# LeanCat: A Benchmark Suite for Formal Category Theory in Lean (Part I: 1-Categories)

**Rongge Xu**[1*]  **Hui Dai**[2]  **Yiming Fu**[3]  **Jiedong Jiang**[4]  **Tianjiao Nie**[5]
**Hongwei Wang**[6]  **Junkai Wang**[1]  **Holiverse Yang**[7]  **Jiatong Yang**[5]  **Zhi-Hao Zhang**[8]

[1]Yau Mathematical Sciences Center, Tsinghua University    [2]Iluvatar CoreX    [3]Department of Mathematics, Southern University of Science and Technology    [4]Westlake Institute for Advanced Study, Westlake University    [5] Qiuzhen College, Tsinghua University    [6] School of Mathematics and Physics, Xi'an Jiaotong-Liverpool University    [7]Department of Physics, The Chinese University of Hong Kong    [8]Yanqi Lake Beijing Institute of Mathematical Sciences and Applications (BIMSA)

## Abstract

Large language models (LLMs) have made rapid progress in formal theorem proving, yet current benchmarks under-measure the kind of abstraction and library-mediated reasoning that organizes modern mathematics. In parallel with FATE's emphasis on frontier algebra, we introduce **LeanCat**[1], a Lean benchmark for category-theoretic formalization—a unifying language for mathematical structure and a core layer of modern proof engineering—serving as a stress test of structural, interface-level reasoning. Part I: 1-Categories contains 100 fully formalized statement-level tasks, curated into topic families and three difficulty tiers via an LLM-assisted + human grading process. The best model solves $8.25\%$ of tasks at pass@1 ($32.50\%/4.17\%/0.00\%$ by Easy/Medium/High) and $12.00\%$ at pass@4 ($50.00\%/4.76\%/0.00\%$). We also evaluate LeanBridge which use LeanExplore to search Mathlib, and observe consistent gains over single-model baselines. LeanCat is intended as a compact, reusable checkpoint for tracking both AI and human progress toward reliable, research-level formalization in Lean.

## 1 Introduction

Recent advances in large language models (LLMs) and agentic training have revived the prospect of *end-to-end* formal theorem proving. On the formal side, systems such as OpenAI's early Lean-based prover (Polu et al., 2022) and DeepMind's AlphaProof (Hubert et al., 2025) demonstrate that reinforcement learning with formal verification feedback can produce non-trivial Lean proofs. More recently, specialized provers such as Seed-Prover 1.5 (Chen et al., 2025) show that large-scale agentic RL and test-time scaling can push formal success rates on established benchmarks substantially higher. These results suggest that formal proof generation is no longer limited to toy domains, and that tight tool feedback loops (retrieve–generate–verify) can be a decisive ingredient.

Despite steady progress in neural theorem proving, current formal benchmarks still do not adequately exercise *library-grounded, abstraction-heavy* reasoning. Widely used datasets such as miniF2F (Zheng et al., 2022) and FIMO Liu et al. (2023) largely draw from olympiad-style problems, while university-level suites like ProofNet Azerbayev et al. (2022) and PutnamBench Tsoukalas et al. (2024) focus on undergraduate competition or textbook material. These benchmarks are valuable, but they often reward short clever tricks or computation rather than sustained work inside rich abstract frameworks. By contrast, modern research mathematics is written at high generality, organized around reusable interfaces, and deeply intertwined with large libraries of definitions and lemmas—so success hinges less on a single key insight and more on navigating abstraction, managing definitions, and coherently composing library knowledge over long horizons.

---

[*]Email: `rongge@tsinghua.edu.cn`
[1]LeanCat (Part I: 1-Categories) is open sourced at `https://github.com/sciencraft/LeanCat`.

Category theory provides a natural stress-test for this capability: as the interface language of modern mathematics—categories, functors, natural transformations, adjunctions, limits/colimits, monads—formal proofs typically rely on *diagrammatic* and *universal-property* reasoning, i.e., constructing morphisms with the right naturality/uniqueness guarantees and showing commutativity across families of structures. Yet existing formal benchmarks rarely target this abstraction layer directly. To bridge this gap, we introduce **LeanCat**, a benchmark of 100 category-theory problems formalized in Lean 4 (mathlib), designed to test whether automated provers can operate *inside* a mature library and compose high-level abstractions rather than solve isolated puzzles. LeanCat complements algebra-focused benchmarks such as the FATE series Jiang et al. (2025) by shifting the frontier from abstract algebra to category theory. Our baseline evaluation reveals a stark abstraction gap: across five strong models, the best achieves only $8.25\%$ at pass@1 and $12\%$ at pass@4, with accuracy dropping sharply from Easy to High once library navigation and long-horizon abstraction management become necessary (Figure 1). We further observe a persistent gap between producing a plausible natural-language argument and generating compilable Lean code, highlighting a pronounced *natural-to-formal bottleneck* (Figure 2).
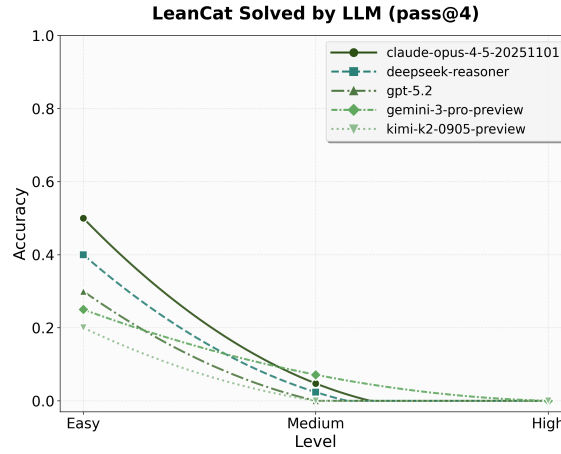


