

# ALCHEMY: AMPLIFYING THEOREM-PROVING CAPABILITY THROUGH SYMBOLIC MUTATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Formal proofs are challenging to write even for experienced experts. Recent progress in Neural Theorem Proving (NTP) shows promise in expediting this process. However, the formal corpora available on the Internet are limited compared to the general text, posing a significant data scarcity challenge for NTP. To address this issue, this work proposes *Alchemy*, a general framework for data synthesis that constructs formal theorems through symbolic mutation. Specifically, for each candidate theorem in Mathlib, we identify all invocable theorems that can be used to rewrite or apply to it. Subsequently, we mutate the candidate theorem by replacing the corresponding term in the statement with its equivalent form or antecedent. As a result, our method increases the number of theorems in Mathlib by an order of magnitude, from 110k to 6M. Furthermore, we perform continual pretraining and supervised finetuning on this augmented corpus for large language models. Experimental results demonstrate the effectiveness of our approach, achieving a 5% absolute performance improvement on *Leandoho* benchmark. Additionally, our synthetic data achieve a 2.5% absolute performance gain on the out-of-distribution miniF2F benchmark. To provide further insights, we conduct a comprehensive analysis of synthetic data composition and the training paradigm, offering valuable guidance for developing a strong theorem prover.

## 1 INTRODUCTION

Nowadays, some pioneer mathematicians are attempting to verify their proofs using the proof assistant Lean (de Moura et al., 2015; Tao, 2023). Writing proofs for formal statements demands mastery of formal language and domain-specific mathematical knowledge. To mitigate the complexity associated with completing proofs, several research efforts (Polu & Sutskever, 2020; Polu et al., 2023; Trinh et al., 2024) seek to automatically generate formalized proof through a neural model, known as Neural Theorem Proving (NTP). NTP represents a long-standing challenge for machine learning-based methods (Li et al., 2024), highlighting the limitations in the reasoning abilities of neural models. Prevalent Large Language Models (LLMs) (Brown et al., 2020; Dubey et al., 2024) still struggle with theorem-proving, despite excelling in related reasoning-intensive scenarios such as math reasoning (Reid et al., 2024) or code generation (Guo et al., 2024).

The key challenge of theorem-proving lies in data scarcity (Li et al., 2024; Trinh et al., 2024). Due to the difficulties associated with the manual formalization of theorems, formal corpora available on the Internet are relatively scarce compared to the general text (Azerbayev et al., 2023). Synthetic data has shown promise in alleviating the data scarcity problem. Some works propose to directly create theorems in symbolic space. For instance, Wang & Deng (2020) attempts to train a neural theorem generator on human-written formal theorems for the low-weighted formal system Metamath. Other efforts focus on generating theorems based on symbolic rules (Wu et al., 2021; Trinh et al., 2024), which are restricted to a specific domain of mathematics, such as inequality theorems and 2D geometry. Additionally, there are endeavors focused on autoformalization (Xin et al., 2024; Ying et al., 2024), which typically translates natural language mathematical problems into formalized statements, samples correct proofs, and retrains the theorem prover iteratively. Autoformalization has yielded promising results in competition-level theorem-proving tasks through the use of large autoformalized datasets. However, the process of formalizing problems and retrieving proofs is labor-intensive and cost-prohibitive. The distribution of formalized theorems is constrained by the pool of human-collected natural language problems and the intrinsic capabilities of the model.

054 Compared to autoformalization, synthesizing theorems in symbolic space is a more direct process  
 055 without intermediate translation, and is also easier to scale up to large, cost-effective CPU units.  
 056

057 Building upon the advanced Lean theorem prover, we introduce a general method that synthesizes  
 058 theorems directly in symbolic space. We analogize theorem synthesis to constructing functions in  
 059 general programming language and adopt an up-to-down approach. Initially, a new statement (func-  
 060 tion declaration) is constructed for each candidate theorem. Specifically, with the mathematical  
 061 library of Lean Mathlib<sup>4</sup> as seed data, we aim to find a symbolic manipulation  $\Phi$  between two  
 062 existing statements. We posit that Lean’s tactics serve as suitable candidates for manipulation  $\Phi$   
 063 because of their efficacy in handling symbolic expressions.  $\{rw, apply\}$  are basic tactics frequently  
 064 used in theorem proving and capable of handling the equality and implication relationship between  
 065 terms. We assign both tactics to the set of manipulations  $\Phi$  and retrieve the invocable theorems for  
 066 each candidate theorem by executing a predefined list of instructions in an interactive Lean environ-  
 067 ment. Then we mutate the candidate statement by replacing its components with their corresponding  
 068 equivalent forms or logical antecedents. Ultimately, we construct the corresponding proof (function  
 069 body) based on the existing proof and verify its correctness using Lean. The worked example shown  
 070 in Fig 1 illustrates the entire procedure of our algorithm. This algorithm is executed on a large CPU-  
 071 only computing unit for several days. Our method increases the number of theorems in Mathlib  
 072 by an order of magnitude from 110,657 to 6,326,679. This significant increase in the number of  
 073 theorems demonstrates the potential of creating theorems in symbolic space.

074 We pre-train the LLMs on the combination of Mathlib theorems and their mutated variants. Then  
 075 we fine-tune the models on the extracted state-tactic pairs, composing both the training split of  
 076 Mathlib and additional synthesized state-tactic pairs. We demonstrate the effectiveness of our  
 077 method by evaluating the theorem-proving capability of these provers on the challenging *Lean-  
 078 dojo* benchmark. Our synthetic data improve the performance by around 5% (over 70 theorems) on  
 079 the novel\_premises split. Furthermore, the synthesized data exhibit promise in enhancing the out-  
 080 of-distribution theorem-proving ability of LLMs, as evidenced by a performance increase of about  
 081 2.5% on the competition-level miniF2F benchmark.

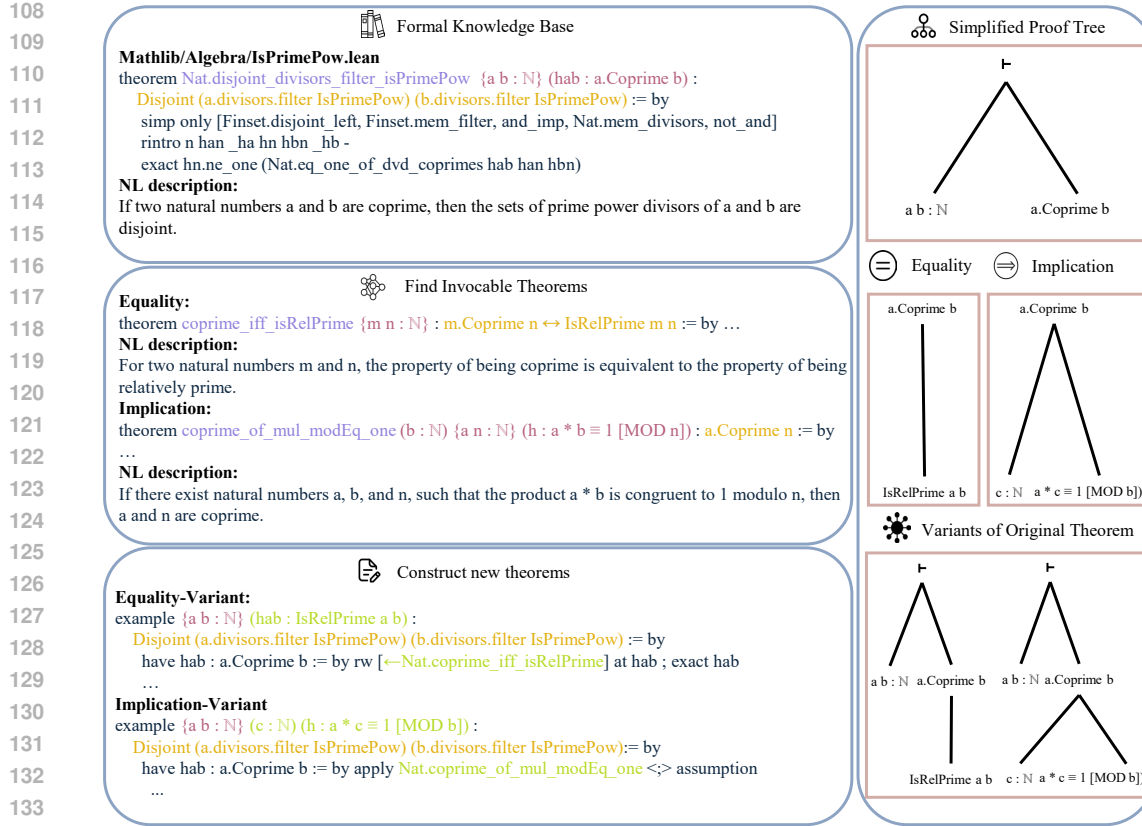
082 Our main contributions are as follows. 1) To the best of our knowledge, this work represents the  
 083 first general data synthesis framework in the symbolic space for the Lean theorem prover, effectively  
 084 complementing mainstream autoformalization-based methods. Notably, our synthesis pipeline in-  
 085 creases the number of theorems in Mathlib4 by an order of magnitude. 2) The synthesized data  
 086 and associated code will be made open-source to facilitate further research in data synthesis for  
 087 formal systems. Also, the synthesized theorems can serve as a valuable supplement to Mathlib. 3)  
 088 We conduct a comprehensive evaluation on both in-distribution and out-of-distribution benchmarks,  
 089 providing empirical insights to enhance the theorem-proving capabilities of LLMs.

## 090 2 RELATED WORK

091 **Neural Theorem Proving** Proof assistants such as Lean (de Moura et al., 2015), Isabelle (Paulson,  
 092 1994) or Coq (Barras et al., 1997) are gaining traction within the mathematical community. These  
 093 tools help mathematicians in interactively formalizing and checking the correctness of proofs (Tao,  
 094 2024). Neural networks have shown promise in lowering the barrier of using a specific formal  
 095 language for mathematicians, serving as a copilot (Song et al., 2024; Welleck & Saha, 2023). Polu  
 096 & Sutskever (2020) propose to prove theorems automatically by training a decoder-only transformer  
 097 to predict the next proofstep and construct the entire proof through a predefined search strategy. Then  
 098 a series of works seek to enhance the efficiency of this framework by incorporating auxiliary training  
 099 objectives (Han et al., 2022), conducting reinforcement learning (Polu et al., 2023; Xin et al., 2024),  
 100 improving proof search strategy (Lample et al., 2022; Wang et al., 2023; Xin et al., 2024), refining  
 101 the premise-selection (Mikula et al., 2023; Yang et al., 2023) and so on. Our work follows the  
 102 framework proposed by Polu & Sutskever (2020), using proofstep prediction as the objective and  
 103 best-first-search as the search strategy.

104 **Synthetic Theorem Creation** Data scarcity is a main challenge for NTP (Li et al., 2024). Syn-  
 105 thetic data can effectively alleviate this problem alongside manual data collection (Wu et al., 2024).  
 106

107 <sup>1</sup><https://github.com/leanprover-community/mathlib4>



135  
136  
137  
138  
139  
140

Figure 1: The overview of our synthesis pipeline. At the theorem level, we find invocable theorems that can be used to rewrite or apply to the assumptions or assertion of the candidate statement, such as the *iff* and implication rules about the *Coprime*. Then, we construct the new statements by replacing the specific component with its equivalent form or antecedent. At the proof tree level, our method merges two existing proof trees.

141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157

The current approach for synthesizing theorems diverges into two pathways. For autoformalization-based methods, the prevalent statement-level autoformalization is to translate a set of natural language problems into formal statements, followed by expert iteration to sample a collection of proofs for these statements (Wu et al., 2022; Xin et al., 2024; Ying et al., 2024). The proof-level autoformalization (Jiang et al., 2023; Huang et al., 2024) leverages LLM to generate a proof sketch, which is completed by symbolic engines such as Sledgehammer (Böhme & Nipkow, 2010). In contrast, the second pathway focuses on synthesizing theorems in formal space. Wang & Deng (2020) propose to train a neural theorem generator to synthesize theorems on a low-weight formal system, Metamath (Megill & Wheeler, 2019) which has only one tactic *substitute*. Wu et al. (2021) sequentially edits the seed expression according to a predefined set of axioms and an axiom order to create a new statement, concatenating the implications from all steps to build a complete proof. This method is used to create theorems on domains grounded in well-established axioms, such as inequality theorems and ring algebra (Polu & Sutskever, 2020). Beyond these works, AlphaGeometry (Trinh et al., 2024) can solve olympiad geometry without human demonstrations by constructing statements and proofs in symbolic space from scratch, using a carefully designed deduction engine and large-scale computing resources. Our method aims to directly synthesize theorems in symbolic space on the advanced Lean theorem prover, fully utilizing the power of computing.

158  
159  
160  
161

**Benchmarks for Theorem Proving** Most neural theorem provers based on Lean are primarily trained on Lean’s mathematical library, Mathlib. It encompasses a broad spectrum of mathematical subjects (e.g., algebra and analysis), composed of over 120,000 theorems along with their respective axioms and definitions. Researchers test the capability of neural models to prove in-distribution theorems on a held-out set of Mathlib (Polu & Sutskever, 2020; Han et al., 2022; Polu et al., 2023).

Yang et al. (2023) creates a challenging data split of Mathlib (*novel\_premise* split) which requires testing proofs to use at least one premises not seen in the training stage and mitigates the over-estimated phenomena in the traditional setting of evaluation (*random* split). Another widely-used benchmark, miniF2F, (Zheng et al., 2022) is a cross-system benchmark and includes competition-level problems as well as IMO-level problems in the domain of algebra and number theory.

### 3 METHOD

Theorems written in Lean can be viewed as a special form of code, where declarations and function bodies possess precise mathematical meanings. The initial step in creating a new theorem involves formulating a theorem statement (function declaration) that defines the essence of the theorem. Then, one must verify its correctness by generating a proof block (function body) and submitting it to the proof assistant for validation. The resulting theorems that pass type checking can serve as supplementary data for training a neural theorem prover.

