

# STP: Self-play LLM Theorem Provers with Iterative Conjecturing and Proving

Kefan Dong  
Stanford University  
kefandong@stanford.edu

Tengyu Ma  
Stanford University  
tengyuma@stanford.edu

## Abstract

A fundamental challenge in formal theorem proving by LLMs is the lack of high-quality training data. Although reinforcement learning or expert iteration partially mitigates this issue by alternating between LLM generating proofs and finetuning them on correctly generated ones, performance quickly plateaus due to the scarcity of correct proofs (sparse rewards). To keep improving the models with limited data, we draw inspiration from mathematicians, who continuously develop new results, partly by proposing novel conjectures or exercises (which are often variants of known results) and attempting to solve them. We design the Self-play Theorem Prover (STP) that simultaneously takes on two roles, conjecturer and prover, each providing training signals to the other. The conjecturer is trained iteratively on previously generated conjectures that are barely provable by the current prover, which incentivizes it to generate increasingly challenging conjectures over time. The prover attempts to prove the conjectures with standard expert iteration. We evaluate STP with both Lean and Isabelle formal verifiers. With 19.8 billion tokens generated during the training in Lean, STP proves 26.3% of the statements in the LeanWorkbook dataset, doubling the previous best result of 13.2% achieved through expert iteration. The final model achieves state-of-the-art performance among whole-proof generation methods on miniF2F-test (61.1%, pass@3200), Proofnet-test (23.1%, pass@3200) and PutnamBench (8/644, pass@64).

## 1 Introduction

The reasoning capability of large language models (LLMs) is critical for various applications, including coding assistants, question-answering, and agents [Plaat et al., 2024, Shinn et al., 2023, Yao et al., 2022, Shao et al., 2024, Li et al., 2023, Nijkamp et al., 2022]. It is also a key criterion for achieving artificial general intelligence (AGI). Automated theorem proving with formal languages by LLMs stands at the forefront of reasoning research, partly because it allows objective and reliable evaluation through classical verifiers such as Lean [Moura and Ullrich, 2021] and Isabelle [Nipkow et al., 2002]. Moreover, it arguably encapsulates the essence of advanced reasoning tasks while abstracting away the ambiguity of natural language, enabling meaningful studies on a relatively smaller scale.

However, a fundamental challenge in improving reasoning performance—whether in natural or formal languages—lies in the lack of high-quality training data. Collecting reasoning data requires domain experts, making it expensive to scale. There are only a limited number of advanced math papers and theorems in existence, orders of magnitude smaller than other data sources.

Reinforcement learning (RL) on datasets *without* solutions (e.g., datasets with theorem statements or reasoning questions and answers) is a prominent approach for improving the reasoning capability, as seen in the recent development of OpenAI o1 [Jaech et al. [2024], DeepSeek-Prover [Xin et al., 2024a] and DeepSeek R1 [Guo et al., 2025]. Often referred to as expert iteration [Anthony et al. [2017], it partially mitigates the data scarcity issue by alternating between LLMs generating proofs and finetuning them on correctly generated ones [Kaliszyk et al., 2018, Wu et al., 2021, AlphaProof, 2024, Xin et al., 2024b, Ying et al., 2024].

However, as Wu et al. [2024] pointed out, RL or expert iteration often saturates at a low pass rate because the number of samples required to generate a correct proof for an unproven theorem grows exponentially. As a result, a massive amount of computation is wasted on generating incorrect proofs that provide no training signal to the model. For instance, in the proof sampling process of Wu et al. [2024], 98.5% of the compute yields no successful proofs,