Figure 1: Formal proof accuracy (pass@4) across LeanCat model baselines.
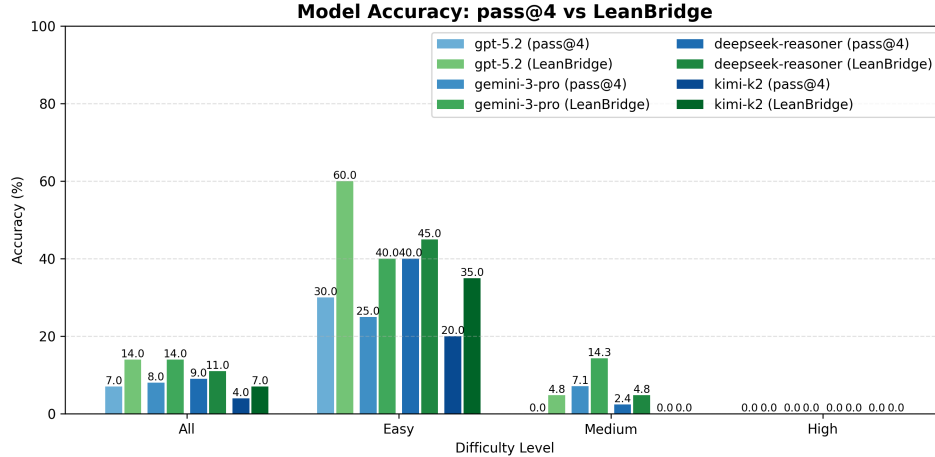


Figure 2: Intermediate natural-language (NL) vs. formal-language (FL) (pass@4) accuracy on LeanCat, highlighting the natural-to-formal gap.

To our knowledge, LeanCat is the first installment of a *category-theory benchmark series*. This paper focuses on on 1-category theory. We envision future extensions targeting richer structures such as monoidal categories, enriched and braided settings, and eventually 2-categories and higher-categorical interfaces. Beyond benchmarking, we view this direction as meaningful for both (i) *human*

*mathematics*, by mapping what parts of abstract library-level reasoning remain hard to formalize and where libraries need strengthening, and (ii) *AI*, by forcing progress on abstraction-aware planning, retrieval of relevant lemmas, and robust compilation-driven refinement.

Our key contributions are summarized as follows:

- **LeanCat benchmark (1-Categories).** We present LeanCat, a set of 100 category-theory problems formalized in Lean 4 (mathlib 4.19.0). The tasks span eight topic clusters (from basic categorical properties to monads), curated to cover reusable abstract interfaces rather than contest-style tricks.

- **Difficulty annotation pipeline.** We propose a grading process combining model-based estimates and expert judgment. Each problem receives multiple 1–10 ratings (from advanced LLM attempts and from human formalizers), aggregated with heavier weight on human scores to assign Easy/Medium/High labels (20/42/38).

- **Baseline evaluation.** We benchmark state-of-the-art provers on LeanCat under uniform conditions. Our evaluation (Section 3-4) includes ChatGPT-5.1 and ChatGPT-5.2 (OpenAI, 2025a;b), Claude 4.5 (Anthropic, 2025), Google's Gemini 3 pro (Gemini Team, Google, 2025), an advanced chain-of-thought reasoner DeepSeek -V3.2-Thinking and DeepSeek V3.2 Speciale (Liu et al., 2025) and ann agentic model Kimi K2 Team et al. (2025). At pass@1, the best model solves **8.25%** of LeanCat; at pass@4, the best reaches **12%**. We provide breakdowns by difficulty and identify dominant failure modes (library gaps, abstraction mismatch, and stalled multi-step plans).

