

# COMBI BENCH: BENCHMARKING LLM CAPABILITY FOR COMBINATORIAL MATHEMATICS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Neurosymbolic approaches integrating large language models with formal reasoning have recently achieved human-level performance on mathematics competition problems in algebra, geometry and number theory. In comparison, combinatorics remains a challenging domain, characterized by a lack of appropriate benchmarks and theorem libraries. To address this gap, we introduce CombiBench, a comprehensive benchmark comprising 100 combinatorial problems, each formalized in Lean 4 and paired with its corresponding informal statement. The problem set covers a wide spectrum of difficulty levels, ranging from middle school to IMO and university level, and span over ten combinatorial topics. CombiBench is suitable for testing IMO solving capabilities since it includes all IMO combinatorial problems since 2000 (except IMO 2004 P3 as its statement contain an images). Furthermore, we provide a comprehensive and standardized evaluation framework, dubbed *Fine-Eval* (for **F**ill-in-the-blank **i**n **L**ean **E**valuation), for formal mathematics. It accommodates not only proof-based problems but also, for the first time, the evaluation of fill-in-the-blank questions. Using *Fine-Eval* as the evaluation method and Kimina Lean Server as the backend, we benchmark several LLMs on CombiBench and observe that their capabilities for formally solving combinatorial problems remain limited. Among all models tested (none of which has been trained for this particular task), Kimina-Prover attains the best results, solving 7 problems (out of 100) under both “with solution” and “without solution” scenarios. We open source the benchmark dataset alongside with the code of the proposed evaluation method at <https://github.com/MoonshotAI/CombiBench/>.

## 1 INTRODUCTION

Automated theorem proving (ATP) has long been a prominent research area at the intersection of artificial intelligence and mathematics (Bibel, 2013); aiming to develop computer programs capable of automatically verifying or discovering mathematical proofs. In recent years, the integration of large language models (LLMs) with formal theorem provers such as Lean Moura and Ullrich (2021), Coq The Coq Development Team (2023), and Isabelle Wenzel et al. (2008) has received increased attention. These formal systems can rigorously verify step-by-step the proof generated by an LLM and provide feedback, thereby enhancing the models’ performance in automated theorem proving.

In this research area, the design of high-quality benchmark tests is essential, as they offer valuable guidance for improving model capabilities. Three well-known competition-based benchmarks in automated theorem proving include miniF2F Zheng et al. (2021), FIMO Liu et al. (2023) and PutnamBench (Tsoukalas et al., 2024). miniF2F primarily features high school-level mathematics problems and competition questions, FIMO consists of algebra and number theory problems from the International Mathematical Olympiad (IMO), and PutnamBench originates from the prestigious North American university-level Putnam competition. Next to these, ProofNet Azerbayev et al. (2023) is a benchmark for autoformalization and formal proving of undergraduate-level mathematics.

However, each of these benchmarks has notable limitations. For instance, miniF2F is limited to high school mathematics, FIMO and ProofNet are written in the now outdated formal proof language Lean 3 (de Moura et al., 2015), and PutnamBench includes only a small number of combinatorial problems. However, it’s worth mentioning that a Lean 4 version of ProofNet is provided in DeepSeek-

Prover-V1.5 (Xin et al., 2024b). We have analyzed the quantity and proportion of combinatorial problems in some existing benchmarks in Table 1, as well as their applicability to Lean 4.

Table 1: Benchmark comparison of combinatorics problems

TYPE	miniF2F	FIMO	PutnamBench	ProofNet	CombiBench
Count	0	0	29	0	<b>100</b>
Ratio	0	0	4.4%	0	<b>100%</b>
Lean 4 friendly	😞	😞	😞	😞	😄

The above formal mathematics benchmarks in Lean 4 scarcely cover combinatorial mathematics. But combinatorial mathematics is also the area with the weakest model capabilities in the field of automated theorem proving at present. In 2024, Google DeepMind’s AlphaProof successfully solved 4 problems (out of 6) in the International Mathematical Olympiad (IMO) using Lean 4, achieving a level comparable to human silver medalists. However, the remaining two unsolved problems by AlphaProof were combinatorial. This also confirms what we mentioned above. This outcome can be attributed to three primary factors:

- The absence of robust benchmarks to assess models’ proficiency in combinatorial mathematics;
- The limited coverage of combinatorial content within Lean’s math library—Mathlib (mathlib Community, 2020).
- In combinatorial mathematics, the gap between informal and formal is larger than in other areas. For a specific problem, it is often necessary to add some definitions that are unique to that problem.

CombiBench aims to address these gaps by providing testing standards for evaluating the capabilities of large language models in this area.

Besides, traditional formal math benchmarks (such as miniF2F, FIMO) primarily consist of theorem-proving tasks, where evaluation methods focus on assessing whether a model can complete a proof. However, PutnamBench introduces a new standardized formalization for fill-in-the-blank problems, which existing evaluation approaches cannot handle effectively. Unlike standard proof-based problems, fill-in-the-blank questions require the model to construct a solution and verify its correctness. Current evaluation methods are insufficient for this problem, focusing solely on proving given propositions without addressing solution generation and validation. To bridge this gap, we developed a novel evaluation framework, named FINE-EVAL, to assess model performance on fill-in-the-blank problems. This approach adapts evaluation techniques to new problem formats and provides a more comprehensive measure of a model’s mathematical reasoning abilities, particularly in solution construction and verification.

The contributions of this work are outlined as follows:

- We introduce CombiBench, a benchmark containing 100 combinatorial formal statements in Lean 4 and their corresponding informal descriptions. The benchmark spans problems of varying difficulty, from middle school to IMO and university level, covering over ten distinct combinatorial topics.
- We propose a standardized and comprehensive evaluation framework for formal mathematics, enabling the assessment of proof-based problems and fill-in-the-blank questions for the first time.

## 2 RELATED WORK

### 2.1 IMO

The International Mathematical Olympiad (IMO) is the most prestigious and competitive global mathematics competition for high school students. Established in 1959, the IMO aims to challenge

108 and inspire young mathematicians worldwide. Each year, teams of up to six students from over 100  
109 countries participate in the event, which consists of solving six highly complex mathematical problems  
110 over two days. The problems span various areas of mathematics, including algebra, geometry, number  
111 theory and combinatorics, requiring not only deep mathematical knowledge but also creativity and  
112 problem-solving skills. Many individuals who have won the Fields Medal—the highest honor in  
113 mathematics, have participated in the IMO. More recently, the annual IMO competition has also  
114 become widely recognized as the ultimate grand challenge for Artificial Intelligence (AI) (Challenge,  
115 2019). Artificial Intelligence Mathematical Olympiad (AIMO) has set up a \$10 million prize to  
116 reward the first AI system that reaches the level of an IMO gold medalist (Prize, 2023).

## 117 118 119 2.2 AUTOMATED THEOREM PROVING 120

121 Early research in automated theorem proving focused on first-order theorem provers designed for  
122 simple logical frameworks Schulz (2002) and theorem provers based on symbolic engines (Chou  
123 et al., 2000). With the advent of deep learning technologies, there has been a significant shift toward  
124 utilizing large language models (LLMs) to automate the theorem-proving process.

125 GPT-f Polu and Sutskever (2020) is an automated theorem prover and proof assistant developed  
126 for the Metamath formalization language. Its introduction marks the first instance that a formal  
127 mathematics community accepts and incorporates proofs generated by deep learning. The emergence  
128 of GPT-f has significantly advanced the field of automated theorem proving. Subsequently, the Draft,  
129 Sketch, and Prove (DSP) Jiang et al. (2022) was introduced, utilizing informal proofs to aid in the  
130 generation of formal ones. LEGO-Prover Wang et al. (2023a) builds a library through a modular  
131 formal proof, allowing large language models (LLMs) to retrieve existing skills and generate new  
132 ones during the proof process. DT-Solver Wang et al. (2023b) introduces a dynamic-tree Monte  
133 Carlo search algorithm. BFS-Prover Xin et al. (2025), using a best-first search approach, attained  
134 state-of-the-art performance among theorem provers based on search algorithms.

