prompt used for informalization is shown in Figure 5.

**User Query Filtering.** We employ GPT-40 to filter out user queries that cannot be addressed by any Lean 4 formal statement (Section 3.1.1). The temperature is set to 0.0 with a maximum of 500 output tokens. The prompted used for query filtering is shown in Figure 6.

**User Intent Clustering.** To cluster user intents (Section 3.1.1), we use the o3 API. The process involves two stages:

- **Bootstrap Categorization.** An initial set of representative user queries is clustered using a temperature of 1.0 and a token limit of 20000. The prompt used for initial cluster generation is shown in Figure 7.
- **Progressive Clustering.** Remaining queries are progressively categorized with the same temperature and token settings. The prompt used for this part is shown in Figure 8.

**User Query Generation.** To synthesize diverse user queries (Section 3.1.2), we use the GPT-4.1 mini API with a temperature of 0.7 and a maximum output of 1000 tokens. The prompt used is shown in Figure 9.

## 660 C. Training Details

We fine-tune our Lean Finder based on the DeepSeek-Prover-V1.5-RL 7B (Xin et al., 2024b), using a query–code contrastive learning objective. We apply LoRA (Hu et al., 2021) with a rank r = 32 and  $\alpha = 64$ . We train Lean Finder for 2 epochs with a per-GPU batch size of 10, gradient accumulation steps of 6, and 4 NVIDIA A100 GPUs. We use the AdamW optimizer with a learning rate of  $2 \times 10^{-5}$  and  $\epsilon = 1 \times 10^{-8}$ . The learning rate is linearly warmed up over the first 1,210 steps, then decayed to 0 using a cosine schedule. The momentum decoder is updated using a momentum coefficient of 0.99, and the temperature  $\tau$  for scaling the similarity distributions is set to 0.7. We set the momentum queue size to K = 640. For soft labeling, the weight  $\beta$  begins at 0 and is linearly increased to 0.4 over the course of training. Due to imbalanced input modalities in our training data (Table 2), we adopt a weighted sampler assigning higher sampling probabilities to examples from the underrepresented modality.

## **D.** Privacy and Data Release

To protect user privacy, we do not release any real user discussion data collected from Lean Zulip. All training of Lean Finder involving Zulip-derived queries are based solely on synthetic user queries that are reverse-engineered from formal Lean statements. We will only release these synthetic queries, along with queries from other input modalities. No identifiable information or original messages from Zulip users are included in our datasets.

## **E.** Prompts

/	You are an expert mathematician and an expert in Lean and Mathlib.
	**Instruction**:
	Your task is to first translate the formal theorem provided below into an informal statement that is more accessible to mathematicians and written in LaTeX. There are
	parts of information as inputs:
	1. Formal name: The formal name of the theorem you need to translate.
	2. Formal statement: The formal statement of the theorem you need to translate.
	3. Docstrings: A list of natural language comments written by humans in the theorem
	4. Neighbor statement: The statement located in the same file and closest in position to the formal theorem you need to translate.
	5. Dependent theorems: These are JSON list of dependent theorems at which the theorem you need to translate uses, and their informal statements and name.
	6. Related statement: This is a JSUN object for the statement from matinible that has related meaning as the current theorem at which you need to translate. The related statement is a statement to matinible that has related meaning as the current theorem at which you need to translate. The related statement is shown as the statement is shown as the statement.
-	statement formal statement, informal statements, and informal name are provided.
	Then create an informal name. Use the provided formal name of the statement according to the naming conventions. Utilize the informal statement written in the first
	Make sure you follow the principles of informal naming when creating informal names. Principles of informal statements should also be followed as much as possible
	**Principles of Informal Statements**:
	The informal statement should be written using human-used mathematical notations and formulas in LaTeX as much as possible, while also explaining the meaning
1	symbols involved.
	Explain more detailed mathematical setup only if the definition appearing in the statement is not commonly accepted. Both the inputs and values of the definition sho
	expressed using mathematical formulas as much as possible.
	*Example*-
	DO NOT use "Real loa":
	Use \$\log\$ instead.
	**Principles of Informal Names**:
	Emphasize the core concepts in the definition. The definition name should not merely list concepts; use words that indicate logical relationships and clearly state the
	conclusion.
	*Example*:
	Use "A equals B" or "A implies B" (or simply "\$A = B\$" and "\$A \to B\$"), instead of "theorem of A and B" when the theorem states the result of A = B or A → B. Both t'
	inputs and values of the definition should be expressed using mathematical formulas as much as possible.
	**Input**:
	"Formal name": <>
	*Formal Statement*: <>
	"Docstring": <>
	Neighbor statement. S
	**Output**:
	Return *only* a JSON object with keys
	"informal name"
	"informal statement"
1	
`	

Figure 5: Prompt used for Lean statement informalization.

)	
7	
/	
/	You are an expert at Lean 4 and an experienced user of Lean Zulip Chat.
	**Instruction**:
	Your task is to first read the excerpt provided below that contains the first 5 messages from a discussion thread in Lean Zulia Chat and then decide to accord on the excerpt based on the following criteria:
	2 unp Characteria and men decide to accept or reject the excerpt based on the following chiefia.
	a Lean 4 statement
	2. You should accept the excerpt if showing a pre-existing Lean 4 statement will move the discussion forward in a
	meaningful way.
	3. You should reject the excerpt if none of the above is true for this excerpt.
	4. When uncertain, use your best judgment.
	Then identify the main question in this discussion thread based on the excerpt provided below. The main question is the
	one that the user who initiated the discussion wants an answer for.
	CEPT:
	• The user is asking "Is there a lemma/theorem/definition that 2
	• The user needs a property that already exists as a lemma (e.g. "`map, prod` preserves multiplication")
	• The user doesn't directly ask for a statement, but providing an example statement to the user can address part of the
	problem.
	REJECT:
	<ul> <li>Library infrastructure, version bumps, build tools, CI, style.</li> </ul>
	{input_biock}
	**Outout**-
	Return *only* a JSON object with keys
	"decision"
	"main question"
	Where the value for decision must be *exactly* one word ACCEPT or REJECT
	<b>Figure 6:</b> Prompt used user query filtering
	rigure of frompt used user query mering.