#### 3.1 STATEMENT GENERATION

**Find invocable theorems** Constructing a new statement is the first step in creating a Lean theorem. The candidate theorem  $t$  has a statement denoted as  $s$ . In the corresponding Lean repository, there exists a set of potentially invocable theorems  $T_p = \{t_i\}_{i=0}^N$ . We assume that the challenge in creating a new theorem involves effectively leveraging the possibly invocable theorem  $t_i$  to mutate the candidate statement  $s$ . This understanding arises from two perspectives. Each theorem in Lean can be represented in the form of a proof tree as presented in Fig 1. The leaf nodes represent the assumptions, and the root node signifies the assertion. At the tree level, the task of generating a new Lean theorem with existing theorems is equivalent to defining operations  $\Phi$  that combine the proof trees of  $t_i$  and  $t$ . To streamline this process, our focus is solely on establishing the connection between the root node of  $t_i$  and the leaf node (or root node) of the candidate theorem  $t$ . From a mathematical standpoint, we can transform a target formula into an equal variant or break it down into multiple subformulas that suffice to prove the original formula, by employing the equality or “only if” relationship between formulas. The mathematical interconnections between formulas provide heuristic insights on how to mutate  $s$  to create a new theorem. Similarly, we can substitute the terms in  $s$  with their equivalent forms or logical antecedents. For instance, consider the statement  $a + b > c + d, m > 0 \rightarrow m(a + b) > m(c + d)$  and the known theorems  $a > b \iff e^a > e^b$  and  $a > c, b > d \implies a + b > c + d$ . From these, we can derive new theorems:  $a + b > c + d, m > 0 \rightarrow e^{m(a+b)} > e^{m(c+d)}$ , and  $a > c, b > d, m > 0 \implies m(a + b) > m(c + d)$ . In summary, identifying operations  $\Phi$  that use  $t_i$  to modify the assumptions or assertion of  $s$  is the primary step in constructing new statements.

With their intrinsic mathematical meanings and proficiency in manipulating terms within Lean, tactics are promising candidates for the operations  $\Phi$ . Following the preceding discussion, we choose two frequently used basic tactics, *rw* and *apply* to formulate  $\Phi$ .

- **rw** The “rewriting” tactic *rw* is mostly used to replace some terms in the target expression with their equivalent forms according to the given identity or *iff* (a.k.a., if and only if) rules<sup>2</sup>. In the presence of an identity  $h : a = b$  or an *iff* rule  $h : P \iff Q$ , *rw [h]* substitutes all occurrences of term on the left side of equality in the proof goal with term on the right side. The direction of substitution can be reversed by adding a back arrow in the bracket (*rw [← h]*). The target of rewriting can also be changed using *at*, e.g. *rw [h] at h<sub>1</sub>*, where  $h_1$  is an arbitrary assumption of the current proof state.
- **apply** The *apply* tactic is a “suffice-to” tactic. Given an implication, it will match the consequent with the proof goal. If matched, it will transform the goal into the antecedent of the implication. With an implication rule  $h : P \implies Q$  and a proof goal  $Q$ , then *apply [h]* will reduce the goal to proving  $P$ , which means that “proving P suffices to prove Q by implication”. Similarly, *apply* can be used to modify the assumption by deducing the implication forward. With assumption  $h_1 : P$ , then *apply [h] at h<sub>1</sub>* will change  $h_1$  into  $Q$ , which means “If P is true, then we can assert Q is true by the implication”.

<sup>2</sup>Strictly speaking, the *rw* tactic is used to handling equality in Lean, the identity and *iff* are just some kinds of equality.

**Algorithm 1** Find invocable theorems

```

216
217
218 Input: candidate statement  $s$ , potential invocable theorems  $T_p$ , instruction templates  $I$ 
219 Output: invocable theorems  $T_i$   $\triangleright T_i : \{(init\_state, next\_state, instruction) \dots\}$ 
220  $(env, init\_state) \leftarrow \text{INIT}(s)$   $\triangleright$  initialize gym-like environment and retrieve initial state
221  $T_i \leftarrow \emptyset$ 
222 for  $t$  in  $T_p$  do
223   for  $i$  in  $I$  do  $\triangleright$  for each instruction template
224     instruction  $inst \leftarrow \text{FORMAT}(t, i)$ 
225      $next\_state \leftarrow \text{RUN\_TAC}(env, init\_state, inst)$   $\triangleright$  run a tactic specified by instruction  $i$ 
226   and theorem  $t$ 
227     if  $\text{VALID}(next\_state)$  then  $\triangleright$  if return a valid proof state
228       Add  $(init\_state, next\_state, inst)$  to  $T_i$ 
229     end if
230   end for
231 end for

```

Table 1: Templates for instructions designed to be executed in a Lean environment. We determine if a theorem is invocable by running the specific instruction.

Tactic	Instruction Template	Description
	<b>Equality</b> invocable_theorem : $a = b$ or $a \iff b$	
<i>rw</i>	rw [invocable_theorem]	replace all $as$ in goal with $b$
	rw [ $\leftarrow$ invocable_theorem]	replace all $bs$ in goal with $a$
	rw [invocable_theorem] at assumption	replace all $as$ in assumption with $b$
	rw [ $\leftarrow$ invocable_theorem] at assumption	replace all $bs$ in assumption with $a$
	<b>Implication</b> invocable_theorem : $a \implies b$	
<i>apply</i>	have assumption := by apply invocable_theorem	set assumption as current proof goal, and try to argue backwards

To generate a new statement, we need to find the relationship between the candidate statement  $s$  and the potentially invocable theorems  $T_p$ . The pseudocode outlined in Algorithm 1 describes the main procedure to find invocable theorems. The process involves initializing a gym-like environment to interact with Lean and extracting the initial proof state for the candidate statement. Then, the algorithm iteratively tests whether one theorem can be used to rewrite or apply to the candidate theorem leveraging the instruction templates shown in Table 1. Suppose the feedback from the interactive environment is deemed valid according to predefined criteria, the algorithm adds the proof states before and after the tactic running together with the respective instruction to the set of invocable theorems  $T_i$ . More information about this process is described in Appendix C.2.

**Mutate statements** After obtaining the initial set of invocable theorems, we applied some filtering rules to  $T_i$  to improve the quality of the data and lower the complexity of mutating statements. With filtered invocable theorems  $T_i$ , we construct new statements by replacing the components with their equivalent forms or antecedents. Since we use tactics in Lean to formulate the operations  $\Phi$ , most symbolic manipulations are bypassed to the Lean proof assistant. What remains is just parsing and replacing. Specifically, for the candidate statement  $s$  and instruction  $i$ , we utilize its abstract syntax tree to pinpoint the exact location within the code that requires modification. Then we replace the corresponding parts with mutants parsing from the subsequent proof state generated by the execution of a specific tactic. The details of our algorithm and helpful source code are described in C.3.

3.2 PROOF GENERATION AND THEOREM VERIFICATION

Mutated statements can serve as useful lemmas for theorem-proving only if we can construct proofs that pass the verification of the proof assistant. We construct the entire proof using symbolic rules. Although neural provers and other automated theorem proving (ATP) tools (e.g., hammer) can generate more natural and diverse proofs than rule-based methods, they are compute-intensive and do

not guarantee the correctness of the generated proofs. The idea of building a proof block is intuitive. Given that we only make a one-step modification to the statement, transforming the original proof state to a mutated proof state, a logical approach is to reverse the mutation and utilize the original proof to complete the remaining proving process. We use *have* tactic to restore the modified part of a statement (the original assumption or assertion) by introducing a lemma.

- **have** The *have* tactic enables users to introduce new assumption into the current proof state if they can prove it. Given an assumption  $h_1 : P$  and an implication rule  $h_2 : P \implies Q$ , a new assumption  $h : Q$  can be added by *have h: Q := by apply h\_2 at h\_1; exact h\_1*. This tactic is usually used to introduce helpful lemmas when proving a theorem.

In addition to its ability to introduce new assumptions into the proof state, *have* can be used in both tactic-style proof and term-style proof, which enlarges the margin for theorems to which our method can be applied. Apart from this, the additional *have* instruction transforms the mutated complex proof state into a canonical proof state. To some extent, this transformation is analogous to constructing an auxiliary point in geometry problems, which we assume will be beneficial for theorem proving in the general domain. Subsequently, we combine the original proof with this lemma to build the proof for the new statement. The details of the implementation of proof generation are depicted in the Appendix C.3. We construct the proof block for each mutated theorem. Then we submit the synthesized theorems to the Lean theorem prover for verification and remove the wrong ones. Details of the verification process are provided in Appendix C.4. Finally, we obtain a set of variants  $V = \{v_i\}_{i=0}^n$  defined by the keyword “example” for each candidate theorem.

### 3.3 MODEL TRAINING

Regarding the synthetic data, we have two observations. At the theorem level, the synthetic data comprises numerous theorems, each with statement distinct from existing theorems. At the state-tactic level, the process of constructing proofs introduces additional state-tactic pairs, primarily centered on *rw* and *apply*. Based on these insights, we assume that the synthetic data can serve as an augmented corpus for continual pretraining and supervised finetuning. Specifically, we fine-tune LLMs using the proofstep prediction objective proposed by Polu & Sutskever (2020), utilizing state-tactic pairs derived from both seed theorems and synthetic theorems. Given the current proof state, the model is required to predict the next tactic sequence that contributes to the proving of the target theorem. We utilize the prompt template used by Welleck (2023), as shown in Fig.2.

```

/- You are proving a theorem in Lean 4.
You are given the following information:
- The current proof state, inside [STATE]...[/STATE]

Your task is to generate the next tactic in the proof. Put the next tactic inside [TAC]...[/TAC] -/
[STATE]
{state}
[/STATE]
[TAC]

```

Figure 2: Prompt template

## 4 EXPERIMENTS

We implement the data-synthesis pipeline described in Section 3 for *rw* and *apply*, constructing a set of variants for each candidate theorem in Mathlib. We train the LLMs on a mixture of human-written theorems and synthetic ones. To examine the effectiveness of synthetic data, we evaluate the theorem prover on two benchmarks that are widely adopted by the research community: 1) **Test split of Mathlib**, which shares the same distributional characteristics as the seed theorems; 2) **miniF2F**, a challenging benchmark focusing on competition-level problems that exhibits a distinct distribution

324 compared to seed data. The experimental results derived from both benchmarks demonstrate the  
 325 potential efficacy of our approach.  
 326

#### 327 4.1 IMPLEMENTATION DETAILS 328

329 **Data-Synthesis** We choose Mathlib<sup>3</sup> which contains 116,695 theorems as the seed data for data-  
 330 synthesis. Our synthesis pipeline is built upon *Leandojo*<sup>4</sup> (Yang et al., 2023), a Python module that  
 331 enables tracing a specific Lean repository, extracting the state-tactic pairs and abstract syntax trees  
 332 (ASTs), and interacting with the Lean environment<sup>5</sup> (*run\_tac API*). Finding invocable theorems is  
 333 the most time-consuming step of our pipeline. For *rw*, the time overhead amounts to 14 days using  
 334 4,096 CPU cores<sup>6</sup>. For *apply*, it takes 7 days at this stage using 2,048 CPU cores with a one-hour  
 335 timeout for each theorem. The substantial time cost is attributed to the  $O(n^2)$  complexity of our  
 336 algorithm and the memory-intensive characteristics of *Leandojo*. We believe this overhead could be  
 337 greatly reduced through a more meticulous implementation. After retrieving the invocable theorems,  
 338 we construct new statements and proofs for the target theorems in approximately an hour using 24  
 339 CPU cores. We then write back the mutated theorems and compile the enlarged repository through  
 340 *lake build*, utilizing 2,048 CPU cores. We retrieve the error messages returned by Lean, which  
 341 can be parsed to locate the wrong theorems. Finally, we trace the enlarged repository on a 96-core  
 342 machine for 3 days, obtaining the additional state-tactic pairs by parsing the AST of each file.

343 **Model Training** We select *Llama-3-8B* (Dubey et al., 2024) and *deepseek-coder-base-v1.5- 7B*  
 344 (Guo et al., 2024) as our base models. We conduct continual pretraining with the next-token predic-  
 345 tion objective for one epoch. Then we fine-tune the models with the proofstep prediction objective  
 346 (Polu & Sutskever, 2020) for two epochs. All experiments are conducted on  $8 \times H100$  GPUS. We  
 347 employ a linear learning rate scheduler with a 3% warm-up period and a maximum learning rate of  
 348  $2e-5$ . We set the global batch size to 256 and the cutoff length to 2,048. All models are trained using  
 349 *Deepspeed ZeRO Stage3* (Rajbhandari et al., 2021) and *Flash-Attention 2* (Dao, 2023). We utilize  
 350 the open-sourced codebase *Llama-Factory* (Zheng et al., 2024) for all training experiments.

351 **Evaluation** We follow the evaluation setting used in Azerbayev et al. (2023). We use the fre-  
 352 quently used best-first-search as our search strategy and set a 10-minute timeout. The search budget  
 353 can be represented as  $N \times S \times T$ , where  $N$  denotes the number of attempts,  $S$  denotes the number  
 354 of generated tactics per iteration, and  $T$  denotes the maximum number of generations. Following  
 355 Azerbayev et al. (2023), we set  $N = 1$ ,  $S = 32$  and  $T = 100$ . Our evaluation script is modified from  
 356 an open-source implementation (Welleck, 2023) which is based on *vLLM* (Kwon et al., 2023) and  
 357 *Leandojo* (Yang et al., 2023). We utilize *Leandojo Benchmark* (Yang et al., 2023) which contains  
 358 2,000 theorems as the test split of Mathlib4 and report the results on both the *random* split and the  
 359 *novel\_premises* split. We remove the subsets of theorems for both splits that can not be initialized  
 360 by *Leandojo*. There remain 1,929 theorems in *random* split and 1,659 theorems in *novel\_premises*  
 361 split. We upgrade the tool-chain version of miniF2F (Zheng et al., 2022) to *v4.6.0 rc1*.  
 362

#### 363 4.2 ANALYSIS OF SYNTHETIC DATA 364

365 Table 2: Number of theorems. Stage one: the number of invocable instructions for all candidate  
 366 theorems. Stage two: the number of theorems that pass the verification of the Lean theorem prover.