135 The aforementioned work employs a search-based approach to predict the next step in the proof.  
136 Recently, progress in this area of research has led to the development of an alternative approach,  
137 where the language model generates the entire proof directly. DeepSeek-Prover Xin et al. (2024a) and  
138 Goedel-Prover Lin et al. (2025) both utilize this method. In particular, DeepSeek-Prover incorporates  
139 both search-based and whole-proof approaches. Despite notable progress, these existing methods face  
140 significant challenges. While LLMs excel at pattern matching and sequence generation, effectively  
141 capturing the deep, structured, and often non-linear reasoning required for complex formal proofs  
142 remains difficult. Kimina-Prover Preview Wang et al. (2025) proposed a novel reasoning-driven  
143 exploration paradigm for formal theorem proving, using reinforcement learning to enhance whole-  
144 proof generation.

## 145 146 2.3 EXISTING BENCHMARKS 147

148 With the rapid advancement of automated theorem proving, several formal mathematical benchmarks  
149 in the Lean language have been introduced in recent years. miniF2F is a benchmark designed  
150 to evaluate automated theorem-proving systems in various formal systems. It includes various  
151 mathematical problems, such as exercises from prestigious math olympiads (AMC, AIME, IMO)  
152 and high school and undergraduate coursework. The main objective of miniF2F is to provide  
153 a standardized benchmark to directly evaluate and compare theorem-proving systems. Initially  
154 implemented in Lean and Metamath, it has since been extended to HOL Light and Isabelle, broadening  
155 its applicability across different formal proof environments. The current state-of-the-art result in the  
156 miniF2F test set is 80.74%, achieved by Kimina-Prover Preview (Wang et al., 2025).

157 For combinatorial problems, miniF2F and FIMO focus exclusively on algebra and number theory  
158 problems from the IMO. PutnamBench includes 26 combinatorial problems, but its coverage of  
159 combinatorial topics (e.g. countability, power sets, discrete structures, games) is limited. LeanComb  
160 Xiong et al. (2025) introduces a data enhancement approach and offers a benchmark specifically  
161 focused on combinatorial identities. The proofs of these identities often involve analytical techniques  
rather than purely combinatorial reasoning.

### 3 COMBIBENCH

CombiBench is the first benchmark focused on combinatorial competition problems written entirely in Lean 4. It is a manually produced benchmark that includes 100 combinatorial mathematics problems of varying difficulty and knowledge levels.

**Range of Topics.** For the selection of topics in combinatorics, we follow Brualdi’s classical combinatorics textbook “Introductory Combinatorics” (Brualdi, 2004). It consists of fourteen chapters and is widely used in undergraduate and graduate courses in combinatorics.

- What Is Combinatorics?
- Permutations and Combinations
- The Pigeonhole Principle
- Generating Permutations and Combinations
- The Binomial Coefficients
- The Inclusion-Exclusion Principle and Applications
- Recurrence Relations and Generating Functions
- Special Counting Sequences
- Systems of Distinct Representatives
- Combinatorial Designs
- Introduction to Graph Theory
- More on Graph Theory
- Digraphs and Networks
- Pólya Counting

The book systematically introduces the fundamental concepts, methods, and applications of combinatorics, covering many important topics such as permutations and combinations, the pigeonhole principle, generating functions, graph theory, and combinatorial design.

**Composition and Diversity.** CombiBench consists of 10 easy problems from <https://www.hackmath.net/>, 42 exercises from Brualdi’s book, 36 IMO problems, and 12 problems from other math competitions, which are shown in Table 2. This composition ensures that CombiBench covers a wide range of difficulty, from easy to difficult, reflecting good diversity in difficulty.

Table 2: Problem source distribution

Source	Hackmath	Brualdi’s book	IMO	APMO	Balticway
Conut	10	42	36	2	1
Source	EGMO	IMO-Shortlist	IZHO	BXMO	USAMO
Conut	1	4	2	1	1

For the IMO problems, we collected all combinatorics problems from the official IMO problems since 2000, totaling 37 problems. However, two problems contain images (Problem 3 from 2004 and Problem 5 from 2023), making their direct formalization difficult. After our analysis, we found that Problem 5 from 2023 could potentially be formalized. The figure is only for illustration and does not contain any additional information that is not already in the text. Then, we removed Problem 3 from 2004 to get 36 IMO problems.

For the exercises from the Brualdi book, we randomly sampled problems from 14 chapters, choosing three problems from each chapter. This ensured that the 42 problems were evenly distributed across all 14 chapters, helping to guarantee the diversity of the topics covered in our selection.

The complete proofs of Problem 3 and Problem 5 from 2024 have already been formalized in Mathlib. Therefore, we directly refer to the statements of these problems from Mathlib in CombiBench, along with the necessary definitions used in the statements.

**Naming Convention.** Depending on the source of the problems, we adopt different naming conventions. For the 10 easy problems, we name them `hackmath_1`, `hackmath_2`, etc. Exercises from Brualdi’s book are named using the format `brualdi_#chapter_#number`. Each competition problem is named as “`#competition_#year_p#number`”. For instance, “`imo_2019_p5`” refers to the IMO Problems 2009 Problem 5, and “`imo_2019_p5_1`” refers to the first sub-question of the IMO Problems 2009 Problem 5. And “`brualdi_ch8_6`” refers to Problem 6 from the exercises section of Chapter 8 in Brualdi’s book.

**Classification of Problem.** In CombiBench, 45% of the problems require first providing a solution to the problem and then proving its correctness, such as the following problem :

- Determine the number of permutations of  $\{1, 2, \dots, 8\}$  in which no even integer is in its natural position.

Therefore, these problems do not directly state a proposition and cannot be directly formalized. Previous benchmarks, such as miniF2F and FIMO, circumvented this issue by modifying the problem statement to require proving that a given solution satisfies the problem’s constraints. However, this modification reduces the overall difficulty of the problem, as a significant portion of the challenge may lie in coming up with the solution itself. PutnamBench introduces a standardized method for formalizing statements of such problems, which more accurately reflects the difficulty of informal problems. The following is an example where we formalize the above combinatorial problem in PutnamBench style:

```
abbrev brualdi_ch6_11_solution : ℕ := sorry

/--
Determine the number of permutations of {1,2,...,8} in which no even
integer is in its natural position.
-/
theorem brualdi_ch6_11
  (sols : Finset (Equiv.Perm (Finset.Icc 1 8)))
  (h_sols : ∀ σ, σ ∈ sols ↔ (∀ i, Even i.1 → σ i ≠ i)) :
  sols.card = brualdi_ch6_11_solution := by sorry
```

**Challenges.** Our formalization team consists of five doctoral students and one master’s student, each with over a year of experience in learning Lean, including a major contributor to Mathlib and a reviewer of Mathlib. During the formalization process, we found that besides simple problems, most problems require more than 30 minutes to formalize. For problems at the International Mathematical Olympiad (IMO) level, almost every problem takes over 3 hours to formalize, with some problems taking more than 8 hours. This indicates that formalizing combinatorial mathematics problems at the IMO level remains challenging and time-consuming.