- **Search-augmented proving via LeanBridge.** We evaluate a retrieve–analyze–generate–verify loop that integrates mathlib retrieval (via LeanExplore) and compiler feedback, illustrating how tool-augmented workflows can improve robustness on a subset of problems.

## 2 LEANCAT BENCHMARK DESIGN

### 2.1 BENCHMARK STRUCTURE AND CONTENT

**Data Source**  The problems in our benchmark are organized into two main parts: Abstract and Concrete:

- **Abstract part**: Problems are primarily drawn from standard, widely used textbooks in category theory, especially Category Theory in Context (Riehl, 2017) and Categories for the Working Mathematician (Mac Lane, 1998), together with a small number of problems adapted from unpublished lecture notes (Kong; Zheng).

- **Concrete part**: Problems are curated mainly from Abstract and Concrete Categories (Adámek et al., 1990), which provides a systematic source of exercises on concretizability, injectivity, and related themes.

- **Others**: Beyond these core sources, we also include selected problems inspired by research papers and advanced community-driven references (Chen, 2021; Adámek et al., 2021).

Each LeanCat problem is self-contained at the statement level: the formal statement of the theorem is provided (often with informal description as in the Problem listing above), and all necessary definitions exist in Lean's environment (either already in Mathlib (mathlib Community, 2020) or introduced as part of the problem setup). Where possible, we drew on known theorems from category theory literature, many tasks were crafted or adapted specifically to probe corner cases and interface interactions that an AI prover might struggle with. In several High-difficulty cases, relevant lemmas were not readily available, forcing human formalizers to devise intermediate results. This aspect of LeanCat, that it sometimes goes beyond the library, makes it a particularly stringent test for automated provers, which cannot rely solely on rote application of existing library facts.

LeanCat comprises 100 theorem statements in category theory, each fully formalized in Lean 4 (i.e. each problem is given as a Lean theorem declaration, with required definitions and context available). The problems are organized into eight topical clusters reflecting key areas of category theory:

- **Basic Category Properties (Problems 1-18)**: Fundamental results about categories and morphisms. These include properties of monomorphisms and epimorphisms, initial/terminal objects, factorization of idempotents, and examples of categorical constructions.
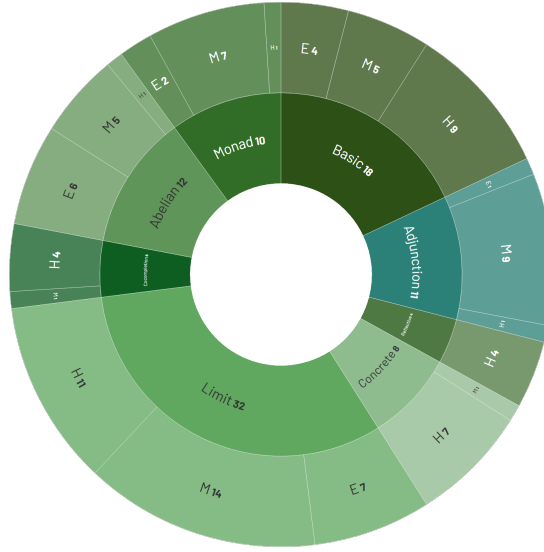
Figure 3: Sunburst diagram showing the distribution of our LeanBridge problem sets by topic and difficulty. The inner ring groups problems into thematic sections (Basic, Adjunction, Reflective, Concrete, Limit, Cocompletion, Abelian, Monad), while the outer ring lists individual problems, labelled by difficulty: E = Easy, M = Medium, H = High.

- **Adjunctions (Problems 19-29)**: Constructions and criteria involving adjoint functors, a central concept in category theory. Problems include proving that familiar functors have left or right adjoints, as well as general adjointness criteria like the comma category condition (Problem 19) andsome concrete examples (Problem 28). These tasks test a prover's ability to manipulate universal properties and to move between element-wise ("pointwise") reasoning and diagrammatic reasoning.

- **Reflective and Coreflective Subcategories (Problems 30-33)**: Abstract properties and concrete examples of a special kind of subcategories (e.g. classifying the reflective subcategories of $\mathcal{S}et$ and $\mathcal{T}\mathrm{op}^{\mathcal{C}H}$).

- **Concrete Categories (Problems 34-41)**: Categories with a faithful forgetful functor to Set and related notions. These tasks are largely concrete and overlap substantially with other areas of mathematics, such as topology, order theory, and set theory. They are designed to test a model's ability to connect abstract concepts with concrete examples.

- **Limits and Colimits (Problems 42-73)**: This is the largest cluster, covering a range of results about limits, colimits, and related categorical constructs. Many of these statements are at the frontier of what Lean's Mathlib currently contains, indeed some (like Problem 46 or 67) required developing new formal definitions. This cluster stresses a prover's ability to chain together multiple categorical facts.