Tactic	Candidate theorems	Stage one	Stage two	Expansion	Conversion Ratio
rw	110,657	5,081,544	2,830,817	$\times 25$	56%
apply	78,871	9,483,504	3,495,832	$\times 44$	37%

372 We separately run the synthesis pipeline for these two tactics. For *rw*, we choose Mathlib theorems  
 373 as candidate theorems. Additionally, candidate theorems for *apply* should have at least one explicit  
 374

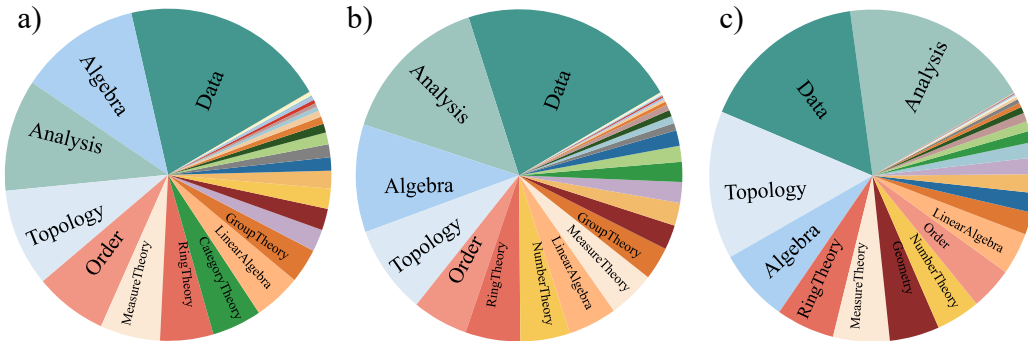
375 <sup>3</sup>commit: 2iufd 3c307701fa7e9acbd0680d7f3b9c9fed9081740

376 <sup>4</sup>version: 1.7.1

377 <sup>5</sup>lean-toolchain: v4.6.0 rc1

<sup>6</sup>512 CPU nodes, each node has 8 cores and 56GB RAM

378 assumption. In practice, the synthesis process is divided into two stages. In the first stage, we find  
 379 the potential invocable theorems for each candidate theorem by running a specific tactic. In the  
 380 second stage, we construct the new theorems and verify their correctness using the Lean theorem  
 381 prover. Table 2 shows the number of theorems of different stages. For both tactics, we increase the  
 382 number of theorems by an order of magnitude ( $\times 25$  for *rw* and  $\times 44$  for *apply*). The conversion ratios  
 383 from the potential invocable theorems to the outcomes are primarily determined by the method used  
 384 to construct the new statements and proofs. We believe that a finer implementation could greatly  
 385 improve the conversion ratio. Figure 3 shows the dynamics of the distribution of mathematical  
 386 subjects. The *rw* tactic increases the percentages of Analysis, Ring Algebra, Number Theory, and so  
 387 on. The *apply* tactic mainly contributes to the fields of Analysis and Topology. Further information  
 388 about synthetic data can be found in the Appendix D.



401 Figure 3: Distribution of mathematical subjects. For each employed tactic, we mix the generated  
 402 variants with the original theorems. a) The distribution of Mathlib. b) The distribution of Mathlib +  
 403 *rw*. c) The distribution of Mathlib + *apply*.  
 404

405  
 406 Our method synthesizes a large collection of new theorems utilizing each tactic. Then we combine  
 407 them with the theorems in Mathlib as the training data for continual pre-training. Our approach also  
 408 introduces new state-tactic pairs during the theorem-construction process. We write the variants to  
 409 corresponding lean files and extract additional state-tactic pairs using *Leandojo*. The synthesized  
 410 data are categorized primarily based on the employed tactic, specifically *rw* and *apply*. Variants  
 411 and their corresponding state-tactic pairs that appear in the test split of the *Leandojo* benchmark are  
 412 removed. Furthermore, the extracted state-tactic pairs are deduplicated according to the invocable  
 413 theorem (i.e., premise) used in the tactic instruction. Finally, we obtain about 30k data points for  
 414 each tactic. We combine them with the training set of *Leandojo* (Mathlib-train) that composes over  
 415 200k data points to form the SFT dataset. A detailed description of the deduplication process and  
 416 training data are presented in the Appendix D.3.

417 4.3 EXPERIMENTAL RESULTS

418 4.3.1 MAIN RESULTS

419  
 420 We conduct continual pretraining on the augmented lean corpus. Then we fine-tune the LLMs on  
 421 the mixture of Mathlib-train and additional state-tactic pairs. The training data are grouped by the  
 422 tactic employed in the additional state-tactic pairs. We evaluate the effectiveness of our method on  
 423 the challenging *Leandojo* benchmark and report results on different mixtures of data. As shown in  
 424 Table 3, our synthetic data consistently improve the theorem-proving capabilities of LLMs. Com-  
 425 pared with solely finetuning on the training split of Mathlib, data augmentation for a single tactic  
 426 demonstrates a beneficial effect on the theorem-proving ability of LLMs. Moreover, the positive  
 427 impacts of each tactic can be cumulative. Training on the combination of *rw* variants and *apply*  
 428 variants results in a significant performance improvement in the challenging novel\_premises split  
 429 of *Leandojo* benchmark, where the model is required to use at least one new premise to prove the  
 430 target theorem (+4.7%, 78 theorems for *Llama3-8b*; +4.22%, 70 theorems for *deepseek-coder-7b-  
 431 base-v1.5*). Our synthetic data still make a certain improvement on the random split, where the  
 performance of models is over-estimated by allowing it to prove many theorems through memoriza-



432 Table 3: Results on Mathlib. tidy: a tactic in Mathlib that uses heuristics to complete a proof. We  
 433 select the performance of each model solely fine-tuned using Mathlib-train as the main baseline.  
 434 Mathlib-train + x: the performance of the model pre-trained and fine-tuned on a mixture of Mathlib-  
 435 train and additional data about x.

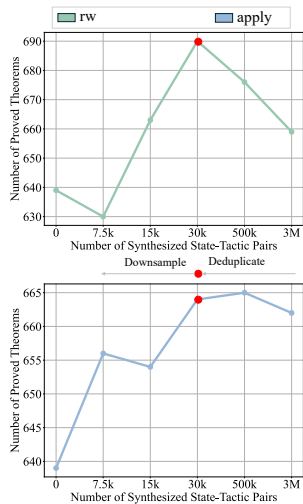
437 Methods	random	novel_premises	Search Budget
438 tidy	23.8	5.3	-
439 GPT-4	29.0	7.4	1 × 35
440 Reprover (Yang et al., 2023)	47.6	23.2	1 × 64
441 w/ retrieval	51.2	26.3	1 × 64
442 llmstep (Pythia 2.8b) (Welleck & Saha, 2023)	47.6	-	1 × 32
443	50.1	-	2 × 32
444 <b>Llama3-8b</b>	58.22	38.52	1 × 32
445 Mathlib-train + rw	59.62 (+1.40)	42.13 (+3.62)	1 × 32
446 Mathlib-train + apply	58.84 (+0.62)	41.29 (+2.77)	1 × 32
447 Mathlib-train + rw + apply	<b>59.82 (+1.60)</b>	<b>43.22 (+4.70)</b>	1 × 32
448			
449 <b>deepseek-coder-7b-base-v1.5</b>	57.7	39.24	1 × 32
450 Mathlib-train + rw	59.25 (+1.55)	42.98 (+3.74)	1 × 32
451 Mathlib-train + apply	58.68 (+0.98)	40.51 (+1.27)	1 × 32
452 Mathlib-train + rw + apply	<b>60.39 (+2.69)</b>	<b>43.46 (+4.22)</b>	1 × 32

453  
 454  
 455 tion. In conclusion, the results of the experiment show that simply mutating the seed theorems and  
 456 introducing state-tactic pairs of a single tactic can relieve the data scarcity problem and enhance the  
 457 theorem-proving ability of LLMs.

458  
 459 4.3.2 EFFECTIVENESS OF CONTINUAL PRETRAINING

460  
 461 Table 4: Effectiveness of continual pre-training. We grouped the  
 462 dataset for CPT and SFT by the tactic employed in the additional  
 463 state-tactic pairs.

464 Methods	random	novel_premises	random	novel_premises
	<b>Llama3-8b</b>		<b>deepseek-coder-base-7b-v1.5</b>	
	<i>sft: mathlib-train</i>			
467 w/o cpt	58.22	38.52	57.70	39.24
468 rw	59.56 (+1.35)	42.56 (+4.04)	58.74 (+1.04)	40.69 (+1.45)
469 apply	58.42 (+0.21)	41.29 (+2.77)	58.58 (+0.88)	40.02 (+0.78)
470 rw + apply	59.72 (+1.50)	42.19 (+3.68)	59.67 (+1.97)	41.65 (+2.41)
	<i>sft: mathlib-train + rw</i>			
471 w/o cpt	57.85	41.59	58.63	41.05
472 rw	59.62 (+1.76)	42.13 (+0.54)	59.25 (+0.62)	42.98 (+1.93)
	<i>sft: mathlib-train + apply</i>			
474 w/o cpt	56.71	40.02	57.96	41.17
475 apply	58.84 (+2.13)	41.29 (+1.27)	58.68 (+0.73)	40.51 (-0.66)
	<i>sft: mathlib-train + rw + apply</i>			
477 w/o cpt	58.53	41.95	58.37	42.92
478 rw + apply	59.82 (+1.30)	43.22 (+1.27)	60.39 (+2.02)	43.46 (+0.54)



479  
 480 Figure 4: Influence of the quantity  
 481 of synthesized data points.

482 To examine the necessity of continual pretraining, we assess and contrast the performance of the  
 483 LLM on *Leandajo* benchmark when the pretraining stage is included versus when it is excluded  
 484 as our baselines and present the results of pretraining on the augmented corpus. As shown in Table 4,  
 485 the continual pretraining stage demonstrates a positive influence on the performance of LLMs across  
 diverse supervised fine-tuning settings. The experimental results indicate that continual pretraining  
 before the supervised finetuning stage is also beneficial to the theorem-proving ability of the LLM.

#### 4.3.3 INFLUENCE OF THE QUANTITY OF SFT DATASET

We deduplicate the synthesized state-tactic pairs of each tactic by the invocable theorem (i.e., premise). Then we obtain about 30k data points for each tactic. To examine the influence of the quantity of the SFT dataset, we compare the performance of *Llama-3-8B*, trained on different quantities of additional data points, on novel\_premises split of *Leandajo* benchmark. As shown in Fig 4, the selected quantity (30k) achieves a relatively optimal compromise between the performance and overhead. The experimental results also reveal that enlarging the quantity of state-tactic pairs of a single tactic tends to lead to rapid saturation. We assume that the key to continually improving the theorem-proving ability lies in keeping the diversity of tactics during the process of scaling the synthetic data. More details are presented in Appendix D.3.4.

#### 4.3.4 ANALYSIS OF OUT-OF-DISTRIBUTION PERFORMANCE

We evaluate *Llama-3-8b* using the competition-level theorem proving benchmark miniF2F. As shown in Table 5, our synthesized data still helps to improve the theorem-proving ability of LLMs on the out-of-distribution benchmark. The magnitude of this improvement is comparatively smaller than that observed on the in-distribution benchmark. We attribute this discrepancy to the divergence between synthesized tactics and the preferred tactics to prove competition-level problems. Through manual inspection of the correct proofs generated by various LLMs trained on Mathlib-train, we identify a tendency to favor advanced and automated tactics (e.g., *simp*, *omega*, *linarith*, *norm\_num*, etc.). Additionally, we analyze the distribution of tactics used in proved theorems across different data compositions and make the following observations. 1) Data augmentation on a single tactic will increase the model’s preference for the specific tactic. 2) Adjusting the distribution of different tactics within the dataset is promising to improve the theorem-proving ability of LLMs. The entire analysis process is illustrated in Appendix E.2.

Table 5: Results on miniF2F. We evaluate the performance across different data compositions and list the ratio of *rw*, *apply*, *norm\_num* and *linarith* used by Llama3-8b to prove these theorems.

Methods	miniF2F-test	Correct/Total	<i>rw</i>	<i>apply</i>	<i>norm_num</i>	<i>linarith</i>
Mathlib-train	34.01	83/244	16.10	0.00	27.12	16.95
Mathlib-train + rw	35.24	86/244	18.75	0.78	14.84	21.88
Mathlib-train + apply	36.07	88/244	8.87	2.42	20.16	15.63
Mathlib-train + rw + apply	<b>36.48</b> (+2.47)	89/244	12.31	0.77	26.92	16.92

## 5 CONCLUSION

We have presented a general data-synthesis framework for the Lean theorem prover, which amplifies the theorem-proving capability of the LLM through symbolic mutation. Our algorithm increases the number of theorems in Mathlib by an order of magnitude and achieves promising results in improving the theorem-proving ability of the LLM. We discuss the limitations of our method in Appendix B. Synthesizing formal theorems is an inherently challenging problem. Our approach, much like ancient alchemy, involves experimenting with a substantial number of theorems in the hope of uncovering valuable “gold”. We aspire for our algorithm and data to serve as a foundation for further research, advancing theorem synthesis from alchemy to chemistry.

## REFERENCES

- Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster, Marco Dos Santos, Stephen McAleer, Albert Q. Jiang, Jia Deng, Stella Biderman, and Sean Welleck. Llemma: An open language model for mathematics. *CoRR*, abs/2310.10631, 2023. doi: 10.48550/ARXIV.2310.10631. URL <https://doi.org/10.48550/arXiv.2310.10631>.
- Bruno Barras, Samuel Boutin, Cristina Cornes, Judicaël Courant, Jean-Christophe Filliatre, Eduardo Gimenez, Hugo Herbelin, Gerard Huet, Cesar Munoz, Chetan Murthy, et al. *The Coq proof assistant reference manual: Version 6.1*. PhD thesis, Inria, 1997.