One reason for this is the scarcity of theorems related to combinatorics in Mathlib, which makes it challenging to formalize corresponding problems and requires us to define many concepts. We analyzed the length of the formalized statements in CombiBench, miniF2F, PutnamBench, and FIMO. We excluded blank lines, comments, `import`’s, and `open`’s, focusing only on the length of the formalized code directly related to the problem statements. See Figure 1. We found that all problem formalizations in FIMO and miniF2F are within 15 lines. In PutnamBench and CombiBench, some problems have formalizations ranging from 16 to 30 lines, and those exceeding 30 lines are almost entirely from CombiBench. This indicates that formalizing combinatorics problems is particularly challenging. We also counted the number of lines of formalization code for all problems (excluding blank lines) in CombiBench. The results showed that more than half of the problems have more than 10 lines of code after formalization, more than a quarter have more than 20 lines, and the most challenging problem reached 67. We present the longest formalization of the statement of the problem in CombiBench in Appendix A.4. This data point indicates that formalizing combinatorial mathematics problems is still highly challenging.

**Data Quality.**

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

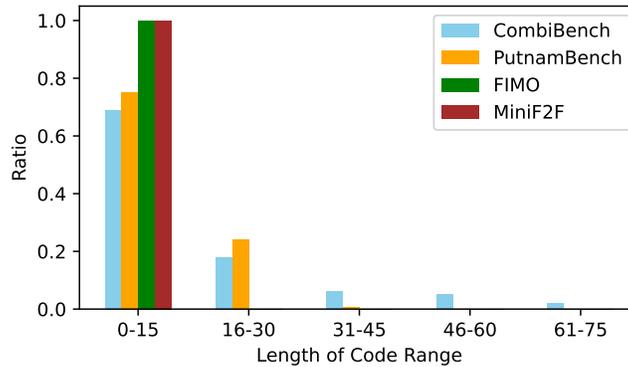


Figure 1: Code length distribution

Correctness is ensured by grammar and semantics. Lean can automatically check the grammatical correctness of our formal statements. Regarding semantic correctness, there’s currently no automated mechanism to verify the accuracy of the formalization. To ensure the quality of our dataset, we assign at least two people to double-check each completed formalization question. We hope that this approach ensures that most questions are correctly formulated. Meanwhile, we used Kimina-Prover-72B to perform Pass@128 sampling on the negations of each problem, and none of them were successfully proven, which in turn helps to ensure the correctness of our problem formalizations. Since human expert performance on the benchmark provides an important reference, we attempted to annotate problems from the benchmark. For example, for a human expert with over one year of Lean experience, producing a complete formal proof of an IMO combinatorics problem requires roughly 20 hours. Therefore, fully annotating the entire benchmark with human proofs would be prohibitively costly and practically infeasible. Nevertheless, we conducted a small-scale validation and successfully completed formal proofs for more than half of the problems. To prevent data leakage, we don’t include any ground truth proofs in our GitHub and huggingface links to prevent automated scripts from incorporating them into training data. And we require the community to provide us with complete proofs when releasing results on CombiBench. We will then upload them to the leaderboard after Lean verification and manual review. If they wish to show proofs on platforms like huggingface or GitHub, we urge them to encrypt them by packaging and to share any necessary information to facilitate reproducibility, such as model weights, API entries, technical reports, and code. In addition to the aforementioned measures, we continuously incorporate community feedback and carry out long-term maintenance and improvements on CombiBench to ensure that it consistently maintains a high level of data quality.

We also run experiments to model the relationship between problem difficulty and model discrimination. We measure the success rate on each problem in the theorem-proving task and analyze the results to demonstrate that CombiBench possesses strong discriminatory power in A.5.

**Licensing and Rules of Engagement.** CombiBench is available under MIT License for Lean 4. We host a public leaderboard at <https://moonshotai.github.io/CombiBench/leaderboard.html> and welcome evaluation results from future works. To prevent data contamination, we do not make the proofs in CombiBench public.

## 4 EVALUATION

We construct experiments to grasp CombiBench’s difficulties with state-of-the-art theorem-proving approaches and general-purpose LLMs. Previous evaluation methods include whole-proof generation and tree search. Whole-proof generation directly queries a language model to generate a complete proof, while tree-search methods gradually build the evidence by searching for possible proof strategies. Both methods work well for verifying proof problems, but cannot verify fill-in-the-blank problems. To address the issue, we propose a new evaluation pipeline to verify both types of problems in Lean 4 named FINE-EVAL (Fill-in-the-blank in Lean Evaluation). We formalize the matching of

the LLM prediction answer and the ground truth into a theorem-proving problem to rigorously verify the fill-in-the-blank question and avoid the complicated rules and format requirements as approaches in the math word problem.

#### 4.1 EVALUATION METHOD

FINE-EVAL interacts with two servers, LLM and Lean. We deployed Kimina Lean server Santos et al. (2025) on a 64-core 512G machine and called it in parallel with 10 workers. LLM takes as input an entire formal statement with the proof and the blanks to be filled in replaced by ‘sorry’s, and writes a snippet of complete Lean 4 code. Since LLMs sometimes cheat by commenting out all the code to pretend it passed compilation, we show an example in the appendix A.6. Therefore, we require the code, after comments are removed, to satisfy the following conditions:

- Does not contain “sorry”s;
- Can not define any new axioms and local instance (Prevent changing the meaning of the statement and generating deceptive proofs).
- Can be compiled by Lean without errors;
- Compared with the input formal statement, except for the replaced “sorry”s, the other parts are exactly matched.

A code that does not meet these requirements is considered a failure. Otherwise, we believe the proof is successful and check whether the solution provided by LLM and the ground truth match exactly. If the solution and ground truth are exactly matched, we believe that the model has successfully solved the problem. If not, we then try to verify that the answer that LLM predicts is equivalent to the ground truth. To avoid unnecessary LLM calls, we construct a formal statement as follows, asserting that “`xx_solution = ground_truth`” and try to prove it using two common tactics: ‘`rfl`’ and ‘`norm_num`’.

```
example : imo_2006_p2_solution = ground_truth := by
  try rfl
  try norm_num
```

If the lean verifier returns that it cannot be proved, we consider that the matching answer is non-trivial and add the statement to the input of the first round of LLM to prove it again. We use the same criterion as the first phase to determine whether the LLM’s proof is successful. Figure 2 demonstrates the process of FINE-EVAL. In the first stage, the model answers “10 / 20” and the corresponding proof at once, which can be verified by the lean 4 server. Next, we try to use “norm\_num” and “rfl” to automatically prove that 10 / 20 is equivalent to 1 / 2. When these attempts fail, we enter the second stage and let the model try to prove this problem, and the model fully demonstrates that the answer it predicts is equivalent to the ground truth. To prevent models from first submitting a trivial solution, then deceptively eliciting the ground truth to subsequently construct the actual proof, we impose a length limit on the second-stage proof. The number of characters in the output proof, after removing spaces and newline characters, must not exceed 42.

In practical applications, the two-stage verification method can often be complex. Consequently, we also provide a simplified evaluation approach, as illustrated in Figure 3. This simplified method initially follows the same procedure as the first stage of the Fine. However, if an exact match fails, we directly employ the `rfl` tactic to verify if the solution generated by the model is definitionally equivalent to the ground truth. This enables the evaluation of fill-in-the-blank problems using only a single Large Language Model (LLM) call. While this approach imposes stricter requirements on the model’s predicted answer, it is more easily embedded into other workflows.

**Metric** For fill-in-the-blank questions, we treat a complete two-stage FINE-EVAL as a single verification. Following PutnamBench, we also conducted a comparative experiment, replacing the corresponding sorry with ground truths and letting the model complete the proof. We report pass@N at different sample budgets as the performance metric.