- **Cocompletions (Problems 74-78)**: This part is based on recent work on cocompletions. It requires the LLM to introduce new definitions and then prove the key theorems built on those definitions that are not currently available in Mathlib.

- **Abelian Categories (Problems 79-90)**: Tasks concerning abelian categories and homological algebra concepts. Abelian categories are highly structured categories (every morphism has a kernel and cokernel, etc.) and generalize the category of modules or abelian groups. These statements mirror standard results in homological algebra, but formalizing them in Lean requires care with categorical abstractions (e.g. kernel objects, exact sequences) that are more complex than their set-based counterparts. Solving them may require the prover to introduce creative auxiliary lemmas about kernels, images, or exactness - a tall order for automated tools.

- **Monads (Problems 91-100)**: The final cluster centers on monads and their associated constructions (Kleisli and Eilenberg-Moore categories). Monads are a high-level concept that encapsulates a form of "computation" or structure; proving properties about them in Lean often demands a two-level reasoning (reasoning about the monad's algebraic laws and about category-theoretic conditions

4

like preservation of coequalizers). This cluster provides a valuable test of an AI's ability to work with highly abstract algebraic structure in a category-theoretic setting.

## 2.2 CURATION WORKFLOW

LeanCat is curated through a three-stage workflow that combines expert selection, LLM-assisted drafting, and rigorous human verification:

1. **Collection.** Three category-theory experts curate candidate problems from established sources (as described above), aiming to cover both core interfaces (e.g., adjunctions, limits/colimits, monads) and representative proof patterns (diagram chasing, universal properties, naturality).

2. **Formalization.** For each selected problem, we first use several LLMs to draft Lean 4 statements. The same three category-theory experts then review these drafts and retain only semantically correct formal statements. For problems where none of the models produced a correct statement, we organized a three-day workshop at Westlake University, bringing together Lean experts to author the missing statements and (when feasible) corresponding proofs.

3. **Review.** Finally, two independent reviewers with strong mathematical backgrounds and Lean expertise perform a full consistency pass, checking compilation, fixing definitional mismatches, and ensuring the formal statement matches the intended mathematical meaning.

**Statement-level tasks.**   LeanCat is a *statement-level* benchmark: each item consists of a single standalone theorem to prove, rather than a scaffolded sequence of intermediate lemmas that gradually builds toward a final goal. This design aims to evaluate general library-grounded proving capabilities—retrieval, definition management, and abstraction navigation—instead of rewarding problem-specific hint engineering.

**Scope and difficulty.**   Overall, LeanCat is *broad* in its topical coverage across category theory, and *deep* in the sense that even short-looking theorems can require layered abstractions and careful use of reusable interfaces, mirroring how mathematicians work inside large formal libraries.

**Formalization standards.**   All benchmark files follow strict, uniform conventions: (i) each Lean file contains exactly one `sorry` after the final theorem; (ii) the natural-language problem description (in LaTeX) is included as a comment immediately preceding the formal statement; (iii) universe levels are explicitly fixed to avoid ambiguity and instability that are common in category-theory developments.

## 2.3 DIFFICULTY ANNOTATION PIPELINE

Rather than relying purely on the problem author's intuition, we implemented a systematic LLM+human rating pipeline to grade problem difficulty on a 10-point scale, then binned the scores into the three categories: Easy, Medium, and High. This approach aimed to capture both human expertise and automated solver perspectives, similar in spirit to FATE's curation process (which combined expert judgment with model feedback for difficulty ranking).

Our pipeline proceeded as follows:

- LLM Difficulty Scoring: For each model, a correct proof contributes a proof score, and (only when that model does not already have a correct proof) a correct statement contributes a smaller statement score. The total score of a problem is the weighted sum of these contributions across models; the difficulty is then defined as Diff=max (0,10-PF score-ST score), so that problems solved by none of the models have difficulty 10, while problems for which all proof columns are green have difficulty 0.

- Human Difficulty Scoring: In parallel, two human mathematicians (with Lean expertise and category theory background) independently rated each problem on the same 1-10 difficulty scale. They factored in things like the length of the proof, intricacy of arguments, and whether any non-obvious lemmas are required. The human scores tended to correlate with intuition: e.g. a trivial diagram chase might be 2/10, whereas a complex construction spanning several definitions could be 9/10.