- 540 Sascha Böhme and Tobias Nipkow. Sledgehammer: Judgement day. In Jürgen Giesl and Reiner  
541 Hähnle (eds.), *Automated Reasoning, 5th International Joint Conference, IJCAR 2010, Edin-*  
542 *burgh, UK, July 16-19, 2010. Proceedings*, volume 6173 of *Lecture Notes in Computer Sci-*  
543 *ence*, pp. 107–121. Springer, 2010. doi: 10.1007/978-3-642-14203-1\_9. URL [https://doi.org/10.1007/978-3-642-14203-1\\_9](https://doi.org/10.1007/978-3-642-14203-1_9).
- 544  
545 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
546 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
547 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 548  
549 Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. *CoRR*,  
550 abs/2307.08691, 2023. doi: 10.48550/ARXIV.2307.08691. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2307.08691)  
551 [48550/arXiv.2307.08691](https://doi.org/10.48550/arXiv.2307.08691).
- 552  
553 Leonardo Mendonça de Moura, Soonho Kong, Jeremy Avigad, Floris van Doorn, and Jakob von  
554 Raumer. The lean theorem prover (system description). In Amy P. Felty and Aart Middeld-  
555 dorp (eds.), *Automated Deduction - CADE-25 - 25th International Conference on Automated*  
556 *Deduction, Berlin, Germany, August 1-7, 2015, Proceedings*, volume 9195 of *Lecture Notes in*  
557 *Computer Science*, pp. 378–388. Springer, 2015. doi: 10.1007/978-3-319-21401-6\_26. URL  
[https://doi.org/10.1007/978-3-319-21401-6\\_26](https://doi.org/10.1007/978-3-319-21401-6_26).
- 558  
559 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
560 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
561 *arXiv preprint arXiv:2407.21783*, 2024.
- 562  
563 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao  
564 Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. Deepseek-coder: When the  
565 large language model meets programming - the rise of code intelligence. *CoRR*, abs/2401.14196,  
566 2024. doi: 10.48550/ARXIV.2401.14196. URL [https://doi.org/10.48550/arXiv.](https://doi.org/10.48550/arXiv.2401.14196)  
[2401.14196](https://doi.org/10.48550/arXiv.2401.14196).
- 567  
568 Jesse Michael Han, Jason Rute, Yuhuai Wu, Edward W. Ayers, and Stanislas Polu. Proof artifact  
569 co-training for theorem proving with language models. In *The Tenth International Conference on*  
570 *Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022.  
571 URL <https://openreview.net/forum?id=rpXJc9j04U>.
- 572  
573 Yinya Huang, Xiaohan Lin, Zhengying Liu, Qingxing Cao, Huajian Xin, Haiming Wang, Zhenguo  
574 Li, Linqi Song, and Xiaodan Liang. MUSTARD: mastering uniform synthesis of theorem and  
575 proof data. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vi-*  
576 *enna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL [https://openreview.net/](https://openreview.net/forum?id=8xliOUg9EW)  
[forum?id=8xliOUg9EW](https://openreview.net/forum?id=8xliOUg9EW).
- 577  
578 Albert Qiaochu Jiang, Sean Welleck, Jin Peng Zhou, Timothée Lacroix, Jiacheng Liu, Wenda Li,  
579 Mateja Jamnik, Guillaume Lample, and Yuhuai Wu. Draft, sketch, and prove: Guiding formal  
580 theorem provers with informal proofs. In *The Eleventh International Conference on Learn-*  
581 *ing Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL  
<https://openreview.net/forum?id=SMa9EAovKMC>.
- 582  
583 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph  
584 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model  
585 serving with pagedattention. In Jason Flinn, Margo I. Seltzer, Peter Druschel, Antoine Kaufmann,  
586 and Jonathan Mace (eds.), *Proceedings of the 29th Symposium on Operating Systems Principles,*  
587 *SOSP 2023, Koblenz, Germany, October 23-26, 2023*, pp. 611–626. ACM, 2023. doi: 10.1145/  
3600006.3613165. URL <https://doi.org/10.1145/3600006.3613165>.
- 588  
589 Guillaume Lample, Timothée Lacroix, Marie-Anne Lachaux, Aurélien Rodriguez, Amaury Hayat,  
590 Thibaut Lavril, Gabriel Ebner, and Xavier Martinet. Hypertree proof search for neural theorem  
591 proving. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh  
592 (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural*  
593 *Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - De-*  
*cember 9, 2022, 2022*. URL [http://papers.nips.cc/paper\\_files/paper/2022/](http://papers.nips.cc/paper_files/paper/2022/hash/a8901c5e85fb8e1823bbf0f755053672-Abstract-Conference.html)  
[hash/a8901c5e85fb8e1823bbf0f755053672-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/a8901c5e85fb8e1823bbf0f755053672-Abstract-Conference.html).

- 594 Zhaoyu Li, Jialiang Sun, Logan Murphy, Qidong Su, Zenan Li, Xian Zhang, Kaiyu Yang, and Xujie  
595 Si. A survey on deep learning for theorem proving. *CoRR*, abs/2404.09939, 2024. doi: 10.48550/  
596 ARXIV.2404.09939. URL <https://doi.org/10.48550/arXiv.2404.09939>.  
597
- 598 Norman Megill and David A Wheeler. *Metamath: a computer language for mathematical proofs*.  
599 Lulu. com, 2019.
- 600 Maciej Mikula, Szymon Antoniak, Szymon Tworowski, Albert Qiaochu Jiang, Jin Peng Zhou,  
601 Christian Szegedy, Lukasz Kucinski, Piotr Milos, and Yuhuai Wu. Magnushammer: A  
602 transformer-based approach to premise selection. *CoRR*, abs/2303.04488, 2023. doi: 10.48550/  
603 ARXIV.2303.04488. URL <https://doi.org/10.48550/arXiv.2303.04488>.  
604
- 605 Lawrence C. Paulson. *Isabelle - A Generic Theorem Prover (with a contribution by T. Nipkow)*,  
606 volume 828 of *Lecture Notes in Computer Science*. Springer, 1994. ISBN 3-540-58244-4. doi:  
607 10.1007/BFB0030541. URL <https://doi.org/10.1007/BFB0030541>.  
608
- 609 Stanislas Polu and Ilya Sutskever. Generative language modeling for automated theorem proving.  
610 *CoRR*, abs/2009.03393, 2020. URL <https://arxiv.org/abs/2009.03393>.
- 611 Stanislas Polu, Jesse Michael Han, Kunhao Zheng, Mantas Baksys, Igor Babuschkin, and Ilya  
612 Sutskever. Formal mathematics statement curriculum learning. In *The Eleventh International  
613 Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenRe-  
614 view.net, 2023. URL <https://openreview.net/forum?id=-P7G-8dmSh4>.  
615
- 616 Samyam Rajbhandari, Olatunji Ruwase, Jeff Rasley, Shaden Smith, and Yuxiong He. Zero-infinity:  
617 Breaking the GPU memory wall for extreme scale deep learning. *CoRR*, abs/2104.07857, 2021.  
618 URL <https://arxiv.org/abs/2104.07857>.
- 619 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-  
620 baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gem-  
621 ini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint  
622 arXiv:2403.05530*, 2024.  
623
- 624 Peiyang Song, Kaiyu Yang, and Anima Anandkumar. Towards large language models as copilots  
625 for theorem proving in lean. *CoRR*, abs/2404.12534, 2024. doi: 10.48550/ARXIV.2404.12534.  
626 URL <https://doi.org/10.48550/arXiv.2404.12534>.  
627
- 628 Terence Tao, 2023. URL <https://teorth.github.io/pfr/blueprint.pdf>.
- 629 Terence Tao. Machine assisted proof. *Notices of the American Mathematical Society, to appear*,  
630 2024.  
631
- 632 Trieu H. Trinh, Yuhuai Wu, Quoc V. Le, He He, and Thang Luong. Solving olympiad geometry with-  
633 out human demonstrations. *Nat.*, 625(7995):476–482, 2024. doi: 10.1038/S41586-023-06747-5.  
634 URL <https://doi.org/10.1038/s41586-023-06747-5>.  
635
- 636 Haiming Wang, Ye Yuan, Zhengying Liu, Jianhao Shen, Yichun Yin, Jing Xiong, Enze Xie,  
637 Han Shi, Yujun Li, Lin Li, Jian Yin, Zhenguo Li, and Xiaodan Liang. Dt-solver: Automated  
638 theorem proving with dynamic-tree sampling guided by proof-level value function. In Anna  
639 Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual  
640 Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023,  
641 Toronto, Canada, July 9-14, 2023*, pp. 12632–12646. Association for Computational Linguistics,  
642 2023. doi: 10.18653/V1/2023.ACL-LONG.706. URL [https://doi.org/10.18653/v1/  
643 2023.acl-long.706](https://doi.org/10.18653/v1/2023.acl-long.706).
- 644 Haiming Wang, Huajian Xin, Chuanyang Zheng, Zhengying Liu, Qingxing Cao, Yinya Huang, Jing  
645 Xiong, Han Shi, Enze Xie, Jian Yin, Zhenguo Li, and Xiaodan Liang. Lego-prover: Neural  
646 theorem proving with growing libraries. In *The Twelfth International Conference on Learning  
647 Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL  
<https://openreview.net/forum?id=3f5PALef5B>.

- 648 Mingzhe Wang and Jia Deng. Learning to prove theorems by learning to generate theo-  
649 rems. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and  
650 Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Con-  
651 ference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12,  
652 2020, virtual*, 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/  
653 d2a27e83d429f0dcae6b937cf440aeb1-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/d2a27e83d429f0dcae6b937cf440aeb1-Abstract.html).
- 654 Sean Welleck. Neural theorem proving tutorial ii. [https://github.com/cmu-13/  
655 ntptutorial-II](https://github.com/cmu-13/ntptutorial-II), 2023.
- 656 Sean Welleck and Rahul Saha. LLMSTEP: LLM proofstep suggestions in lean. *CoRR*,  
657 abs/2310.18457, 2023. doi: 10.48550/ARXIV.2310.18457. URL [https://doi.org/10.  
658 48550/arXiv.2310.18457](https://doi.org/10.48550/arXiv.2310.18457).
- 660 Yuhuai Wu, Albert Q. Jiang, Jimmy Ba, and Roger Baker Grosse. INT: an inequality benchmark  
661 for evaluating generalization in theorem proving. In *9th International Conference on Learning  
662 Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL  
663 <https://openreview.net/forum?id=O6LPudowNQm>.
- 664 Yuhuai Wu, Albert Qiaochu Jiang, Wenda Li, Markus N. Rabe, Charles Staats, Mateja Jamnik,  
665 and Christian Szegedy. Autoformalization with large language models. In Sanmi Koyejo,  
666 S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neu-  
667 ral Information Processing Systems 35: Annual Conference on Neural Information Process-  
668 ing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9,  
669 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/  
670 d0c6bc641a56bebee9d985b937307367-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/d0c6bc641a56bebee9d985b937307367-Abstract-Conference.html).
- 671 Zijian Wu, Jiayu Wang, Dahua Lin, and Kai Chen. Lean-github: Compiling github LEAN reposi-  
672 tories for a versatile LEAN prover. *CoRR*, abs/2407.17227, 2024. doi: 10.48550/ARXIV.2407.  
673 17227. URL <https://doi.org/10.48550/arXiv.2407.17227>.
- 674 Huajian Xin, Z. Z. Ren, Junxiao Song, Zhihong Shao, Wanxia Zhao, Haocheng Wang, Bo Liu,  
675 Liyue Zhang, Xuan Lu, Qiushi Du, Wenjun Gao, Qihao Zhu, Dejian Yang, Zhibin Gou, Z. F.  
676 Wu, Fuli Luo, and Chong Ruan. Deepseek-prover-v1.5: Harnessing proof assistant feedback  
677 for reinforcement learning and monte-carlo tree search. *CoRR*, abs/2408.08152, 2024. doi: 10.  
678 48550/ARXIV.2408.08152. URL <https://doi.org/10.48550/arXiv.2408.08152>.
- 679 Kaiyu Yang, Aidan M. Swope, Alex Gu, Rahul Chalamala, Peiyang Song, Shixing Yu,  
680 Saad Godil, Ryan J. Prenger, and Animashree Anandkumar. Leandojo: Theorem prov-  
681 ing with retrieval-augmented language models. In Alice Oh, Tristan Naumann, Amir  
682 Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neu-  
683 ral Information Processing Systems 36: Annual Conference on Neural Information  
684 Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16,  
685 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/  
686 hash/4441469427094f8873d0fecb0c4e1cee-Abstract-Datasets\\_and\\_  
687 Benchmarks.html](http://papers.nips.cc/paper_files/paper/2023/hash/4441469427094f8873d0fecb0c4e1cee-Abstract-Datasets_and_).
- 688 Huaiyuan Ying, Zijian Wu, Yihan Geng, Jiayu Wang, Dahua Lin, and Kai Chen. Lean workbook:  
689 A large-scale lean problem set formalized from natural language math problems. *arXiv preprint  
690 arXiv:2406.03847*, 2024.
- 691 Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. minif2f: a cross-system benchmark  
692 for formal olympiad-level mathematics. In *The Tenth International Conference on Learning  
693 Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL  
694 <https://openreview.net/forum?id=9ZPegFuFTFv>.
- 695 Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyuan Luo, and Yongqiang Ma. Lla-  
696 mafactory: Unified efficient fine-tuning of 100+ language models. *CoRR*, abs/2403.13372, 2024.  
697 doi: 10.48550/ARXIV.2403.13372. URL [https://doi.org/10.48550/arXiv.2403.  
698 13372](https://doi.org/10.48550/arXiv.2403.13372).
- 699  
700  
701

702	CONTENTS	
703		
704	<b>A Background on Lean</b>	<b>15</b>
705		
706	<b>B Limitations</b>	<b>15</b>
707		
708		
709	<b>C Detailed Information of Synthesizing Algorithms</b>	<b>16</b>
710	C.1 Overview . . . . .	16
711	C.2 Find Invocable Theorems . . . . .	16
712	C.3 Construct New Theorems . . . . .	17
713	C.3.1 <i>rw</i> tactic . . . . .	17
714	C.3.2 <i>apply</i> tactic . . . . .	19
715	C.4 Verify the Theorems . . . . .	20
716	C.5 Limitations of Synthesis Pipeline . . . . .	21
717		
718		
719		
720		
721	<b>D Deeper Analysis of Synthetic Dataset</b>	<b>22</b>
722	D.1 Numerical Analysis . . . . .	22
723	D.2 Examples . . . . .	22
724	D.3 Details of Training Data . . . . .	22
725	D.3.1 Examples of Training Data . . . . .	22
726	D.3.2 Preprocessing . . . . .	23
727	D.3.3 Classification of Extracted Tactics . . . . .	23
728	D.3.4 Influence of the Quantity of SFT Dataset . . . . .	23
729		
730		
731		
732		
733	<b>E Additional Experiments</b>	<b>26</b>
734	E.1 Effectiveness of Different Tactics . . . . .	26
735	E.2 Analysis of the Tactics to Prove miniF2F Theorems . . . . .	26
736	E.2.1 Preference in Used Tactics . . . . .	26
737	E.2.2 Influence of Additional Tactics . . . . .	26
738		
739		
740		
741		
742		
743		
744		
745		
746		
747		
748		
749		
750		
751		
752		
753		
754		
755		

## A BACKGROUND ON LEAN

Lean is a functional programming language and interactive theorem prover based on dependent type theory. As one of the most popular formal systems, Lean aids mathematicians in formalizing statements and proofs in a semi-auto style and enables them to verify the correctness of each proof step through rigorous type-checking.