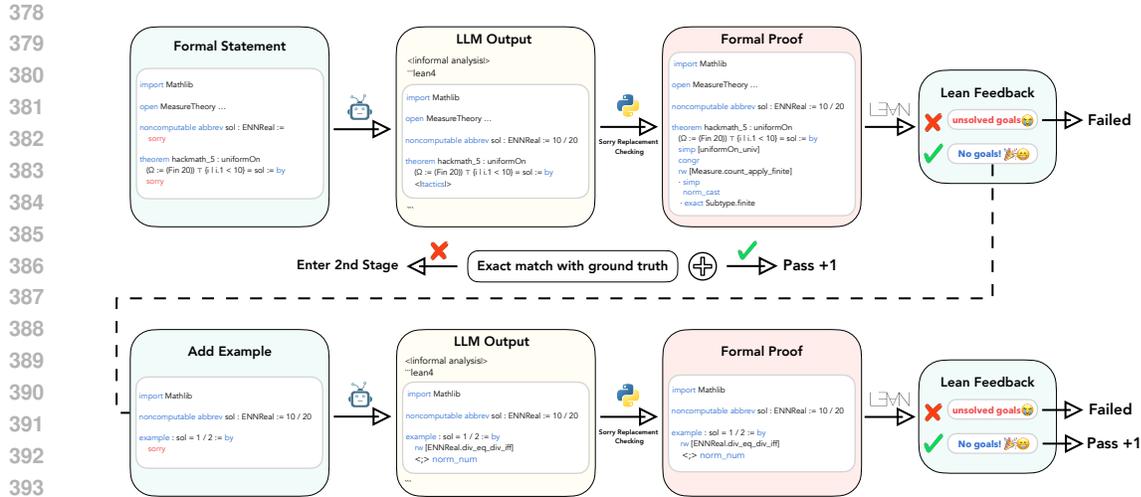


Figure 2: Pipeline of the two-stage FINE-EVAL.

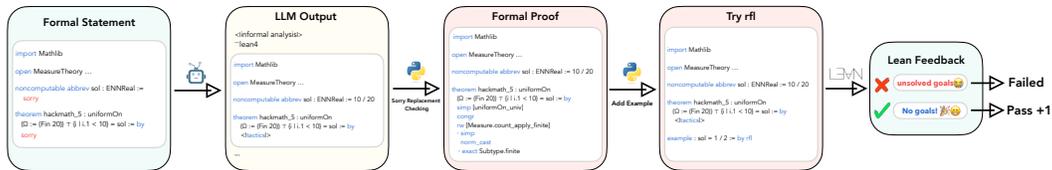


Figure 3: Pipeline of the one-stage FINE-EVAL.

## 4.2 BASELINE RESULTS

411 To demonstrate the significant challenge that CombiBench poses to LLMs, we evaluate it on different  
412 LLMs using various computational budgets. We consider two types of models, one is the model  
413 fine-tuned on the automated theorem proving task (indicated by "theorem prover"), the size of such  
414 models is usually 7B or less (except Kimina-Prover Preview) and the other is the general reasoning  
415 model (indicated by "reasoning model"). We evaluate pass@16 of these models on CombiBench  
416 using whole-proof generation. Considering the economic cost and reproducibility, we do not evaluate  
417 with a larger sample budget and tree-search method. We set temperature 1.0 during evaluation, and  
418 use a maximum tokens length equals 16,000. The system prompt and user prompt are shown in A.3.

419 The experimental results are presented in Table 3. We found that all models solved only a few  
420 problems, likely due to the absence of formal combinatorial question libraries and the inherent  
421 difficulty of such questions. None of the theorem-proving models with 7B parameters or fewer  
422 were able to solve all problems, which we attribute to the limited generalization capabilities of  
423 smaller models when facing out-of-distribution data. The general reasoning model answered a few  
424 questions correctly, and its performance showed minimal difference between the with\_solution and  
425 without\_solution settings, likely due to its strong informal mathematical reasoning ability. In the  
426 without\_solution case, both evaluation methods mentioned earlier produced identical results. This  
427 observation implies that, due to the model's current capability limitations, the solutions predicted are  
428 typically either completely correct or entirely wrong. As a result, sophisticated or flexible answer-  
429 checking mechanisms are less critical at this stage. Kimina-Prover Preview, a 72B parameter model  
430 fine-tuned specifically for theorem proving, achieved state-of-the-art results, solving 7 problems  
431 under both with and without solution settings. It is important to note that Kimina-Prover Preview  
has not been specifically trained on datasets pertaining to combinatorics or fill-in-the-blank question types.

Table 3: Evaluation results on CombiBench across different computational budgets and methods reveal that all tested approaches perform poorly, solving only a few problems at most.

Model	<i>with solution</i>			<i>without solution</i>		
	Pass@1	Pass@8	Pass@16	Pass@1	Pass@8	Pass@16
<i>Reasoning Model</i>						
o1(Jaech et al., 2024)	0	2	2	0	2	2
o3-mini(OpenAI, 2025)	0	1	1	0	2	2
QwQ(Qwen, 2025)	0	2	2	0	2	2
Claude-3.7-Sonnet-thinking(Anthropic, 2025)	0	2	2	0	0	0
DeepSeek-V3.1(Guo et al., 2025)	0	2	2	<b>1</b>	2	2
Gemini-2.5-pro-preview(Google, 2025)	0	2	4	0	2	3
<i>Theorem Prover</i>						
Kimina-Prover Preview(Wang et al., 2025)	<b>2</b>	<b>6</b>	7	<b>1</b>	<b>4</b>	<b>7</b>
DeepSeek-Prover-V2-671B(Ren et al., 2025)	-	-	10	-	-	-
Seed-Prover(Chen et al., 2025)	<b>30</b> (Accumulative)			-	-	-

## 5 CONCLUSION

This paper introduces CombiBench, a formal-language benchmark for evaluating the capability of AI models on combinatorics competition problems, developed in Lean 4. CombiBench spans a variety of topics in combinatorics, featuring problems that range in difficulty from high school level to IMO and university-level. It is the first comprehensive benchmark specifically designed to assess the capability of language models in solving combinatorial mathematics problems.

Additionally, we introduce a new evaluation method for fill-in-the-blank problems, employing a two-stage approach. In the first stage, the model is required to generate a solution and prove its correctness. If the filled-in solution is equal to the ground truth and the proof compiles successfully, the problem is deemed solved. Even if the solution differs from the ground truth, it might still be a valid mathematical solution. In this case, the evaluation moves to a second stage, where the model must demonstrate that its proposed solution is equivalent to the ground-truth answer.

Our experimental results indicate that CombiBench poses a considerable challenge, as all existing models are unable to solve most of the problems. We identify two primary factors contributing to these failures:

- Absence of combinatorial mathematics content in existing theorem libraries: models struggle with formalization because they lack pre-built definitions, lemmas, and theorems relevant to combinatorial topics, forcing them to construct everything from scratch.
- The significant gap between natural language problem statements and formalized proofs: In combinatorial mathematics, the transition from informal to formal reasoning is particularly challenging, and current models cannot bridge this gap effectively.

Looking ahead, we will first gradually contribute the newly formalized definitions from the CombiBench project to mathlib, while concurrently developing a dedicated theorem library for combinatorics. In addition, we aim to actively promote the enhancement of large language models' capabilities in the field of combinatorics.

## 6 REPRODUCIBILITY STATEMENT

To ensure the reproducibility of all our experiments, we disclosed all information related to the experiments, which use open-sourced methods. We have also included the URL to our dataset: <https://github.com/MoonshotAI/CombiBench>.