- Aggregation: We combined the scores giving 50% weight to human ratings and 50% to LLM ratings. Finally, we mapped the numeric scores to difficulty categories. We defined thresholds informed by the distribution of scores: roughly, Easy for scores $\leq 6$, High for scores $\geq 8.5$, and Medium for the rest. These cut-offs neatly split the set into 20 Easy, 42 Medium, and 38 High, see Table 4.

| Sections | Order | Statement | | | | Proof | | | | | | | Difficulty | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Models | | ChatGPT 5.1 | Gemini 3 pro | Deepseek | Kimi k2 0905 | ChatGPT 5.1 | ChatGPT 5.2 | Gemini 3 pro | Deepseek-reasoner | Deepseek-v3.2-speciale | Kimi k2 0905 | Grok 4.1 | LLM | Human | Average |
| Weight | | 0.4 | 0.4 | 0.4 | 0.4 | 2 | | 2 | 2 | | 2 | 2 | 0.5 | 0.5 | |
| Basic | 1 | | | | | | | | | | | | 1.6 | 3 | 2.3 | Easy |
| | 2 | | | | | | | | | | | | 8 | 3 | 5.5 | Easy |
| | 3 | | | | | | | | | | | | 9.6 | 2 | 5.8 | Easy |
| | 4 | | | | | | | | | | | | 8 | 4 | 6 | Easy |
| | 5 | | | | | | | | | | | | 7.2 | 5 | 6.1 | Medium |
| | 6 | | | | | | | | | | | | 10 | 3 | 6.5 | Medium |
| | 7 | | | | | | | | | | | | 9.2 | 5 | 7.1 | Medium |
| | 8 | | | | | | | | | | | | 10 | 5 | 7.5 | Medium |
| | 9 | | | | | | | | | | | | 8.8 | 7 | 7.9 | Medium |
| | 10 | | | | | | | | | | | | 10 | 7 | 8.5 | High |
| | 11 | | | | | | | | | | | | 10 | 7 | 8.5 | High |
| | 12 | | | | | | | | | | | | 10 | 7 | 8.5 | High |
| | 13 | | | | | | | | | | | | 10 | 8 | 9 | High |
| | 14 | | | | | | | | | | | | 10 | 8 | 9 | High |
| | 15 | | | | | | | | | | | | 10 | 8 | 9 | High |
| | 16 | | | | | | | | | | | | 10 | 8 | 9 | High |
| | 17 | | | | | | | | | | | | 10 | 9 | 9.5 | High |
| | 18 | | | | | | | | | | | | 10 | 9 | 9.5 | High |
| | 19 | | | | | | | | | | | | 5.2 | 5 | 5.1 | Easy |

Figure 4: We use LeanBridge to benchmark LeanCat across several LLMs. The table records, for each problem and model on natural-language statements, whether the model produced a correct statement or proof. A green cell indicates that the model produced a correct Lean proof of the intended statement. A blue cell indicates that the model produced a correct statement but without a correct proof. An olive cell indicates that the model's statement is not phrased in the same natural language as the original problem. A grey cell indicates that the output is not a complete Lean script (for example, it contains syntax or type errors).

This joint annotation procedure yields richer insight than a single-expert classification. It effectively uses large models as "second-opinion graders." The resulting difficulty labels have already proved useful in analysis: for instance, all of the seven problems solved by the best model (Section 4) came from the Easy set, and problems with score $\geq 9$ (the "hardest of High") saw zero successes across all models - quantitative evidence that our difficulty rankings align with actual solvability.

## 3 EXPERIMENTS AND RESULTS

### 3.1 EVALUATION PROTOCOL

We evaluate prover performance on LeanCat using a standardized *pass@k* protocol, adapted from prior work in code generation and automated theorem proving . Concretely, for each model–problem pair, we sample up to $k$ independent proof attempts under the same prompting and tool settings, and count the problem as solved if *any* attempt compiles and passes verification. We report both pass@1 and pass@4: pass@1 reflects single-try reliability, while pass@4 captures the benefit of limited sampling and iterative diversity. Unless otherwise noted, all evaluations are run under identical

conditions (same model settings, context budget, and verification pipeline) to ensure comparability across models.

**Environment and Input**: Each LeanCat problem is provided to the model in a uniform format: we give the statement in Lean (including the precise theorem name, hypotheses, and conclusion) along with context such as imported libraries and definitions. Models thus see the formal goal as it would appear to a human using Lean. No informal hints or decomposed lemmas are given - the model must find the proof from the statement and standard library knowledge. This setup mimics a realistic scenario where an AI prover is tasked with solving a new theorem given the definitions.

**Automated Proof Generation**: Language-model-as-prover: for API-based LLMs (GPT-5.2, Claude, Gemini), we prompt the model to generate a Lean proof script directly, we adopt the same prompt template as in FATE-Eval Jiang et al. (2025) (shown in Listing 1) for consistency across evaluations. The model produces a proof term or tactic script, which we then feed into Lean for verification.

```
"""
    You are an expert in Lean 4 and Mathematics. Please finish the
following proof in Lean4 code.
    Do not change the original statement. Copy the final statement
to prove exactly.
    Please include the complete header (including imports and
namespaces) so that your code can
    pass the Lean4 compiler. Please solve the statement step by
step and provide your complete
    Lean4 code between ```lean4 and ``` after careful reasoning.
    The statement for you to complete is:

```lean4
    {formal_statement}
```
"""
```

Listing 1: Prompt template for proof generation

**Verification**: A proof attempt is considered successful if the Lean theorem prover accepts it as a valid proof of the given statement. We instrument Lean to check the model outputs automatically. If the model output is incomplete or incorrect (does not typecheck), that attempt counts as a failure. Models do not get to "see" the result of the verification; in pass@k evaluation, each attempt is separate.