**Theorem in Lean** To some extent, theorems in Lean can be seen as a special variant of functions in general-purpose programming languages. A theorem consists of a statement and corresponding proof. In Lean, the keyword “theorem”, “example” or “lemma” is used to define the “function”, sometimes followed by a specific function name. The assumption of a statement can be formatted as implicit or explicit arguments, while the assertion of the statement specifies the return type of the function. The proof of the statement can be viewed as the function body, which constructs a proof term with the type specified by the assertion. There are two main proof styles in Lean: term-style and tactic-style. In term-style proofs, theorems are proven using constructive methods. On the other hand, tactic-style proofs sequentially decompose the proof goal using specific tactics. Although tactic-style proofs are less readable, they tend to have shorter proof lengths. Most machine learning-based theorem-proving systems focus on tactic-style proof. The synthesis method proposed by our paper can be applied to both styles.

**Tactic** Lean offers various advanced tactics for theorem proving, which set it apart from other formal systems (e.g., Coq, Isabelle). In handwritten proofs, authors tend to guide the reader on building the proof through instructions such as “apply the previous lemma”, “invoke the principle of mathematical induction”, or “simplify the expression”. Similarly, tactics in Lean are used to describe how to construct a proof term incrementally. They help users decompose the proof goal step by step, allowing users to focus on only one proof goal at a time.

**Mathlib** Mathlib<sup>7</sup> is a comprehensive mathematical library for Lean, largely maintained by the community, which encompasses a broad spectrum of mathematical subjects (e.g., algebra and analysis) and consists of over 120,000 theorems along with their respective axioms and definitions. This extensive knowledge base serves as the primary corpus for neural theorem provers.

## B LIMITATIONS

Our method exhibits some limitations that remain to be addressed in future endeavors.

**Data Diversity and Quality** We only define two symbolic rules (using two tactics) to synthesize new theorems. The implementation of the synthesis pipeline is over general and utilizes little domain knowledge, which affects the diversity and quality of synthetic data.

**The Cost of Synthesizing** Despite the CPU-only nature of our algorithm, the cost of synthesizing remains huge. We believe the overhead can be significantly reduced with a finer implementation and more specialized tools to interact with the Lean theorem prover.

**Single-Round v.s. Multi-Round** Theoretically speaking, our algorithms can be iteratively executed by adding the synthesized theorems into seed theorems. Conversely, the synthesized repository is very heavy, which makes it hard to interact with Lean using *Leandajo* and deploy our algorithm on existing hardware.

**Theorem-level or Term-level** Our method synthesizes theorems from top to bottom and introduces additional state-tactic pairs of specific tactics. Synthesizing formal data at the theorem level is not efficient and not consistent with the step-by-step nature of theorem-proving. Ideally, we anticipate that we can synthesize formal data directly at the term level, which aligns with the characteristics of interactive theorem proving.

<sup>7</sup><https://github.com/leanprover-community/mathlib4>

**Up-to-down v.s. Down-to-up** We synthesize theorems in an up-to-down fashion. We construct the new statements first and then retrieve the correct proofs. The up-to-down fashion depends on a specific set of seed theorems, which restricts the diversity of synthetic data. A more fundamental idea is that we can sample some terms in the symbolic space directly, merge them using symbolic manipulations, and then find the corresponding goals for this new theorem. This *AlphaGeometry*-style idea is hard to implement in Lean and requires a large amount of domain knowledge and engineering endeavors.

**Symbolic Synthesis in Conjunction with Other Techniques** Our proposed method demonstrates significant potential for integration with other techniques to enhance the theorem-proving capabilities of LLMs. We posit that theorem synthesis in the symbolic space serves as a valuable complement to prevailing auto-formalization methods. For instance, it may contribute to the expansion of autoformalized datasets. Besides, our approach generates a substantial quantity of new proven statements which can be utilized as a comprehensive database for Retrieval-Augmented Generation (RAG) (Yang et al., 2023; Wang et al., 2024). Our objective is to amalgamate these methodologies to develop a robust theorem prover in the future.

## C DETAILED INFORMATION OF SYNTHESIZING ALGORITHMS

### C.1 OVERVIEW

As discussed in Section 3, the entire algorithm is composed of four steps. 1) Find invocable theorems for the candidate theorem by executing a specific tactic and retrieving the resulting proof state. 2) Construct new statements, where we parse the resulting proof state and mutate the old statement with the help of AST. 3) Establish the entire proof by inserting a *have* tactic and integrating it with the old proof to build the whole proof for this new statement. 4) Verify the correctness of generated theorems in Lean theorem prover. In practice, we separately run the time-consuming first step on hundreds of 8-core CPU nodes and unify step 2) and step 3) together to construct the new theorem. Then we will write back synthetic theorems and run “lake build” to verify the generated theorems.

### C.2 FIND INVOCABLE THEOREMS

For each candidate theorem, we check whether other theorems can be used to rewrite or apply to it by executing tactics. We use the *run\_tac* API provided by *Leandajo* to run a specific tactic and extract the valid proof state according to predefined criteria. The instruction templates for each tactic are listed in Table1. Here is the code snippet that illustrates this process.

```

845 1 '''args:
846 2     dojo: interactive environment
847 3     init_state: initial proof state of target theorem
848 4     theorem: a possible invocable theorem
849 5     hypos: the assumptions of the target theorem (extracted by parsing
850 6     the AST)
851 7 '''
852 8 def is_invocable_theorem(
853 9     dojo, init_state, theorem, hypos, mode="rw"
854 10 ):
855 11     name = theorem.full_name
856 12     if mode == "rw":
857 13         # e.g. rw [name] at hypo_name
858 14         insts = get_rw_insts(name, hypos)
859 15     elif mode == "apply":
860 16         # e.g. have hypo_str := by apply name
861 17         insts = get_apply_insts(name, hypos)
862 18     res = []
863 19     for i, inst in enumerate(insts):
864 20         try: next_state = dojo.run_tac(init_state, inst)
865 21         except Exception as e: ...
866 22         else:
867 23             state_info = {
868 24                 "init_state": init_state.pp, # pp means pretty-printed

```



```

864 24         "next_state": next_state.error if isinstance(next_state,
865 LeanError) else next_state.pp,
866 25         "rule": inst
867 26     }
868 27         if isinstance(next_state, LeanError):
869 28             if mode == "implication" \
870 29                 and "unsolved goals" in next_state.error:
871 30                 res.append(state_info)
872 31             elif isinstance(next_state, TacticState):
873 32                 res.append(state_info)
874 33     return res

```

Listing 1: Find invocable theorems by running tactics.

We set different validation criteria for each tactic. For the *rw* tactic, if the resulting state is a *TacticState*, we annotate this theorem as invocable. In contrast, for the *apply* tactic, the resulting state should be “unsolved goals”. Additionally, we filter the resulting invocable theorems to simplify the problem of constructing new theorems. Specifically, we remove the invocable theorems whose *next\_state* contains meta-variables (e.g., *?a*, *?m123*) for the *rw* tactic and unnamed meta-variables (e.g., *?e12384*) for the *apply* tactic. Ultimately, we retrieve the invocable theorems for each candidate theorem. One example of invocable theorems is shown in Fig 5.

The experiments run on a large collection of CPUs ( $512 \times 8$ -core for the *rw* tactic and  $256 \times 8$ -core for *apply*). The substantial CPU requirement is largely due to the memory-intensive nature of *Leandajo*, which hinders multiprocessing on a single node. We anticipate a significant reduction in the cost of our experiments by implementing a lighter interface for Lean interaction. The operation of *apply* is more complex and time-consuming than *rw*. We set a one-hour timeout for each *dojo* environment to reduce the time cost. When running a specific tactic, we do not add additional imports to the *dojo* environment to avoid introducing human preferences in the process of synthesis. This setting may narrow the scope of theorems that the tactic can access and lower the variety of invocable theorems.

In summary, finding invocable theorems constitutes the most time-consuming and computationally intensive stage of our algorithm, entailing trade-offs among cost, time, and generated data volume.

### C.3 CONSTRUCT NEW THEOREMS

To create a new theorem, we construct the new statement using the invocable theorems returned by Section C.2 and then establish the entire proof through *have* tactic. Our symbolic engine is built upon *Leandajo* API, utilizing the extracted AST and some string manipulations. To facilitate the detailed explanation of algorithms, we will delineate the implementation of these two tactics separately in the following pseudocode or source code.

#### C.3.1 *rw* TACTIC

The logic of constructing a new statement for *rw* tactic is simple. We just identify whether a specific assumption or assertion has been rewritten by parsing invocable instructions with regular expressions. Then we parse the AST node of the candidate statement to locate the corresponding part that should be mutated. Finally, we extract the new assumption or assertion from the next proof state and replace the old one with the new one. The main procedure is shown in Algorithm 2.

---

#### Algorithm 2 Construct new statement for *rw* tactic

---

```

910 Input: candidate statement  $s$ , invocable theorem  $i$ 
911 Output: mutated statement  $s_m$ 
912  $node \leftarrow \text{EXTRACT\_AST}(s)$  ▷ extract the AST of candidate statement
913  $-, next\_state, inst \leftarrow i$  ▷ get the next state and instruction
914  $flag \leftarrow \text{IDENTIFY}(i)$  ▷ flag specifies whether the assumption or assertion should be mutated
915  $location\ l \leftarrow \text{PARSE}(node, i, flag)$  ▷ parse AST node and locate the corresponding part that
916  $m \leftarrow \text{CONSTRUCT}(next\_state)$  ▷ parse the next proof state and construct the target string
917  $new\ statement\ s_m \leftarrow \text{REPLACE}(s, m, l)$ 

```

---

```

918
919 theorem_name: Char.ofNat_toNat
920 rule: have h : isValidCharNat c.toNat := by apply List.rel_of_pairwise_cons
921 init_state:
922   c : Char
923   h : isValidCharNat (toNat c)
924   ⊢ ofNat (toNat c) = c
925 next_state:
926   unsolved goals
927   case hp
928     c : Char
929     h : isValidCharNat (toNat c)
930     ⊢ Std.RBNode.All isValidCharNat ?t
931   case H
932     c : Char
933     h : isValidCharNat (toNat c)
934     ⊢ ∀ {x : ℕ}, x ∈ ?lb → isValidCharNat x
935   case a
936     c : Char
937     h : isValidCharNat (toNat c)
938     ⊢ Std.RBNode.lowerBound? ?cut ?t ?lb = some (toNat c)
939   case lb
940     c : Char
941     h : isValidCharNat (toNat c)
942     ⊢ Option ℕ
943   case cut
944     c : Char
945     h : isValidCharNat (toNat c)
946     ⊢ ℕ → Ordering
947   case t
948     c : Char
949     h : isValidCharNat (toNat c)
950     ⊢ Std.RBNode ℕ
951
952
953
954
955
956
957
958
959
960
961
962

```

Figure 5: Examples of invocable theorems for *apply*

After creating a new statement, we should insert a *have* tactic to construct the whole proof. If the assumption is modified, then we just restore it to the old one by reversing the direction of *rw* within a *have* instruction and then concatenate it with the original proof. If the assertion is mutated, the *have* tactic can be used to prove the original assertion with initial proof block. Then we just rewrite the old proof goal to the new one to construct the whole proof. Here is a simplified code snippet that illustrates this process.

```

963 1 def proof_generation_rw(
964   2     invocable_inst,
965   3     flag,
966   4     proof_str,
967   5     conc_or_hypo_old=None,
968   6     is_tactic_style=False
969   7   ):
970   8     inst = invocable_inst["rule"]
971   9     if flag == "hypo":
972   10       hypo_name = parse(inst, flag)
973   11       # find the delimiter for proof str(e.g. := by or :=) (simplified
974   version)