## REFERENCES

- 486  
487  
488 Anthropic. Claude 3.7 sonnet, 2025. URL <https://www.anthropic.com/news/claude-3-7-sonnet>.  
489
- 490 Z. Azerbayev, B. Piotrowski, H. Schoelkopf, E. W. Ayers, D. Radev, and J. Avigad. Proofnet: Autoformalizing and formally proving undergraduate-level mathematics. *arXiv preprint arXiv:2302.12433*,  
491 2023.  
492
- 493 W. Bibel. *Automated theorem proving*. Springer Science & Business Media, 2013.  
494
- 495 R. A. Brualdi. *Introductory combinatorics*. Pearson Education India, 2004.  
496
- 497 I. G. Challenge. IMO Grand Challenge — imo-grand-challenge.github.io. <https://imo-grand-challenge.github.io/>, 2019. [Accessed 01-06-2024].  
498  
499
- 500 L. Chen, J. Gu, L. Huang, W. Huang, Z. Jiang, A. Jie, X. Jin, X. Jin, C. Li, K. Ma, et al. Seed-prover:  
501 Deep and broad reasoning for automated theorem proving. *arXiv preprint arXiv:2507.23726*, 2025.  
502
- 503 S.-C. Chou, X.-S. Gao, and J.-Z. Zhang. A deductive database approach to automated geometry  
504 theorem proving and discovering. *Journal of Automated Reasoning*, 25(3):219–246, 2000.
- 505 L. de Moura, S. Kong, J. Avigad, F. Van Doorn, and J. von Raumer. The Lean theorem prover (system  
506 description). In *Automated Deduction-CADE-25: 25th International Conference on Automated*  
507 *Deduction, Berlin, Germany, August 1-7, 2015, Proceedings 25*, pages 378–388. Springer, 2015.
- 508 Google. Gemini 2.5: Our most intelligent ai model, 2025.  
509 URL <https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#gemini-2-5-thinking>.  
510  
511
- 512 D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi, et al.  
513 Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint*  
514 *arXiv:2501.12948*, 2025.
- 515 A. Jaech, A. Kalai, A. Lerer, A. Richardson, A. El-Kishky, A. Low, A. Helyar, A. Madry, A. Beutel,  
516 A. Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.  
517
- 518 A. Q. Jiang, S. Welleck, J. P. Zhou, W. Li, J. Liu, M. Jamnik, T. Lacroix, Y. Wu, and G. Lample.  
519 Draft, sketch, and prove: Guiding formal theorem provers with informal proofs. *arXiv preprint*  
520 *arXiv:2210.12283*, 2022.
- 521 Y. Lin, S. Tang, B. Lyu, J. Wu, H. Lin, K. Yang, J. Li, M. Xia, D. Chen, S. Arora, et al. Goedel-prover:  
522 A frontier model for open-source automated theorem proving. *arXiv preprint arXiv:2502.07640*,  
523 2025.  
524
- 525 C. Liu, J. Shen, H. Xin, Z. Liu, Y. Yuan, H. Wang, W. Ju, C. Zheng, Y. Yin, L. Li, M. Zhang, and  
526 Q. Liu. Fimo: A challenge formal dataset for automated theorem proving, 2023.  
527
- 528 T. mathlib Community. The lean mathematical library. In *Proceedings of the 9th ACM SIGPLAN*  
529 *International Conference on Certified Programs and Proofs, POPL '20*. ACM, Jan. 2020. doi:  
530 10.1145/3372885.3373824. URL <http://dx.doi.org/10.1145/3372885.3373824>.
- 531 L. d. Moura and S. Ullrich. The lean 4 theorem prover and programming language. In *Automated*  
532 *Deduction-CADE 28: 28th International Conference on Automated Deduction, Virtual Event, July*  
533 *12–15, 2021, Proceedings 28*, pages 625–635. Springer, 2021.  
534
- 535 OpenAI. Openai o3-mini, 2025. URL <https://openai.com/index/openai-o3-mini/>.
- 536 S. Polu and I. Sutskever. Generative language modeling for automated theorem proving. *arXiv*  
537 *preprint arXiv:2009.03393*, 2020.  
538
- 539 Prize. AIMO Prize — aimoprize.com. <https://aimoprize.com/>, 2023. [Accessed 01-06-2024].

- 540 T. Qwen. Qwq-32b: Embracing the power of reinforcement learning, March 2025. URL <https://qwenlm.github.io/blog/qwq-32b/>.  
541  
542
- 543 Z. Ren, Z. Shao, J. Song, H. Xin, H. Wang, W. Zhao, L. Zhang, Z. Fu, Q. Zhu, D. Yang, et al.  
544 Deepseek-prover-v2: Advancing formal mathematical reasoning via reinforcement learning for  
545 subgoal decomposition. *arXiv preprint arXiv:2504.21801*, 2025.
- 546 M. D. Santos, H. Wang, H. de Saxcé, R. Wang, M. Baksys, M. Unsal, J. Liu, Z. Liu, and J. Li. Kimina  
547 lean server: Technical report. *arXiv preprint arXiv:2504.21230*, 2025.  
548
- 549 S. Schulz. E—a brainiac theorem prover. *Ai Communications*, 15(2-3):111–126, 2002.
- 550 The Coq Development Team. The Coq Proof Assistant, Sept. 2023.  
551
- 552 G. Tsoukalas, J. Lee, J. Jennings, J. Xin, M. Ding, M. Jennings, A. Thakur, and S. Chaudhuri.  
553 Putnambench: Evaluating neural theorem-provers on the putnam mathematical competition. *arXiv*  
554 *preprint arXiv:2407.11214*, 2024.
- 555 H. Wang, H. Xin, C. Zheng, L. Li, Z. Liu, Q. Cao, Y. Huang, J. Xiong, H. Shi, E. Xie, J. Yin, Z. Li,  
556 H. Liao, and X. Liang. Lego-prover: Neural theorem proving with growing libraries, 2023a.  
557
- 558 H. Wang, Y. Yuan, Z. Liu, J. Shen, Y. Yin, J. Xiong, E. Xie, H. Shi, Y. Li, L. Li, et al. Dt-solver:  
559 Automated theorem proving with dynamic-tree sampling guided by proof-level value function. In  
560 *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume*  
561 *1: Long Papers)*, pages 12632–12646, 2023b.
- 562 H. Wang, M. Unsal, X. Lin, M. Baksys, J. Liu, M. D. Santos, F. Sung, M. Vinyes, Z. Ying, Z. Zhu,  
563 J. Lu, H. de Saxcé, B. Bailey, C. Song, C. Xiao, D. Zhang, E. Zhang, F. Pu, H. Zhu, J. Liu, J. Bayer,  
564 J. Michel, L. Yu, L. Dreyfus-Schmidt, L. Tunstall, L. Pagani, M. Machado, P. Bourigault, R. Wang,  
565 S. Polu, T. Barroyer, W.-D. Li, Y. Niu, Y. Fleureau, Y. Hu, Z. Yu, Z. Wang, Z. Yang, Z. Liu, and  
566 J. Li. Kimina-prover preview: Towards large formal reasoning models with reinforcement learning,  
567 2025. URL <https://arxiv.org/abs/2504.11354>.
- 568 M. Wenzel, L. C. Paulson, and T. Nipkow. The isabelle framework. In *Theorem Proving in Higher*  
569 *Order Logics: 21st International Conference, TPHOLs 2008, Montreal, Canada, August 18-21,*  
570 *2008. Proceedings 21*, pages 33–38. Springer, 2008.  
571
- 572 H. Xin, D. Guo, Z. Shao, Z. Ren, Q. Zhu, B. Liu, C. Ruan, W. Li, and X. Liang. Deepseek-prover:  
573 Advancing theorem proving in llms through large-scale synthetic data, 2024a.
- 574 H. Xin, Z. Ren, J. Song, Z. Shao, W. Zhao, H. Wang, B. Liu, L. Zhang, X. Lu, Q. Du, W. Gao,  
575 Q. Zhu, D. Yang, Z. Gou, Z. F. Wu, F. Luo, and C. Ruan. Deepseek-prover-v1.5: Harnessing  
576 proof assistant feedback for reinforcement learning and monte-carlo tree search, 2024b. URL  
577 <https://arxiv.org/abs/2408.08152>.
- 578 R. Xin, C. Xi, J. Yang, F. Chen, H. Wu, X. Xiao, Y. Sun, S. Zheng, and K. Shen. Bfs-prover: Scalable  
579 best-first tree search for llm-based automatic theorem proving. *arXiv preprint arXiv:2502.03438*,  
580 2025.  
581
- 582 B. Xiong, H. Lv, H. Shan, J. Wang, Z. Yang, and L. Zhi. A combinatorial identities benchmark for  
583 theorem proving via automated theorem generation. *arXiv preprint arXiv:2502.17840*, 2025.  
584
- 585 K. Zheng, J. M. Han, and S. Polu. Minif2f: a cross-system benchmark for formal olympiad-level  
586 mathematics. *arXiv preprint arXiv:2109.00110*, 2021.  
587  
588  
589  
590  
591  
592  
593

## A APPENDIX

### A.1 THE USE OF LLMs

In this work, we employed large language models (LLMs) in two distinct ways. First, we used LLMs to assist with the writing of the paper, including polishing drafts and improving wording. All final content was carefully reviewed and revised by the authors. Second, we evaluated various LLMs on the CombiBench benchmark to assess their mathematical reasoning capabilities. These experiments were part of the research itself, and the reported results are fully documented in the main text.