**Pass@k Calculation**: We compute pass@1 as the fraction of problems for which the model produced a correct proof in a single attempt. Pass@4 is the fraction of problems where at least one out of 4 attempts succeeded. With 100 problems, these percentages translate directly to number of problems solved. We note that LeanCat's problems are all of roughly equal weight (each is a standalone theorem), so simple pass rates are a meaningful indicator of overall capability. We also track pass@1 within each difficulty category (Easy/Med/High) to see how performance degrades as difficulty increases.

We used a uniform evaluation setup: each attempt had an output budget of 50,000 tokens and a Lean verification time limit of 5 minutes, and all models ran against the same Lean environment (Lean 4 + Mathlib 4.19.0) to ensure consistency. If a model exceeded the token budget or failed to finish verifying within the time limit, the attempt was counted as a failure. In practice, however, these resource limits were rarely the deciding factor: most attempts either found a proof quickly (typically within 30 seconds, except for reasoning model like DeepSeek) or got stuck almost immediately (often after only a few dozen tokens).

We emphasize that pass@4 is not meant as a realistic usage scenario (where one would run a model four times per theorem); rather, it provides an optimistic upper bound on success if we could perfectly pick from a handful of model tries. In an ideal case where attempts are independent, pass@4 could be significantly higher than pass@1, but as we will see, the increase was modest for LeanCat. This suggests that when a model fails on a problem in one attempt, trying again usually yields a similar outcome unless guided by a different strategy.

Table 1: LeanCat FL Solved by LLM (Pass@1/Pass@4)

| Model | Easy | Medium | High | All |
|---|---|---|---|---|
| claude-opus-4-5 | 32.50% / 50.00% | 4.17% / 4.76% | 0.00% / 0.00% | 8.25% / 12.00% |
| gpt-5.2 | 27.50% / 30.00% | 0.00% / 0.00% | 0.00% / 0.00% | 5.50% / 7.00% |
| gemini-3-pro | 11.25% / 25.00% | 2.38% / 7.14% | 0.00% / 0.00% | 3.25% / 8.00% |
| deepseek-reasoner | 18.75% / 40.00% | 0.60% / 2.38% | 0.00% / 0.00% | 4.00% / 9.00% |
| kimi-k2-0905 | 10.00% / 20.00% | 0.00% / 0.00% | 0.00% / 0.00% | 2.00% / 4.00% |

Preliminary data indicates only a 1-2 problem increase from pass@1 to pass@4 for the top models, reinforcing the difficulty of LeanCat tasks.

**LeanBridge:** LeanBridge implements a retrieve–analyze–generate–verify loop that augments an LLM with mathlib retrieval and verifier feedback. Given a problem, we first use the natural-language statement as a query to LeanExplore to retrieve relevant mathlib entities (e.g., definitions, lemmas). The retrieved snippets are then provided to the model as context for analysis and Lean code generation. Each generated proof script is checked in a clean Lean environment; a candidate is accepted only if it typechecks and contains no `sorry` or `admit`. To guard against superficial "green" proofs that do not match the intended meaning, all accepted candidates are additionally reviewed by a human expert for semantic correctness with respect to the original problem statement. If verification fails, LeanBridge parses the compiler error messages to decide whether additional retrieval is needed; it then augments the context with both the retrieved information and the verifier feedback, and prompts the model to revise the proof. Unless otherwise stated, we cap the loop at a maximum of 4 iterations for both (i) natural-language-to-statement formalization and (ii) natural-language theorem proving.

## 3.2 BASELINE RESULTS AND ANALYSIS

We evaluated five state-of-the-art models on LeanCat under the protocol described above. We summarize the main findings below.

- **Overall success remains low.** On pass@1 (first attempt), the best model (Claude Opus 4.5) solved $8.25\%$ problems; GPT-5.2 solved $5.5\%$, DeepSeek Reasoner solved $4\%$, while Gemini 3 Pro at $3.25\%$ and Kimi at $2\%$. Across all models, only **10 distinct** problems were solved by at least one model on the first attempt, leaving $91.75\%$ unsolved at pass@1. Allowing up to four attempts per problem improves results but does not change the overall picture: pass@4 rises to $12\%$ for Claude Opus 4.5, $9\%$ for DeepSeek Reasoner and $8\%$ for Gemini 3 Pro, $7\%$ for GPT-5.2, and $4\%$ for Kimi. In total, **14 distinct** problems are solved by at least one model at pass@4.