```

```

972 12   if is_tactic_style:
973 13       delimiter = " := by"
974 14   else:
975 15       delimiter = " :="
976 16   splits = proof_str.split(delimiter)
977 17   proof_seqs = delimiter.join(splits[1:])
978 18   if flag == "hypo":
979 19       rev_inst = reverse_rw(invocable_inst)
980 20       have_template = "have {subgoal} := by {proof_seqs}"
981 21       have_inst = have_template.format(
982 22           subgoal=conc_or_hypo_old,
983 23           proof_seqs=rev_inst)
984 24       have_inst += f';exact {hypo_name}'
985 25       end_inst = proof_seqs
986 26   elif flag == "conclusion":
987 27       have_template = "have : {subgoal} {delimiter} {proof_seqs}"
988 28       have_inst = have_template.format(
989 29           subgoal=conc_or_hypo_old,
990 30           delimiter=delimiter,
991 31           proof_seqs=proof_seqs)
992 32       head = "by " if not is_tactic_style else ""
993 33       _suffix = " at this;exact this"
994 34       end_inst = head + inst + _suffix
995 35   # do indentation
996 36   have_inst = indent_code(delimiter, proof_str, have_inst, indent_level
997 37   =...)
998 38   end_inst = indent_code(delimiter, proof_str, end_inst, indent_level
999 39   =...)
1000 40   # concat the different parts of proof
1001 41   prefix = splits[0] + delimiter + '\n'
1002 42   suffix = end_inst if end_inst.startswith('\n') else '\n' + end_inst
1003 43   new_proof = prefix + have_inst + suffix
1004 44   return new_proof

```

Listing 2: Build the whole proof for *rw* tactic

### C.3.2 *apply* TACTIC

---

**Algorithm 3** Construct new statement for *apply* tactic

---

```

1006 Input: candidate statement  $s$ , invocable instruction  $i$ 
1007 Output: mutated statement  $s_m$ 
1008  $H \leftarrow \emptyset$  ▷ initialize the set of new assumptions
1009  $node \leftarrow \text{EXTRACT\_AST}(s)$  ▷ extract the AST of candidate statement
1010  $-, next\_state, inst \leftarrow i$  ▷ get the next state and instruction
1011  $M, G \leftarrow \text{PARSE}(next\_state)$  ▷ get the set of metavaribales and other subgoals
1012 for  $m \in M$  do ▷ Assigning metavariables
1013     Add  $\text{ASSIGN}(m, next\_state)$  to  $H$ 
1014 end for
1015 for  $g \in G$  do ▷ Fill the other subgoals depending on meta-varibales
1016     Add  $\text{ASSIGN}(g, next\_state, M)$  to  $H$ 
1017 end for
1018  $H \leftarrow \text{HANDLE\_NAMING\_CONFLICTS}(H)$ 
1019 new assumption  $h_m \leftarrow \text{CONCAT}(H)$ 
1020 location  $l \leftarrow \text{PARSE}(node, i)$  ▷ parse AST node and locate the old assumption that needs to be mutated
1021  $s_m \leftarrow \text{REPLACE}(s, h_m, l)$ 

```

---

Constructing new statements for *apply* tactic is more complex than *rw*. Applying a theorem may introduce some metavariables and new subgoals into the local context for the resulting proof state as shown in Fig 5. We assign values to the metavariables by parsing the next\_state and then retrieve all subgoals containing metavariables as new assumptions. For each new assumption, we can extract

its name and type from the proof state. To avoid naming conflicts, we define a set of rules to rename the variable according to the naming conversion of Mathlib<sup>8</sup>. Ultimately, we concatenate all new assumptions and replace the old assumption with them. This procedure is shown in Algorithm 3.

Similarly, we can construct the entire proof for the new statement by inserting a *have* lemma. The simplified code snippet illustrates this process.

```

1032 1 def proof_generation_apply(cases_goals, inst, proof_str, is_tactic_style)
1033   :
1034 2   if len(cases_goals) == 1:
1035 3       lemma = inst + "; assumption"
1036 4   elif len(cases_goals) > 1:
1037 5       lemma = inst + "<> assumption"
1038 6   else:
1039 7       raise Exception("no available case and corresponding goal")
1040 8
1041 9   if is_tactic_style:
1042 10      delimiter = " := by"
1043 11   else:
1044 12      delimiter = " :="
1045 13
1046 14   splits = proof_str.split(delimiter)
1047 15   proof_seqs = delimiter.join(splits[1:])
1048 16   lemma = indent_code(delimiter, proof_str, lemma, indent_level=...)
1049 17   prefix = splits[0] + delimiter + '\n'
1050 18   suffix = proof_seqs if proof_seqs.startswith('\n') else '\n' +
1051 19   proof_seqs
1052 20   new_proof = prefix + lemma + suffix
1053   return new_proof

```

Listing 3: Build the whole proof for *apply* tactic

#### C.4 VERIFY THE THEOREMS

Our method creates a set of variants for each candidate theorem in Mathlib4. We write the variants back to the original file and execute *lake build* for verification. We remove the wrong lines for each file by parsing the error message returned by Lean. Then, we will rebuild the repo to ensure the effectiveness of verification. We remove the files that cause errors in the rebuilding process. Specifically, for each 8-core CPU node, we only build one “.lean” file each time to speed up this process and simplify the logic of parsing. The whole experiment runs on 2,048 CPUs (256 × 8-core). The code snippets illustrate the procedure for each CPU node. After verifying the correctness of the synthesized theorem, we extract the state-tactic pairs from our augmented Mathlib repository using *Leandojo*. For *rw* or *apply*, it takes three days for a 96-core CPU machine to trace the enlarged repository. In practice, we split the modified lean files into several portions, separately write them into multiple lean repositories, and trace the repos on several 96-core CPU machines.

```

1067 1 # A single 8-core CPU node
1068 2 res = []
1069 3 for idx, file in enumerate(files):           # for each modified file
1070 4     '''file {
1071 5         file_name: "name of the lean file",
1072 6         text: "the content of this file after writing synthesized
1073 7         variants into this file"
1074 8         "loc": {"theorem_name": [(start_line_nb, end_line_nb)...]}
1075 9     }'''
1076 9     tmp = {
1077 10         'loc': file['loc'],
1078 11         'file_name': file['file_name'],
1079 12         'text': file['text']
1079 13     }

```

<sup>8</sup><https://leanprover-community.github.io/contribute/naming.html>

```

1080 14 file_name = file['file_name']
1081 15 file_path = os.path.join(mathlib_package_path, file_name)
1082 16 # extract the old content of this file
1083 17 with open(file_path, "r") as f:
1084 18     old_str = f.read()
1085 19 # replace the old content with new content
1086 20 with open(file_path, "w") as f:
1087 21     f.write(file['text'])
1088 22 # change the build target to current file
1089 23 with open(LIBRARY_ROOT_FILE, 'w') as f: # LIBRARY_ROOT_PATH:
1090 24     Mathlib.lean
1091 25     module_name = file_name.replace('/', '.').replace('.lean', '')
1092 26     f.write(f"import {module_name}")
1093 27     if have_variants(file):
1094 28         ## lake build the new mathlib project
1095 29         wd = os.getcwd()
1096 30         result = lake_build(mathlib_package_path) #a helper function
1097 31         os.chdir(wd)
1098 32         ## parse the output
1099 33         # subprocess error
1100 34         if result == None: tmp['valid_loc'] = ["No variants"]
1101 35         elif result == 0:
1102 36             tmp['valid_loc'] = tmp['loc']
1103 37             print('successful build')
1104 38         # timeout error
1105 39         elif result == -1: tmp['valid_loc'] = ["No variants"]
1106 40         else:
1107 41             # find the error locations(line numbers)
1108 42             pattern = fr"({file_name}):(\d+):(\d+): error:"
1109 43             errors = re.findall(pattern, result)
1110 44             if len(errors) == 0: tmp['valid_loc'] = ["No variants"] #
1111 45         parse exception
1112 46         else:
1113 47             # extract line numbers from errors
1114 48             error_line_nbs = ...
1115 49             # get the locations of all variants
1116 50             intervals = ...
1117 51             # drop the error ones and write back
1118 52             valid_locs = diff(intervals, error_line_nbs)
1119 53             write_back(valid_locs, file['text'])
1120 54             ## rebuilt the project if causes error then remove this
1121 55         file
1122 56             wd = os.getcwd()
1123 57             result = lake_build(mathlib_package_path)
1124 58             os.chdir(wd)
1125 59             if result != 0: tmp['valid_loc'] = ["No variants"] #
1126 60         rebuild error
1127 61         else: # pass the rebuilding process
1128 62             tmp['valid_loc'] = valid_locs
1129 63     else:
1130 64         tmp['valid_loc'] = ['No variants']
1131 65 # write back the original content
1132 66 with open(file_path, "w") as f:
1133 67     f.write(old_str)
1134 68 res.append(tmp)

```

Listing 4: Verify the correctness of generated theorems

## C.5 LIMITATIONS OF SYNTHESIS PIPELINE

Our synthesis pipeline is mainly based on the advanced *Leandojo* tool. We use it to interact with Lean, parse abstract syntax trees and trace state-tactic pairs. However, this tool has the following weaknesses. 1) It will generate a significant number of temporary files that consume substantial disk space when initializing a “dojo” environment. The memory-intensive nature of this tool hinders our

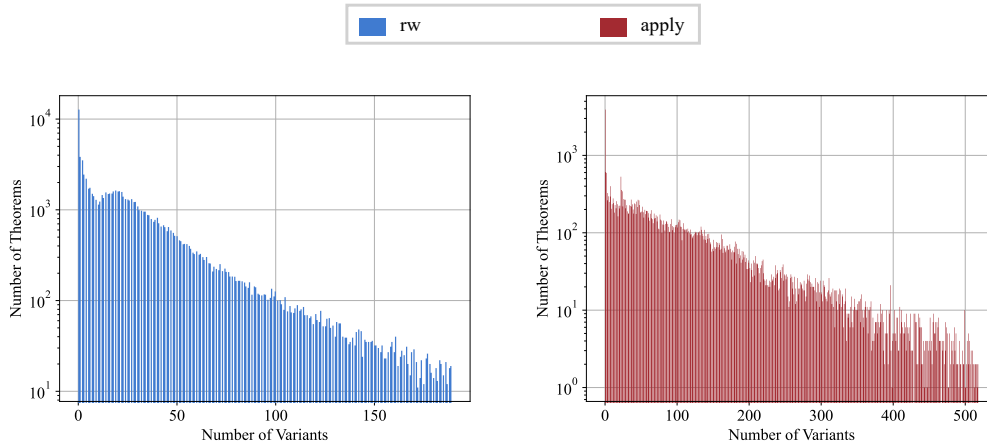
1134 ability to effectively implement multiprocessing. 2) Moreover, it lacks native support for tracing a  
 1135 local Lean repository, so we must first upload our data to GitHub. 3) We encounter challenges when  
 1136 tracing a repository of a scale significantly larger than that of Mathlib, which makes it hard to do  
 1137 multi-round synthesis. We aspire to enhance the functionality of the *Leandojo* tool to tackle more  
 1138 demanding scenarios in our forthcoming endeavors.

1139 In addition, the process of constructing statements and proofs plays an important role in data volume  
 1140 and diversity. Our implementation involves parsing the abstract syntax tree for localization and  
 1141 conducting various string manipulations, which is straightforward but struggles with sophisticated  
 1142 situations such as coercion, naming conflicts, and other corner cases. We are looking forward to  
 1143 refactoring our modification logic with the metaprogramming API of lean<sup>9</sup> in the future, which is  
 1144 more robust and easier to extend.

## 1146 D DEEPER ANALYSIS OF SYNTHETIC DATASET

### 1148 D.1 NUMERICAL ANALYSIS

1150 The histogram of the number of variants synthesized by each tactic is shown in Figure 6.



1167 Figure 6: The distribution of the number of variants (only 99% of the data are visualized).

1170 For each tactic, we also list the top 20 theorems with the highest number of variants in Figure 7.

### 1172 D.2 EXAMPLES

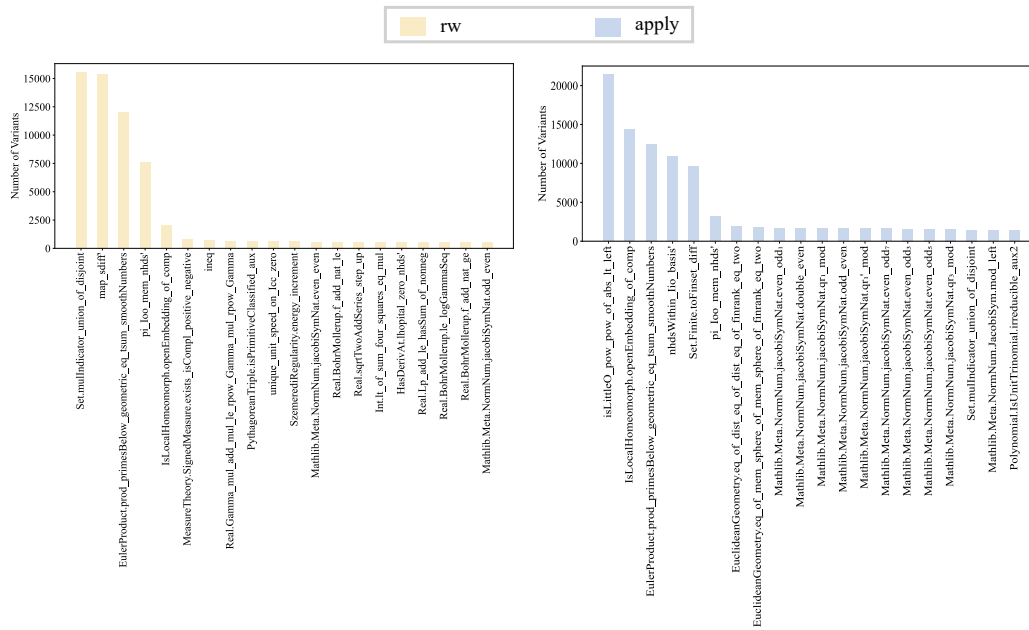
1174 Due to the large volume of synthetic data, it is challenging to display all the data in the appendix.  
 1175 We only display a subset of demo theorems for reference. The proof lengths of these theorems range  
 1176 from 1 to 3 lines. To explore further details, please examine our dataset. The synthesized theorems  
 1177 of *rw* tactic are displayed in Fig 8. The synthesized theorems of *apply* are displayed in Fig 9.