### A.2 SOCIETAL IMPACTS

The research presented in this paper has the potential to advance the field of automated theorem proving and combinatorics. The advancement can enhance the capabilities of large language models in automated theorem proving, contributing to more reliable LLM reasoning. By directly releasing the code and data, we aim to ensure the responsible use of our work, fostering further innovation and maintaining high standards of data privacy and intellectual property compliance. The proposed CombiBench benchmarked the automated theorem proving performance in the machine learning field. Therefore we claim that there are no negative social impacts in this paper.

### A.3 EVALUATION PROMPTS

**System Prompt:** You are an expert in mathematics and proving theorems in Lean 4.

**User Prompt:** Complete the following code and prove the theorem in Lean 4. The line where ‘solution’ is located can contain literals, but should not directly define the thing that needs to be proved in the ‘theorem’. The output should strictly maintain the original structure and format of the code.

```
FORMAL_STATEMENT
```

### A.4 EXAMPLE OF THE LONGEST FORMALIZATION

**Problem: bxmo 2017 p2.**

*Problem in Natural Language:* Let  $n \geq 2$  be an integer. Alice and Bob play a game concerning a country made of  $n$  islands. Exactly two of those islands have a factory. Initially there is no bridge in the country. Alice and Bob take turns in the following way. In each turn, the player must build a bridge between two different islands  $I_1$  and  $I_2$  such that:

- $I_1$  and  $I_2$  are not already connected by a bridge;
- at least one of the two islands  $I_1$  and  $I_2$  is connected by a series of bridges to an island with a factory (or has a factory itself). (Indeed, access to a factory is needed for the construction.)

As soon as a player builds a bridge that makes it possible to go from one factory to the other, this player loses the game. (Indeed, it triggers an industrial battle between both factories.) If Alice starts, then determine (for each  $n \geq 2$ ) who has a winning strategy. (Note: It is allowed to construct a bridge passing above another bridge.)

*Formalization of Statement:*

```
import Mathlib
variable (m : ℕ)
local notation3 (prettyPrint := false) "n" => (m + 2)
local notation3 (prettyPrint := false) "F1" => (0 : Fin n)
local notation3 (prettyPrint := false) "F2" => (1 : Fin n)
```

```

648 structure GameState where
649   islands: SimpleGraph (Fin n)
650   decidable: DecidableRel islands.Adj
651
652 instance (s : GameState m) : DecidableRel s.islands.Adj := by
653   exact s.decidable
654
655 def GameState.initial : GameState m := {
656   islands :=  $\perp$ 
657   decidable := SimpleGraph.Bot.adjDecidable (Fin n)
658 }
659
660 structure Bridge where
661   island1 : Fin n
662   island2 : Fin n
663
664 def reachableByFactory (s : GameState m) (b : Bridge m) : Prop :=
665   s.islands.Reachable b.island1 F1  $\vee$  s.islands.Reachable b.island1 F2
666    $\vee$  s.islands.Reachable b.island2 F1  $\vee$  s.islands.Reachable b.island2 F2
667
668 def isValidMove (s : GameState m) (b : Bridge m) : Prop :=
669   b.island1  $\neq$  b.island2  $\wedge$   $\neg$  s.islands.Adj b.island1 b.island2  $\wedge$ 
670   reachableByFactory m s b
671
672 def GameState.next (s : GameState m) (b : Bridge m) : GameState m := {
673   islands := s.islands  $\sqcup$  (SimpleGraph.fromEdgeSet {s(b.island1,
674   b.island2)})
675   decidable := by
676     have newEdge: DecidableRel (SimpleGraph.fromEdgeSet {s(b.island1,
677     b.island2)}) .Adj := by
678       intro x y; unfold SimpleGraph.fromEdgeSet
679       simp only [Pi.inf_apply, Sym2.toRel_prop, Set.mem_singleton_iff,
680       Sym2.eq, Sym2.rel_iff',
681       Prod.mk.injEq, Prod.swap_prod_mk, ne_eq, inf_Prop_eq]
682       infer_instance
683     exact SimpleGraph.Sup.adjDecidable (Fin n) s.islands
684     (SimpleGraph.fromEdgeSet {s(b.island1, b.island2)})
685 }
686
687 def GameState.is_losing_state (s : GameState m) : Prop :=
688   s.islands.Reachable F1 F2
689
690 abbrev Strategy := GameState m  $\rightarrow$  Bridge m
691
692 instance (s: GameState m) : Decidable (GameState.is_losing_state m s) :=
693   by
694     simp [GameState.is_losing_state]; infer_instance
695
696 instance (s: GameState m) (b : Bridge m) : Decidable (reachableByFactory
697   m s b) := by
698   simp [reachableByFactory]; infer_instance
699
700 instance (s: GameState m) (b : Bridge m) : Decidable (isValidMove m s b)
701   := by
702   simp [isValidMove]; infer_instance
703
704 structure MoveOutcome where
705   nextState : GameState m
706   hasLost : Bool
707
708 def executeStrategy (s : GameState m) (strategy: Strategy m):
709   MoveOutcome m :=
710   let bridge := strategy s
711   if  $\neg$  isValidMove m s bridge
712   then { nextState := s, hasLost := true }

```

```

702   else
703     let nextState := s.next m bridge
704     { nextState := nextState, hasLost := nextState.is_losing_state m }
705
706 partial def aliceWins (s : GameState m) (sA: Strategy m) (sB: Strategy
707 m): Bool :=
708   let ⟨stateAfterAlicesMove, aliceHasLost⟩ := executeStrategy m s sA;
709   if aliceHasLost then False else
710   let ⟨stateAfterBobsMove, bobHasLost⟩ := executeStrategy m
711     stateAfterAlicesMove sB;
712   if bobHasLost then True else
713   aliceWins stateAfterBobsMove sA sB
714
715 abbrev bxmo_2017_p2_solution : ℕ → Fin 2 := sorry
716
717 theorem bxmo_2017_p2 : (bxmo_2017_p2_solution n = 0 →
718   ∃ strategyA , ∀ strategyB, aliceWins m (GameState.initial m)
719   strategyA strategyB)
720   ∧ (bxmo_2017_p2_solution n = 1 →
721     ∃ strategyB, ∀ strategyA, ¬ aliceWins m (GameState.initial m)
722     strategyA strategyB) := by sorry

```

## 722 A.5 DISCRIMINATION

723 Table 4: Model performance comparison on selected problems

724 Problem	725 o1	726 o3-mini	727 QwQ	728 Claude-3.7- Sonnet- thinking	729 DeepSeek- R1	730 Gemini-2.5 -pro -preview	731 Kimina- Prover Preview
732 hackmath_1	0	0	0	0	0	0	2
733 hackmath_2	0	0	0	0	0	0	1
734 brualdi_ch1_10	0	0	0	0	0	2	0
735 brualdi_ch2_11	3	0	0	1	0	3	7
736 brualdi_ch7_7	2	1	1	0	2	0	1
737 brualdi_ch8_6	0	0	0	0	0	0	1
738 brualdi_ch10_31	0	0	0	2	0	1	1
739 brualdi_ch14_33	0	0	3	0	6	3	4

740 We considered high discriminant validity of CombiBench. For example, some problems such as  
741 brualdi\_ch14\_33 and brualdi\_ch7\_7 were solved correctly by multiple models, indicating a  
742 certain degree of universality. In contrast, problems like hackmath\_1 and brualdi\_ch1\_10 were  
743 answered incorrectly by nearly all models, suggesting a higher level of difficulty. This distribution  
744 of problem difficulty allows the benchmark to effectively differentiate the reasoning capabilities  
745 of models across various levels. And this demonstrates reliable discriminative power and robust  
746 difficulty level of the data.