#### Lean Finder: Semantic Search for Mathlib That Understands User Intents

825	
826	
827	
828	
820	
027	
021	
831	
832	
833	
834	
835	
836	
837	
838	You are an expert at Lean 4 and an experienced user of Lean Zulip Chat. **Instructions**:
839	You will be provided below with 25 excerpts that contain the first few messages from discussion threads in Lean Zulip Chat and the main user question that is being asked in that discussion threads and the main user question that is being asked in that the main user question can be assured at least particular back and the main user question threads as that the main user question can be assured at least particular back as the the main user question threads as the the main user question
840	statement will move the discussion forward in a meaningful way. You will also be provided with the current clusters of queries based on previous excerpts. Each cluster represents a
841	aisunct type of query based on the underlying intent or perspective of the main question in the excerpt. You need to do the following: *Task A Fitting new excerpts* :
842	1. For each input excerpt, try to fit the excerpt into an existing cluster based on the cluster's description and representative examples provided for queries in that cluster. 2. If new excerpts fit into an existing cluster of queries and the new excerpts reveal nuances that the current description of the cluster does not capture, you should append at most one
843	new sentence to an existing cluster description.
844	add at most one new example to an existing cluster examples.
845	<ol> <li>DO NOT remove any description and examples in the existing cluster. You are only allowed to add new ones when needed.</li> <li>If the excerpt DOES NOT fit into an existing cluster, do *Task B*, otherwise skip *Task B*</li> </ol>
846	*Task B – New cluster proposal (if needed)*
847	1. If the input excerpt doesn't fit into any existing cluster, then you should propose one new cluster. 2. The new cluster must have a descriptive name and a detailed description of what the cluster represents. The description should include the property of the query in this new cluster
848	and the description of underlying intent or perspective of the query made by the user. This must be detailed enough, such that a person reading this description can create new queries for a label and a detailed enough.
849	3. The new cluster must include at most 5 representative example queries that best define this cluster. The chosen example queries must be the main_question provided in the input for
850	the excerpt. Each included example must illustrate different situations within the cluster, so that a reader can generalise and invent many other queries that would still belong here.
851	**Important**: DO NOT remove existing clusters, even if they are not used by the new excercts.
852	DO NOT remove existing description and examples from existing clusters, You are only allowed add new clusters or add new description/examples to existing clusters.
853	The input will be a list of excerpts and a list of clusters.
854	1. channel: This contains the name of the channel that the excerpt is collected from in the Lean Zulip Chat.
Q55	<ol> <li>topic: This contains the name of the discussion thread the excerpt is collected from.</li> <li>main_question: This contains the main user question that is being asked in the excerpt.</li> </ol>
055	4. messages: This contains a list of initial messages from the discussion thread. Each message in the list contains the full name of the sender and the content of the message.
050	Each cluster will have the following contents:
050	2. cluster_description: the detailed description of what this cluster represents
828	3. examples: representative examples in this cluster
859	**Inputs**: List of excerpts:
860	{list_of_excerpts}
861	{list_of_clusters}
862	**Output**:
863	Return only a JSON object with a key "clusters", and the value for the key is a list of cluster objects. Each cluster object must have keys: "cluster_name", "cluster_description", and "examples".
864	
865	
866	
867	Figure 8: Prompt used for progressive clustering of user queries.
868	
869	
870	
871	
872	
873	

881	
882	
883	
884	
885	
005	
000	
887	
888	
889	
890	
891	
892	
893	You are an expert Lean 4 user and an experienced participant in the Lean Zulip chat. You know how people
894	/ usually phrase questions when they are searching for existing Lean statements.
895	**Instructions**:
896	You will be given:
807	1. The name, description, and 10 canonical examples of a specific cluster of queries from Lean Zulip or
808	GitHub. The queries in this cluster share the same underlying intent or perspective, precisely captured by the
070	description.
099	2. A single Lean 4 formal statement, formal name, informal name, and the natural language description of the
900	statement.
901	Write one realistic user query that:
902	1. Clearly belongs in this cluster the intent and perspective of the query should mirror the description and
903	examples.
904	2. Could be answered at least partially by showing the Lean 4 statement you were given, or would be moved
905	forward in a meaningful way by that statement.
906	3. Does not explicitly reveal the statement.
907	4. Feels like a genuine message a mathematician or Lean 4 user would post (natural tone, sensible level of
908	detail).
909	5. Stands alone – no references to "the cluster" or to these instructions.
910	**
911	Input . Cluster
012	
012	{cluster}
915	Statement:
914	Statement
915	(statement)
916	**Output**·
917	Return only a JSON object with a single key "generated guery". The value for the key is the generated guery.
918	for the input statement that belongs to the input cluster.
919	
920	
921	
922	<b>Figure 9:</b> Prompt used for generating user queries
923	i gare > i l'ompt asca foi generating asci quertos.
924	
025	
026	
720 0 <b>07</b>	
921	
928	
929	