### 1179 D.3 DETAILS OF TRAINING DATA

#### 1181 D.3.1 EXAMPLES OF TRAINING DATA

1182 As shown in Fig 10, we synthesize a series of variants for each candidate theorem by employing  
 1183 different tactic instructions to mutate existing theorems. We simply combine these additional theo-  
 1184 rems with the original theorems in Mathlib and train LLMs on this augmented corpus. In addition  
 1185 to synthesizing variants for each candidate theorem, symbolic manipulations to construct new theo-  
 1186 rems also introduce some new state-tactic pairs. What should be noted is that the state-tactic pairs

1187 <sup>9</sup><https://leanprover-community.github.io/lean4-metaprogramming-book/>

Figure 7: The top20 theorems for *rw* and *apply*.

are extracted by *Leandjo* rather than manually designed symbolic rules. We have not performed any post-processing on the extracted state-tactic pairs. We group the extracted theorems by the employed tactics (*rw*, *apply*, *have*). The examples of *rw* and *apply* are shown in Fig 11. The examples of *have* are shown in Fig 12.

### D.3.2 PREPROCESSING

The synthesized variants of theorems and corresponding state-tactic pairs appearing in the test split of *Leandjo* benchmark are removed. During the data synthesis process, an invocable theorem may be used to rewrite or apply to different candidate theorems. Thus, many data points extracted from the augmented Mathlib repository share the same tactic and invocable theorem (i.e., premise), such as premise A in “*rw* [A]” or “*apply* A”. These data points have similar changes in the proof state. We keep one state-tactic pair for each used premise in the synthesized state-tactic pairs and obtain about 30k data points for each tactic.

### D.3.3 CLASSIFICATION OF EXTRACTED TACTICS

The types of extracted state-tactic pairs are mainly determined by the symbolic manipulations to construct the theorems. We construct the proof by inserting a *have* instruction and integrating it with the original proof. As a result, we manually introduce tactics centered on *rw*, *apply* or *have*. The traced data predominantly features these tactics. The style of the seed theorem (tactic-style or term-style) and the implementation of the tracing tool are also key factors for the traced data. To see more details of this process, it is a good choice to trace the synthesized repository in person. Being familiar with the tracing process will offer some valuable guidance in designing symbolic rules to modify the proof. The extracted state-tactic pairs can also be post-processed (e.g., split the chained tactics into single ones), which has not been explored by our work.

### D.3.4 INFLUENCE OF THE QUANTITY OF SFT DATASET

We assess the impact of varying quantities of additional state-tactics pairs for each tactic under several conditions. 1) Mathlib-train with no additional data points; 2) Downsampling with a ratio of 0.25, resulting in 7.5k additional data points; 3) Downsampling with a ratio of 0.5, resulting in 15k additional data points; 4) Our setting with a deduplication threshold of 1, resulting in 30k additional data points; 5) Deduplication with a threshold of 50, resulting in 500k additional data points; and

---

```

1242                                     Finset.multiplicativeEnergy_mono_right
1243
1244 theorem multiplicativeEnergy_mono_right (ht : t1 ⊆ t2) :
1245   multiplicativeEnergy s t1 ≤ multiplicativeEnergy s t2 :=
1246   multiplicativeEnergy_mono Subset.rfl ht
1247
1248 example (ht : t1 ∩ t2 = t1) :
1249   multiplicativeEnergy s t1 ≤ multiplicativeEnergy s t2 :=
1250   have ht : t1 ⊆ t2 := by rw [Finset.inter_eq_left] at ht; exact ht
1251   multiplicativeEnergy_mono Subset.rfl ht
1252
1253 example (ht : t1.val ⊆ t2.val) :
1254   multiplicativeEnergy s t1 ≤ multiplicativeEnergy s t2 :=
1255   have ht : t1 ⊆ t2 := by rw [←Finset.subset_def] at ht; exact ht
1256   multiplicativeEnergy_mono Subset.rfl ht
1257
1258 example (ht : t1 ⊆ t2) :
1259   max (multiplicativeEnergy s t2) (multiplicativeEnergy s t1) = multiplicativeEnergy s t2 :=
1260   have : multiplicativeEnergy s t1 ≤ multiplicativeEnergy s t2 :=
1261     multiplicativeEnergy_mono Subset.rfl ht
1262   by rw [←max_eq_left_iff] at this; exact this

```

---

```

1261                                     Multiset.card_le_card
1262
1263 theorem card_le_card {s t : Multiset α} (h : s ≤ t) : card s ≤ card t :=
1264   leInductionOn h Sublist.length_le
1265
1266 example {s t : Multiset α} (h : s ≤ t) : ∀ {c : ℕ}, card t < c → card s < c :=
1267   have : card s ≤ card t :=
1268     leInductionOn h Sublist.length_le
1269   by rw [←forall_lt_iff_le'] at this; exact this
1270
1271 example {s t : Multiset α} (h : s ≤ t) : card s ∩ card t = card s :=
1272   have : card s ≤ card t :=
1273     leInductionOn h Sublist.length_le
1274   by rw [←inf_eq_left] at this; exact this
1275
1276 example {s t : Multiset α} (h : s ≤ t) : card s = card t ∨ card s < card t :=
1277   have : card s ≤ card t :=
1278     leInductionOn h Sublist.length_le
1279   by rw [le_iff_eq_or_lt] at this; exact this

```

---

```

1277                                     Nat.one_lt_pow'
1278
1279 theorem one_lt_pow' (n m : ℕ) : 1 < (m + 2) ^ (n + 1) :=
1280   one_lt_pow (n + 1) (m + 2) n.succ_ne_zero (Nat.lt_of_sub_eq_succ rfl)
1281
1282 example (n m : ℕ) : (m + 2) ^ (n + 1) ≠ 0 ∧ (m + 2) ^ (n + 1) ≠ 1 :=
1283   have : 1 < (m + 2) ^ (n + 1) :=
1284     one_lt_pow (n + 1) (m + 2) n.succ_ne_zero (Nat.lt_of_sub_eq_succ rfl)
1285   by rw [Nat.one_lt_iff_ne_zero_and_ne_one] at this; exact this
1286
1287 example (n m : ℕ) : (m + 2) ^ (n + 1) < (m + 2) ^ (n + 1) * (m + 2) ^ (n + 1) :=
1288   have : 1 < (m + 2) ^ (n + 1) :=
1289     one_lt_pow (n + 1) (m + 2) n.succ_ne_zero (Nat.lt_of_sub_eq_succ rfl)
1290   by rw [←Nat.lt_mul_self_iff] at this; exact this

```

---

Figure 8: Examples of synthesized theorems for *rw*

6) No deduplication, resulting in 3M additional data points. We fine-tune Llama-3-8b on these different mixtures of data and evaluate their performance on *random* split of *Leandajo* Benchmark. The experimental results are shown in Fig 4, demonstrating that our setting achieves a relatively optimal balance between overhead and performance.



---

```

1296
1297 StrictMonoOn.mapsTo_Ioc
1298 lemma StrictMonoOn.mapsTo_Ioc (h : StrictMonoOn f (Icc a b)) :
1299   MapsTo f (Ioc a b) (Ioc (f a) (f b)) :=
1300   fun _c hc ↦ (h (left_mem_Icc.2 <| hc.1.le.trans hc.2) (Ioc_subset_Icc_self hc) hc.1,
1301     h.monotoneOn (Ioc_subset_Icc_self hc) (right_mem_Icc.2 <| hc.1.le.trans hc.2) hc.2)
1302
1303 example (h : StrictMonoOn f (Icc a b) ↔ True) :
1304   MapsTo f (Ioc a b) (Ioc (f a) (f b)) :=
1305   have h : StrictMonoOn f (Icc a b) := by apply of_iff_true; assumption
1306   fun _c hc ↦ (h (left_mem_Icc.2 <| hc.1.le.trans hc.2) (Ioc_subset_Icc_self hc) hc.1,
1307     h.monotoneOn (Ioc_subset_Icc_self hc) (right_mem_Icc.2 <| hc.1.le.trans hc.2) hc.2)
1308
1309 example (H : ∀ (b_1 : Prop), (StrictMonoOn f (Icc a b) → b_1) → StrictMonoOn f (Icc a b)) :
1310   MapsTo f (Ioc a b) (Ioc (f a) (f b)) :=
1311   have h : StrictMonoOn f (Icc a b) := by apply peirce'; assumption
1312   ...
1313 example (h : Icc a b ∈ {x | StrictMonoOn f x}) :
1314   MapsTo f (Ioc a b) (Ioc (f a) (f b)) :=
1315   have h : StrictMonoOn f (Icc a b) := by apply Membership.mem.out; assumption
1316   ...
1317
1318 PNat.XgcdType.reduce_a
1319 theorem reduce_a {u : XgcdType} (h : u.r = 0) : u.reduce = u.finish := by
1320   rw [reduce]
1321   exact if_pos h
1322
1323 example {u : XgcdType} (h : 0 | r u) : u.reduce = u.finish := by
1324   have h : u.r = 0 := by apply Nat.eq_zero_of_zero_dvd; assumption
1325   rw [reduce]
1326   exact if_pos h
1327
1328 example {u : XgcdType} (H : u.bp + 1 | u.ap + 1) : u.reduce = u.finish := by
1329   have h : u.r = 0 := by apply Nat.mod_eq_zero_of_dvd; assumption
1330   ...
1331 example {u : XgcdType} (n : N) (H : Nat.gcd (r u) n = 0) : u.reduce = u.finish := by
1332   have h : u.r = 0 := by apply Nat.eq_zero_of_gcd_eq_zero_left<;> assumption
1333   ...
1334
1335 Ordnode.not_le_delta
1336 theorem not_le_delta {s} (H : 1 ≤ s) : ¬s ≤ delta * 0 :=
1337   not_le_of_gt H
1338
1339 example {s} (h : 0 < s) (a : 1 | s) : ¬s ≤ delta * 0 :=
1340   have H : 1 ≤ s := by apply Nat.le_of_dvd<;> assumption
1341   not_le_of_gt H
1342
1343 example {s} (n : N) (H1 : s | n) (H2 : 0 < n) : ¬s ≤ delta * 0 :=
1344   have H : 1 ≤ s := by apply Nat.pos_of_dvd_of_pos<;> assumption
1345   ...
1346 example {s} (l : List N) (p : List.Pairwise LE.le (l :: l)) (a : s ∈ l) : ¬s ≤ delta * 0 :=
1347   have H : 1 ≤ s := by apply List.rel_of_pairwise_cons<;> assumption
1348   ...
1349

```

---

Figure 9: Examples of synthesized theorems for *apply*

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

```

Variant of rw
theorem_name: CategoryTheory.Limits.Multicofork.sigma_condition_variant_0
file_path: Mathlib/CategoryTheory/Limits/Shapes/Multiequalizer.lean
text:
example : MultispanIndex.fstSigmaMap I >> Sigma.desc (π K) ∈ [MultispanIndex.sndSigmaMap I >> Sigma.desc (π K)] := by
  have : I.fstSigmaMap >> Sigma.desc K.π = I.sndSigmaMap >> Sigma.desc K.π := by
    ext
    simp
  rw [←List.mem_singleton] at this; exact this
meta: https://github.com/leanprover-community/mathlib4/commit/3c307701fa7e9acbd0680d7f3b9c9fed9081740

Variant of apply
theorem_name: UniformInducing.equicontinuous_iff_variant_26
file_path: Mathlib/Topology/UniformSpace/Equicontinuity.lean
text:
example {F : ι → X → α} {u : α → β} (B : Set (Set (α → β))) (s : Set (α → β)) (hB : TopologicalSpace.IsTopologicalBasis B)
(hs : IsOpen s) (h : ∀ U ∈ B, U ⊆ s → U ⊆ UniformInducing) (a : u ∈ s) :
  Equicontinuous F ↔ Equicontinuous ((u ∘ ·) ∘ F) := by
  have hu : UniformInducing u := by apply TopologicalSpace.IsTopologicalBasis.subset_of_forall_subset<> assumption
  congrm ∀ x, ?_
  rw [hu.equicontinuousAt_iff]
meta: https://github.com/leanprover-community/mathlib4/commit/3c307701fa7e9acbd0680d7f3b9c9fed9081740

```

Figure 10: Examples of data for pretraining

## E ADDITIONAL EXPERIMENTS

### E.1 EFFECTIVENESS OF DIFFERENT TACTICS

We evaluate the effectiveness of different tactics by combining additional state-tactic pairs of a specific tactic with Mathlib-train and fine-tuning the LLMs using this mixture. The experimental results are shown in Table 6. We observe that state-tactic pairs of *rw* and *apply* are beneficial for the theorem-proving ability of the LLM. And the highest improvement is achieved by the combination of these two tactics. For the state-tactic pairs of *have*, we assume that these data will teach the model to introduce lemmas in the process of proving a theorem, helping them to prove the theorems in multiple steps. However, experimental data show that *have* has complex effects on the proving capacity of LLMs. The performance on a mixture of “have” and other tactics shows poorer results compared to that on a single tactic. We hope to investigate the effectiveness of *have* tactic soon.

### E.2 ANALYSIS OF THE TACTICS TO PROVE *miniF2F* THEOREMS

#### E.2.1 PREFERENCE IN USED TACTICS

To see the preference for the tactics used to prove competition-level problems, we perform a comprehensive analysis of the theorems proved by different LLMs. Specifically, we fine-tune different LLMs with the random train-split of *Leandodo* benchmark and gather all theorems proved by these models. The collection of these models proves 100 theorems out of 244 theorems (41%) on the test split of *miniF2F* benchmark. The average length of the proofs generated by these models is 1.38. And the distribution of these proved theorems is shown in Fig 14. We have the following observations. 1) About half of the theorems in the *miniF2F* test split can be proven with only 1-2 line proofs. 2) Most of the theorems are proved with advanced and automatic tactics in Lean (e.g., *norm\_num*, *linarith*, *omega*, *simp*, etc.). We assume that these tactics play an important role in the theorem-proving ability of LLMs to prove competition-level problems. From the above observations, we assume that synthesizing advanced tactic data points rather than basic data points featuring *rw* and *apply* is promising to improve the performance of proving competition-level problems.