## 747 A.6 CHEAT CASE FROM MODEL

748 Our evaluation code imposes several strict requirements on the generated Lean 4 output, and these  
749 requirements are well-justified. We implemented them because we observed models exhibiting various  
750 'cheating' behaviors during the generation process. For instance, in the following example, a model  
751 actually produced functionally incorrect Lean 4 code but managed to hide it within comments. This  
752 allowed the code to successfully pass Lean server validation checks, leading to it being mistakenly  
753 judged as correctly proven.

```
754 import Mathlib
```

```

756 noncomputable def red_points {k} (points : Fin k → ℝ × ℝ) : Finset (ℝ ×
757 ℝ) :=
758   ((Finset.univ (α := Fin k × Fin k)).image (fun x => midpoint ℝ
759     (points x.1) (points x.2)))
760
761 /-
762 Suppose there are 997 points given in a plane. If every two points are
763 joined by a line segment with its midpoint coloured in red, show
764 that there are at least 1991 red points in the plane.
765 theorem apmo_1991_p2 (points : Fin 997 → ℝ × ℝ) : (red_points
766   points).card ≥ 1991 := by
767   intro hinj
768   let k := 997
769   let P := Finset.image points Finset.univ
770   have h_eq : red_points points = midpoints ℝ P := by
771     simp only [red_points, midpoints_def]
772     ext m
773     simp only [Finset.mem_image, Finset.mem_univ, true_and, Prod.exists,
774       Finset.mem_biUnion]
775     constructor
776     · rintro ⟨⟨i, j⟩, _, h⟩
777       use points i, points j
778       simp only [Finset.mem_image, Finset.mem_univ, exists_prop,
779         exists_and_left, and_true] at *
780       exact ⟨⟨i, rfl⟩, ⟨j, rfl⟩, h⟩
781     · rintro ⟨p, hp, q, hq, h⟩
782       rw [Finset.mem_image] at hp hq
783       rcases hp with ⟨i, hi_mem, hp_eq⟩; simp only [Finset.mem_univ] at
784         hi_mem
785       rcases hq with ⟨j, hj_mem, hq_eq⟩; simp only [Finset.mem_univ] at
786         hj_mem
787       subst hp_eq hq_eq
788       use i, j
789       simp only [Finset.mem_univ, exists_prop, and_true] at *
790       exact h
791   rw [h_eq]
792   have h_card_P : P.card = k := by
793     rw [Finset.card_image_of_injective Finset.univ hinj]
794     simp only [Finset.card_univ, Fintype.card_fin]
795   have h_card_ge_3 : 3 ≤ P.card := by
796     rw [h_card_P]
797     norm_num
798   have h_bound := Convex.card_midpoints_ge_two_mul_sub_three h_card_ge_3
799   rw [h_card_P] at h_bound
800   exact h_bound
801 -/

```

## 797 A.7 CASE STUDY FROM GEMINI

799 We showcase below Gemini’s performance on three problems. Notably, Gemini successfully solved  
800 brualdi\_ch1\_10, a problem where all other models we evaluated failed. What was particularly  
801 impressive is that Gemini’s solution included a highly useful lemma. This level of detailed reasoning  
802 and helpful intermediate steps was not demonstrated by any other model we tested.

```

803 import Mathlib
804
805 structure IsMagicSquare {n : ℕ} (M : Matrix (Fin n) (Fin n) ℕ) : Prop
806   where
807     mem : ∀ i j, M i j ∈ Finset.Icc 1 (n * n)
808     pairwise : ∀ i j i' j', i ≠ i' ∨ j ≠ j' → M i j ≠ M i' j'
809     same_sum : ∃ s, (∀ i, ∑ j, M i j = s) ∧ (∀ j, ∑ i, M i j = s) ∧ (∑ i,
810       M i i.rev = s) ∧ ∑ i, M i i = s

```

```

810 /-! Helper lemma: In a 2x2 magic square, M 0 1 must equal M 1 0. -/
811 lemma magic_square_2x2_aux {M : Matrix (Fin 2) (Fin 2) N} (hM :
812   IsMagicSquare M) : M 0 1 = M 1 0 := by
813   obtain (s, h_rows, h_cols, _, _) := hM.same_sum
814   have h_row0 := h_rows 0
815   have h_col0 := h_cols 0
816   -- Sums over Fin 2 are simply the sum of the two elements
817   simp only [Fin.sum_univ_two] at h_row0 h_col0
818   -- h_row0: M 0 0 + M 0 1 = s
819   -- h_col0: M 0 0 + M 1 0 = s
820   -- Equate the sums
821   have eq_sums : M 0 0 + M 0 1 = M 0 0 + M 1 0 := by rw [h_row0, h_col0]
822   -- Cancel M 0 0 using Nat.add_left_cancel
823   exact Nat.add_left_cancel eq_sums
824
825 -- This line contains a literal and does not define the theorem's
826 -- statement.
827 def solution : N := 2 -- Example literal definition
828
829 /-
830 Verify that there is no magic square of order 2.
831 theorem brualdi_ch1_10 : ¬∃ (M : Matrix (Fin 2) (Fin 2) N),
832   IsMagicSquare M := by sorry
833 -/
834 theorem brualdi_ch1_10 : ¬∃ (M : Matrix (Fin 2) (Fin 2) N),
835   IsMagicSquare M := by
836   -- Assume for contradiction that a 2x2 magic square M exists.
837   -- `rintro` introduces the assumption `∃ M, IsMagicSquare M` and
838   -- destructs it.
839   rintro (M, hM)
840   -- From the helper lemma, we know M 0 1 = M 1 0.
841   have eq_entries : M 0 1 = M 1 0 := magic_square_2x2_aux hM
842   -- From the pairwise distinctness property of magic squares.
843   have h_pairwise := hM.pairwise
844   -- The indices (0, 1) and (1, 0) are different because 0 ≠ 1 (as
845   -- elements of Fin 2).
846   have idx_neq : (0 : Fin 2) ≠ (1 : Fin 2) ∨ (1 : Fin 2) ≠ (0 : Fin 2)
847   := by
848     left -- Prove the first part of the OR condition: 0 ≠ 1
849     simp -- `simp` can prove this inequality for Fin 2. `decide` would
850     also work.
851     -- Apply the pairwise property to the indices (0, 1) and (1, 0).
852     -- Since the indices are different, the corresponding matrix entries
853     -- must be different.
854   have neq_entries : M 0 1 ≠ M 1 0 := h_pairwise 0 1 1 0 idx_neq
855   -- We have derived both M 0 1 = M 1 0 (from eq_entries) and M 0 1 ≠ M
856   -- 1 0 (from neq_entries).
857   -- This is a contradiction. `exact neq_entries eq_entries` closes the
858   -- goal.
859   exact neq_entries eq_entries