- **Clear Easy–Medium–High gap.** Performance decreases monotonically with our difficulty annotations. For example, Claude Opus 4.5 achieves $32.5\%$ on Easy, $4.17\%$ on Medium, and $0\%$ on High at pass@1 (and $50\%$, $4.76\%$, $0\%$ at pass@4). GPT-5.2 shows a similar trend ($27.5\%$, $0\%$, $0\%$ at pass@1). Even within the Easy subset, the absolute success rates remain far from saturation, suggesting that LeanCat's "baseline difficulty" is already beyond what current models reliably handle once nontrivial abstraction and library navigation are required.

- **Case study (typical success): Problem 82.** This problem is consistently solved by models that effectively translate the categorical notion of simplicity into concrete linear algebra. The successful approach recognizes that a simple object in $\mathrm{Vect}_k$ must be one-dimensional, then constructs an explicit isomorphism with k by leveraging a nonzero vector and the simplicity condition. The proof elegantly bridges abstract category theory with elementary vector space properties, demonstrating a clear understanding of how structural definitions manifest in concrete categories.

- **Retries help selectively, indicating high variance and brittle search.** Moving from pass@1 to pass@4 yields only modest absolute gains for the strongest model (+5 for Claude Opus 4.5), but it substantially benefits some weaker models (e.g., DeepSeek Reasoner improves from 3 to 8). This pattern is consistent with high-variance behavior: many problems

are either solved quickly or not approached effectively at all, and additional attempts mainly help when a model happens to sample a viable strategy or recall the right library lemma.

- **Error analysis: library knowledge gaps dominate, followed by abstraction mismatch and incomplete plans.** Manual inspection of failed runs suggests three recurring failure modes: (i) *library knowledge gaps*: models often fail to recall the correct Mathlib definitions/lemmas or their usable forms, leading to dead ends or hallucinated lemma names; (ii) *abstraction mismatch*: some attempts switch to element-wise reasoning when the intended proof is categorical/structural, which is typically unproductive in Lean without substantial setup; (iii) *incomplete multi-step planning*: models may discharge a few local goals but then stall, unable to connect intermediate facts into a coherent end-to-end proof. Pure syntax-level errors occur, but are less frequent than these semantic/strategic failures.

Overall, these baselines confirm that LeanCat is substantially harder for current LLM-based provers than earlier Lean benchmarks. Even with multiple attempts, success remains sparse on Medium/High problems, pointing to a need for improved library retrieval, better abstraction-aware proof planning, and more reliable strategy exploration.

## 4 DISCUSSION AND FUTURE WORK

**LeanCat as a benchmark (and a series).** LeanCat is intended as a reusable checkpoint for LLM-based theorem proving in *abstract* mathematics. This paper introduces **LeanCat-1** (1-category theory); we view it as the first installment of a broader *LeanCat series*, with planned extensions targeting richer categorical interfaces such as monoidal categories and higher-categorical structures (e.g., bicategories / strict 2-categories) that are already represented in Mathlib's ecosystem.

**Library integration.** All LeanCat statements are formalized in Lean 4; as solutions are found, they can be merged into Mathlib, creating a feedback loop: benchmark → solutions → stronger libraries/solvers → harder remaining frontier.

**What LeanCat stresses.** Our results highlight three persistent bottlenecks for current automated provers: (i) *library awareness* (finding and applying the right Mathlib lemmas), (ii) *abstraction control* (staying at the correct categorical level rather than drifting to element-wise reasoning), and (iii) *long-horizon coherence* (maintaining a consistent plan across many dependent steps).

**Future work and broader impact.** On the benchmark side, we will expand LeanCat beyond 1-categories to include additional clusters and multi-theorem tasks, and to progressively cover higher abstraction layers—e.g., a monoidal track and a 2-categorical track (where monoidal categories can be viewed through the lens of a single-object bicategory), enabling finer-grained diagnosis of where provers fail as abstraction increases. On the solver side, promising directions include stronger retrieval over Mathlib, hierarchical decomposition into auxiliary lemmas, and multi-agent pipelines (planner/verifier/lemma-suggester). For *human mathematics*, we expect LeanCat-style checkpoints to help prioritize missing interfaces and reusable lemmas in the library; for *AI*, they provide a concrete target for improving abstraction-aware planning and library-grounded reasoning. Finally, porting LeanCat to other proof assistants (e.g., Coq or Isabelle) would enable cross-system comparisons and encourage transfer of proof-engineering techniques.

### REFERENCES

Jiří Adámek, Horst Herrlich, and George Strecker. *Abstract and concrete categories*. Wiley-Interscience, 1990.