#### E.2.2 INFLUENCE OF ADDITIONAL TACTICS

We analyze the distribution of used tactics in proven *miniF2F* problems across different data compositions. The dynamics of distribution changes are shown in Fig. 15. We assume that increasing

Table 6: The effectiveness of different tactics

Methods	random	novel_premises	Search Budget
<b>Llama3-8b</b>			
Mathlib-train	58.22	38.52	1 × 32
<b>rw tactic</b>			
Mathlib-train + rw	57.85 (-0.37)	41.59 (+3.07)	1 × 32
Mathlib-train + have	58.27 (+0.05)	41.29 (+2.77)	1 × 32
Mathlib-train + rw + have	57.96 (-0.26)	41.53 (+3.01)	1 × 32
<b>apply tactic</b>			
Mathlib-train + apply	56.71 (-1.51)	40.02 (+1.51)	1 × 32
Mathlib-train + have	57.44 (-0.78)	39.24 (+0.72)	1 × 32
Mathlib-train + apply + have	57.23 (-0.99)	38.34 (-0.18)	1 × 32
<b>both tactic</b>			
mathlib-train + rw + apply	58.53 (+0.31)	41.95 (+3.44)	1 × 32
<b>deepseek-coder-7b-base-v1.5</b>			
Mathlib-train	57.7	39.24	1 × 32
<b>rw tactic</b>			
Mathlib-train + rw	58.63 (+0.93)	41.05 (+1.81)	1 × 32
Mathlib-train + have	58.11 (+0.41)	39.06 (-0.18)	1 × 32
Mathlib-train + rw + have	58.74 (+1.04)	40.57 (+1.33)	1 × 32
<b>apply tactic</b>			
Mathlib-train + apply	57.96 (+0.26)	41.17 (+1.93)	1 × 32
Mathlib-train + have	57.02 (-0.68)	39.66 (+0.42)	1 × 32
Mathlib-train + apply + have	58.16 (+0.46)	39.78 (+0.54)	1 × 32
<b>both tactic</b>			
Mathlib-train + rw + apply	58.37 (+0.67)	42.92 (+3.68)	1 × 32

Table 7: The results of miniF2F for different LLMs. We fine-tune each model with the random train-split of *Leandoho* benchmark and evaluate their performance on miniF2F benchmark.

Methods	miniF2F-test
Llama-3-8b	34.01
deepseek-coder-base-7b-v1.5	37.70
deepseek-math-7b-base	34.42
llemma-7b	32.38
mistral-7b	32.38
internlm2-math-7b	36.06
Combination	<b>40.98</b>

the diversity of synthesized tactics and adjusting the tactic distribution will be beneficial to enhance the theorem-proving ability of LLMs.

1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

```

instruction:
/- You are proving a theorem in Lean 4.
You are given the following information:
- The current proof state, inside [STATE]...[/STATE]

Your task is to generate the next tactic in the proof.
Put the next tactic inside [TAC]...[/TAC]
-/
rw tactic
input:
[STATE]
 $\alpha$  : Type u_1
 $\beta$  : Type u_2
R : Type u_3
s t : Finset  $\alpha$ 
f :  $\alpha \rightarrow \beta$ 
n :  $\mathbb{N}$ 
inst† : DecidableEq  $\alpha$ 
this : s.card = 2  $\leftrightarrow$   $\exists x y, x \neq y \wedge s = \{x, y\}$ 
 $\vdash$  s.card = 2  $\leftrightarrow$   $\exists a \in \text{Set.range fun } x y \Rightarrow x \neq y \wedge s = \{x, y\}$ , Exists a
[/STATE]
[TAC]

output:
rw [←Set.exists_range_iff] at this[/TAC]
apply tactic
input:
[STATE]
n :  $\mathbb{N}$ 
inst† : NeZero n
a b : ZMod n
w : val b  $\in$  HAdd.hAdd (val a)  $^{-1}$  Nat.le n
 $\vdash n \leq \text{val } a + \text{val } b$ 
[/STATE]
[TAC]

output:
apply Set.mapsTo_preimage[/TAC]

```

Figure 11: Examples of *rw* and *apply* data points for finetuning

1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565

```

instruction:
/- You are proving a theorem in Lean 4.
You are given the following information:
- The current proof state, inside [STATE]...[/STATE]

Your task is to generate the next tactic in the proof.
Put the next tactic inside [TAC]...[/TAC]
-/

have tactic-rw-related
Case one
input:
[STATE]
m n a b : ℕ
cop : Coprime m n
ha : a ≠ 0
hb : b ≠ Int.toNat 0
⊢ a * m + b * n ≠ m * n
[/STATE]
[TAC]

output:
have hb : b ≠ 0 := by rw [Int.toNat_zero] at hb; exact hb[/TAC]

Case two
input:
[STATE]
α : Type u_1
inst† : DecidableEq α
A : Finset (Finset α)
u v a† : Finset α
r : ℕ
huv : Equiv.swap u.card v.card = Equiv.refl N
a : Finset α
⊢ (compress u v a).card = a.card
[/STATE]
[TAC]

output:
have huv : u.card = v.card := by rw [Equiv.swap_eq_refl_iff] at huv; exact huv[/TAC]

have tactic-apply-related
Input:
[STATE]
α : Type u
β : Type v
inst†¹ : DecidableEq α
inst† : Fintype α
σ : Perm α
n : ℕ
h : orderOf σ ∈ {x | Nat.Coprime n x}
⊢ support (σ ^ n) = support σ
[/STATE]
[TAC]

output:
have h : Nat.Coprime n (orderOf σ) := by apply Membership.mem.out; assumption[/TAC]

```

Figure 12: Examples of *have* data points for finetuning

1566  
 1567  
 1568  
 1569  
 1570  
 1571  
 1572  
 1573  
 1574  
 1575  
 1576  
 1577  
 1578  
 1579  
 1580  
 1581  
 1582  
 1583  
 1584  
 1585  
 1586  
 1587  
 1588  
 1589  
 1590  
 1591  
 1592  
 1593  
 1594  
 1595  
 1596  
 1597  
 1598  
 1599  
 1600  
 1601  
 1602  
 1603  
 1604  
 1605  
 1606  
 1607  
 1608  
 1609  
 1610  
 1611  
 1612  
 1613  
 1614  
 1615  
 1616  
 1617  
 1618  
 1619

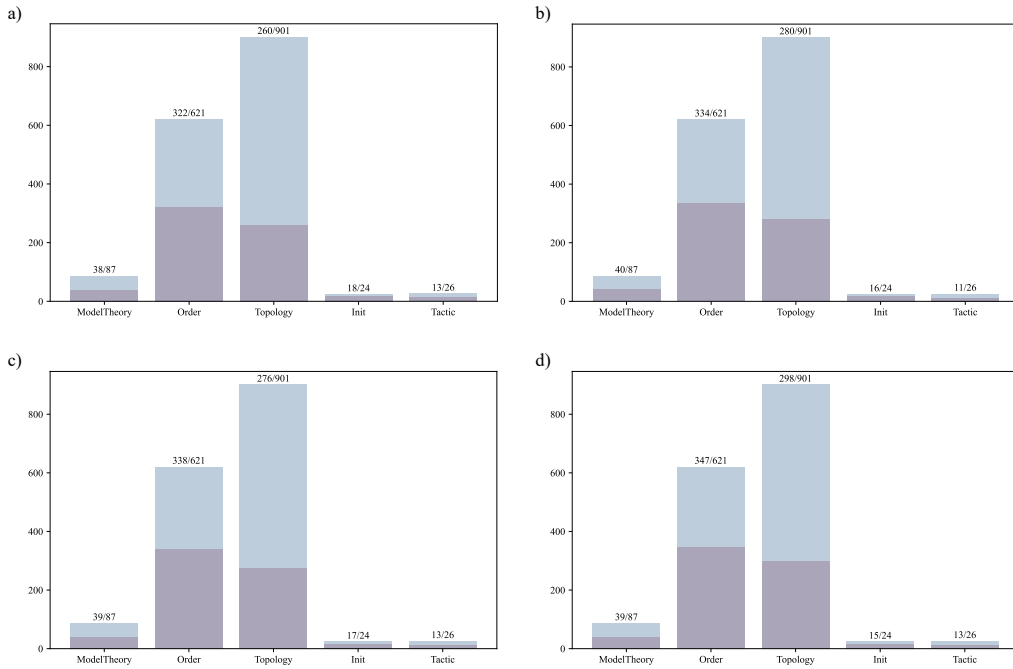


Figure 13: The performance of models fine-tuned on different SFT datasets on novel\_premises split.  
 a) Mathlib-train; b) Mathlib-train + rw; c) Mathlib-train + apply; d) Mathlib-train + rw + apply.

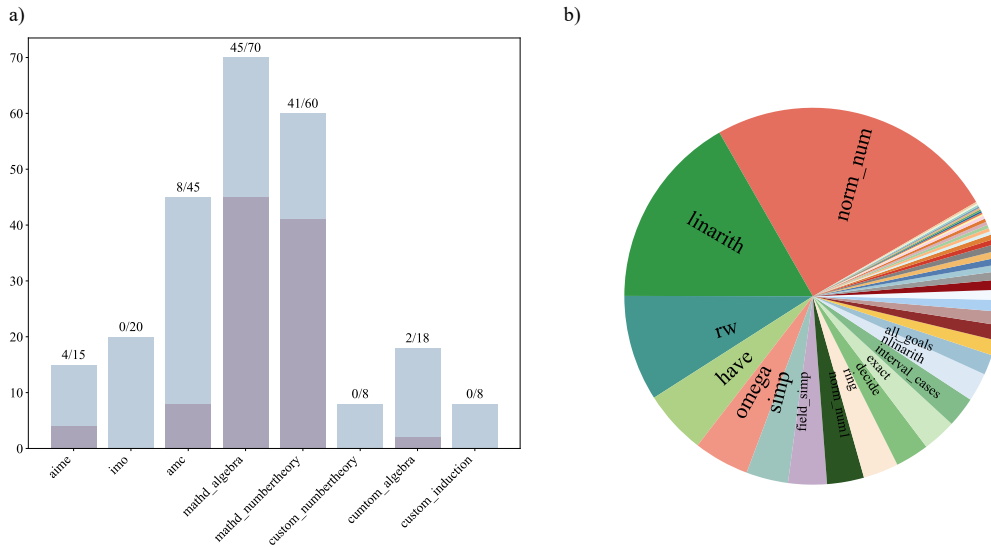


Figure 14: a) The distribution of theorems proved by different LLMs; b) The distribution of tactics used in the proved theorems.

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

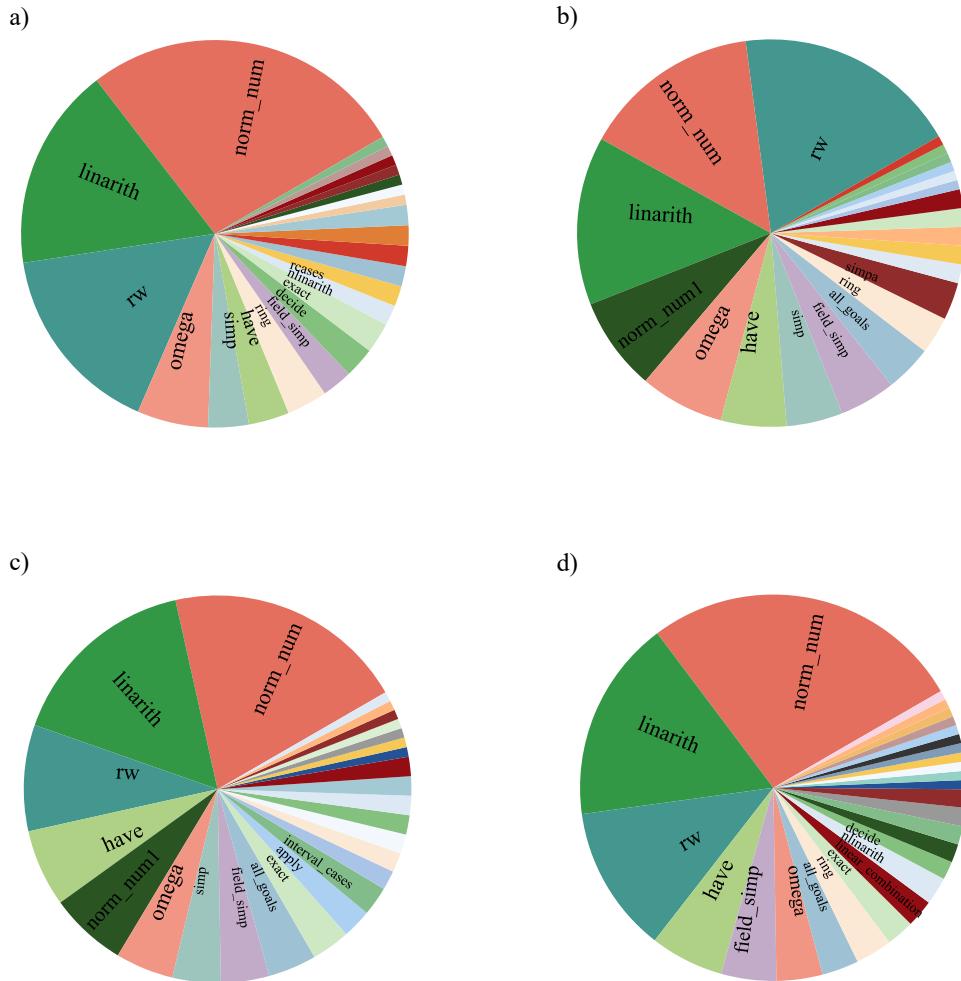


Figure 15: The distribution of used tactics for Llama-3-8b fine-tuned on different SFT datasets to prove miniF2F. a) Mathlib-train; b) Mathlib-train + *rw*; c) Mathlib-train + *apply*; d) Mathlib-train + *rw* + *apply*.