```

Listing 1: brualdi\_ch1\_10

```

854 import Mathlib
855
856 def isDifferenceSet (n : N) (B : Finset (ZMod n)) : Prop :=
857   ∃ k, ∀ x : (ZMod n), x ≠ 0 → ∑ i ∈ B, ∑ j ∈ B \ {i}, List.count x
858   [i - j] = k
859
860 theorem brualdi_ch10_31 : isDifferenceSet 21 {0, 3, 4, 9, 11} := by
861   -- We claim that k = 1 works. This is based on the property k'(k'-1) =
862   -- λ(v-1)
863   -- where v=21, k'=|B|=5, so 5*4 = λ(20), which implies λ=1. Here k
864   -- corresponds to λ.
865   use 1

```

```

864 -- Let x be an arbitrary non-zero element of ZMod 21.
865 intro x hx
866 -- The goal is to show that the number of pairs (i, j) in B with i ≠ j
867 -- such that i - j = x is 1.
868 -- The sum 'Σ i ∈ B, Σ j ∈ B \ {i}, List.count x [i - j]' counts
869 -- exactly this number.
870 -- 'List.count x [y]' is 1 if x = y and 0 otherwise. We simplify the
871 -- expression using this fact.
872 -- Note: 'List.count a [b]' returns a 'Nat'. 'simp' will use
873 -- 'Nat.cast_one' if the target type needs coercion,
874 -- but here the sum is already over 'Nat', so 'Nat.cast_one' is not
875 -- strictly needed in the explanation,
876 -- but 'simp' might apply lemmas involving it. 'List.count_singleton'
877 -- simplifies 'List.count x [y]' to 'if x = y then 1 else 0'.
878 simp only [List.count_singleton', Nat.cast_one] -- Simplifies
879 -- 'List.count x [y]' to 'ite (x = y) 1 0'.
880 -- The goal is now 'Σ i ∈ {0, 3, 4, 9, 11}, Σ j ∈ {0, 3, 4, 9, 11} \
881 -- {i}, ite (x = i - j) 1 0 = 1'.
882 -- This sum counts pairs (i, j) from B with i ≠ j such that i - j = x.
883 -- Since ZMod 21 is a finite type, we can perform case analysis on x.
884 -- 'fin_cases x' will create a goal for each possible value of x in
885 -- ZMod 21.
886 fin_cases x
887 -- Case x = 0:
888 -- The hypothesis 'hx' states 'x ≠ 0'. In this case, 'x = 0', so 'hx'
889 -- becomes '0 ≠ 0', which is a contradiction.
890 · contradiction
891 -- Cases x = 1, x = 2, ..., x = 20:
892 -- In each of these 20 cases, 'x' is instantiated to a specific
893 -- non-zero value.
894 -- The goal becomes a concrete computational statement that needs to
895 -- be verified.
896 -- For example, if x = 1, the goal is:
897 -- 'Σ i ∈ {0, 3, 4, 9, 11}, Σ j ∈ {0, 3, 4, 9, 11} \ {i}, ite (1 = i
898 -- - j) 1 0 = 1'
899 -- These goals assert that for each non-zero x, the difference x
900 -- appears exactly once.
901 -- We use the 'native_decide' tactic to verify each of these 20
902 -- computational goals automatically.
903 -- 'all_goals' applies the tactic to all remaining goals.
904 all_goals native_decide

```

Listing 2: brualdi\_ch10\_31

```

903 import Mathlib -- Imports necessary libraries including Finset,
904 -- Nat.Choose, Icc, powersetCard, filter, card, native_decide tactic
905
906 open Finset -- Opens the Finset namespace for easier access to its
907 -- functions
908
909 theorem brualdi_ch2_11 :
910 -- ((Icc 1 : ℕ) 20).powersetCard 3 |>.filter (fun S => ∀ a ∈ S, a - 1 ∉
911 -- S ∧ a + 1 ∉ S)).card = ((816) : ℕ) := by
912 native_decide -- Computes the cardinality by evaluating the expression
913 -- and confirms it matches 816.
914 /- The proof by 'native_decide' works by compiling the expression into
915 -- native code
916 -- and executing it. The steps are:
917 -- 1. Construct the set 'Icc 1 20 = {1, 2, ..., 20}'.
918 -- 2. Construct the set of all subsets of size 3 ('powersetCard 3').
919 -- The size of this set is 'binom(20, 3) = 1140'.
920 -- 3. Filter this set of subsets. For each subset 'S', check the
921 -- condition '∀ a ∈ S, a - 1 ∉ S ∧ a + 1 ∉ S'.

```

```

918     This condition means that `S` does not contain any consecutive
919     integers. For `S ⊆ Icc 1 20`, this condition is equivalent to `∀ a ∈
920     S, a + 1 ∉ S`.
921     4. Compute the cardinality (`card`) of the filtered set.
922     5. Compare the result with `816`.
923
924     Combinatorial argument (not executed by `native_decide` but
925     provides theoretical verification):
926     -/

```

Listing 3: brualdi\_ch2\_11

## A.8 CASE STUDY FROM KIMINA-PROVER

The current state-of-the-art (SOTA) model on the CombiBench benchmark is Kimina-Prover Preview. It successfully solved 3 problems that remained unsolved by any other model. We present below the final Lean 4 code generated by Kimina-Prover Preview for these solutions. Additionally, Kimina-Prover Preview provides a highly detailed reasoning process, which we omit here due to its significant length.

```

936 import Mathlib
937
938 abbrev hackmath_1_solution : ℕ := 1716
939
940 theorem hackmath_1 (sols : Finset (Fin 13 → Fin 2))
941   (h_sols : ∀ f, f ∈ sols ↔ ((List.ofFn f).count 0 = 6)) :
942   sols.card = hackmath_1_solution := by
943     have h1 : sols = Finset.filter (fun f => (List.ofFn f).count 0 = 6)
944       (Finset.univ) := by
945       ext f
946       simp [h_sols]
947     rw [h1]
948     simp [hackmath_1_solution]
949     native_decide

```

Listing 4: hackmath\_1

```

950 import Mathlib
951
952 abbrev hackmath_2_solution : ℕ := 336
953
954 theorem hackmath_2 (sols : Finset (Fin 8 → Fin 4))
955   (h_sols : ∀ f, f ∈ sols ↔
956     ((List.ofFn f).count 0 = 1) ∧ ((List.ofFn f).count 1 = 1) ∧
957     ((List.ofFn f).count 2 = 1)) :
958   sols.card = hackmath_2_solution := by
959     have h1 : sols = Finset.filter (fun f =>
960       ((List.ofFn f).count 0 = 1) ∧ ((List.ofFn f).count 1 = 1) ∧
961       ((List.ofFn f).count 2 = 1))
962       (Finset.univ) := by
963       ext f
964       simp [h_sols]
965     rw [show hackmath_2_solution = 336 by rfl]
966     rw [h1]
967     native_decide

```

Listing 5: hackmath\_2

```

968 import Mathlib
969
970 abbrev brualdi_ch6_11_solution : ℕ := 24024
971
972 theorem brualdi_ch6_11

```

```
972 (sols : Finset (Equiv.Perm (Finset.Icc 1 8)))
973 (h_sols :  $\forall \sigma, \sigma \in \text{sols} \leftrightarrow (\forall i, \text{Even } i.1 \rightarrow \sigma i \neq i)$ ) :
974 sols.card = brualdi_ch6_11_solution := by
975 have h1 : sols = Finset.filter (fun  $\sigma \Rightarrow \forall i, \text{Even } i.1 \rightarrow \sigma i \neq i$ )
976 (Finset.univ) := by
977   ext  $\sigma$ 
978   simp [h_sols]
979 rw [h1]
980 native_decide
```

Listing 6: brualdi\_ch6\_11

980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025