Jiří Adámek, Liang-Ting Chen, Stefan Milius, and Henning Urbat. Reiterman's theorem on finite algebras for a monad. *ACM Transactions on Computational Logic (TOCL)*, 22(4):1–48, 2021.

Anthropic. System card: Claude opus 4.5. Anthropic, November 2025. URL `https://www.anthropic.com/claude-opus-4-5-system-card`. Released November 24, 2025. Accessed on January 1, 2026.

Zhangir Azerbayev, Bartosz Piotrowski, and Jeremy Avigad. ProofNet: A benchmark for autoformalizing and formally proving undergraduate-level mathematics problems. In *Second MATH-AI Workshop*, 2022.

Jiangjie Chen, Wenxiang Chen, Jiacheng Du, Jinyi Hu, Zhicheng Jiang, Allan Jie, Xiaoran Jin, Xing Jin, Chenggang Li, Wenlei Shi, et al. Seed-prover 1.5: Mastering undergraduate-level theorem proving via learning from experience. *arXiv preprint arXiv:2512.17260*, 2025.

Ruiyuan Chen. On sifted colimits in the presence of pullbacks. *arXiv preprint arXiv:2109.12708*, 2021.

Gemini Team, Google. Gemini 3 pro model card. November 2025. URL `https://storage.googleapis.com/deepmind-media/Model-Cards/Gemini-3-Pro-Model-Card.pdf`.

Thomas Hubert, Rishi Mehta, Laurent Sartran, Miklós Z Horváth, Goran Žužić, Eric Wieser, Aja Huang, Julian Schrittwieser, Yannick Schroecker, Hussain Masoom, et al. Olympiad-level formal mathematical reasoning with reinforcement learning. *Nature*, pp. 1–3, 2025.

Jiedong Jiang, Wanyi He, Yuefeng Wang, Guoxiong Gao, Yongle Hu, Jingting Wang, Nailing Guan, Peihao Wu, Chunbo Dai, Liang Xiao, et al. Fate: A formal benchmark series for frontier algebra of multiple difficulty levels. *arXiv preprint arXiv:2511.02872*, 2025.

Liang Kong. Category theory. Unpublic.

Aixin Liu, Aoxue Mei, Bangcai Lin, Bing Xue, Bingxuan Wang, Bingzheng Xu, Bochao Wu, Bowei Zhang, Chaofan Lin, Chen Dong, et al. Deepseek-v3. 2: Pushing the frontier of open large language models. *arXiv preprint arXiv:2512.02556*, 2025.

Chengwu Liu, Jianhao Shen, Huajian Xin, Zhengying Liu, Ye Yuan, Haiming Wang, Wei Ju, Chuanyang Zheng, Yichun Yin, Lin Li, et al. FIMO: A challenge formal dataset for automated theorem proving. *arXiv preprint arXiv:2309.04295*, 2023.

Saunders Mac Lane. *Categories for the working mathematician*, volume 5. Springer Science & Business Media, 1998.

The mathlib Community. The Lean mathematical library. In *Proceedings of the 9th ACM SIGPLAN International Conference on Certified Programs and Proofs*, 2020.

OpenAI. Gpt-5.1 instant and gpt-5.1 thinking system card addendum. OpenAI, November 2025a. URL `https://openai.com/index/gpt-5-system-card-addendum-gpt-5-1/`. Released November 12, 2025. Accessed on January 1, 2026.

OpenAI. Update to gpt-5 system card: Gpt-5.2. OpenAI, December 2025b. URL `https://openai.com/index/gpt-5-system-card-update-gpt-5-2/`. Released December 11, 2025. Accessed on January 1, 2026.

Stanislas Polu, Jesse Michael Han, Kunhao Zheng, Mantas Baksys, Igor Babuschkin, and Ilya Sutskever. Formal mathematics statement curriculum learning. *arXiv preprint arXiv:2202.01344*, 2022.

Emily Riehl. *Category theory in context*. Courier Dover Publications, 2017.

Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, et al. Kimi k2: Open agentic intelligence. *arXiv preprint arXiv:2507.20534*, 2025.

George Tsoukalas, Jasper Lee, John Jennings, Jimmy Xin, Michelle Ding, Michael Jennings, Amitayush Thakur, and Swarat Chaudhuri. PutnamBench: Evaluating neural theorem-provers on the putnam mathematical competition. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 11545–11569. Curran Associates, Inc., 2024. URL `https://proceedings.neurip s.cc/paper_files/paper/2024/file/1582eaf9e0cf349e1e5a6ee453100aa 1-Paper-Datasets_and_Benchmarks_Track.pdf`.

Hao Zheng. Category note. Unpublic.

Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. MiniF2F: A cross-system benchmark for formal olympiad-level mathematics. In *The Tenth International Conference on Learning Representations*, 2022.