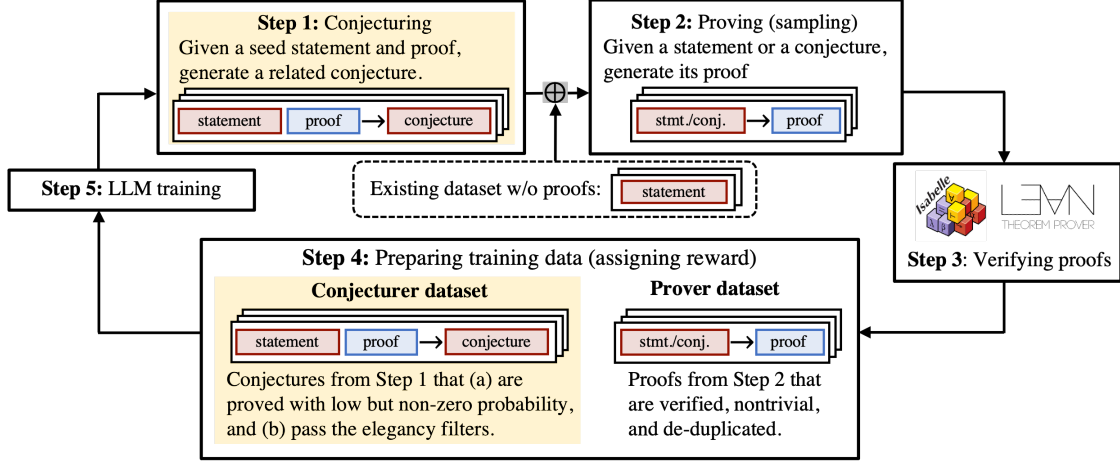


Figure 1: Self-play Theorem Prover (STP). Our model simultaneously takes on two roles — the conjecturer that generates new, related conjecture given a seed theorem with proof (Step 1), and the prover that attempts to prove the statements in an existing dataset and the generated conjectures (Step 2). Step 4 selects the correct, approachable, elegant, yet challenging conjectures to train the conjecturer, and the verifier selects correct proofs in Step 3 to train the prover. The main difference between STP and expert iteration is the conjecturer role highlighted with a yellow background.

despite the pass rate being only 13.2% on LeanWorkbook. In other words, after a few rounds of expert iteration, re-training the model becomes much less effective due to the limited number of new successful proofs.

In addition, RL’s capability is fundamentally bounded by the difficulty level of the theorems in the training dataset—it is unlikely, in principle, for a model to learn college-level proof techniques solely by working on high school-level problems or to solve open math problems using RL on graduate-level problems. Moreover, there are likely not enough open problem statements available for RL training to generalize to other open problems, particularly more advanced ones. In other words, RL or expert iteration algorithms cannot train indefinitely without continuously collecting more theorem statements or math problems.

We need an algorithm that can run and self-improve indefinitely *without more data*. To this end, we draw inspiration from how mathematicians learn and develop advanced mathematics; they refine their understanding and sharpen their proof skills by working on synthesized exercises—variants, extensions, or combinations of existing theorems. Additionally, they frequently propose and publish conjectures, a process widely regarded as just as important, if not more so, than solving them. In other words, unlike the current training of LLMs, mathematicians engage with far more exercises and conjectures (referred to collectively as conjectures in this paper) than the polished, published results found in academic papers and books. Moreover, the continuous generation of new conjectures keeps mathematical fields dynamic and moving forward.

In this paper, we design Self-play Theorem Prover (STP), which mimics how mathematicians learn and develop mathematics. It simultaneously assumes two roles—conjecturer and prover—providing training signals to each other.

As illustrated in Fig. 1, the conjecturer, given a seed theorem with proof, proposes a new, related conjecture (Step 1), while the prover attempts to prove conjectures and statements from an existing dataset (Step 2). Then, the verifier selects correct proofs (Step 3) to train the prover using standard RL and identifies correct, approachable, elegant, yet challenging conjectures to supervise the training of the conjecturer (Step 4). More concretely, in each iteration, the conjecturer is trained on previously generated conjectures that: (a) are barely provable by the current prover (i.e., the prover’s success probability with respect to its random seed is positive but low), and (b) pass certain elegance filters. This iterative process gradually increases the difficulty of conjectures and proofs without requiring additional data. Our method can be viewed either as a self-play algorithm between conjectures and provers or as automated curriculum learning [Portelas et al., 2020] with a self-generated adaptive curriculum (via conjecturers).

We empirically evaluate our method with both Lean [Moura and Ullrich, 2021] and Isabelle [Nipkow et al., 2002].

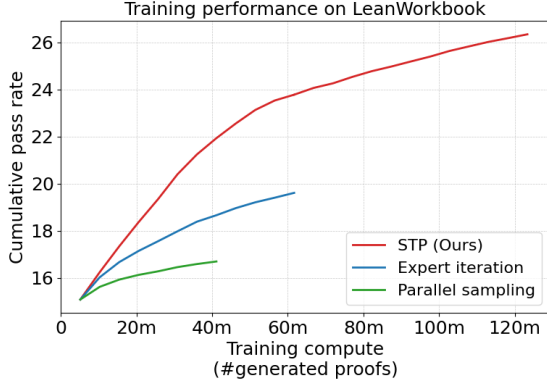


Figure 2: The cumulative pass rates of STP, expert iteration, and parallel sampling on LeanWorkbook shows that STP achieves a much better scaling in terms of the performance vs number of generated proofs. The compute for generating conjectures and training the conjecturer in STP is negligible because the number of generated proofs during training is 64 times the number of conjectures.

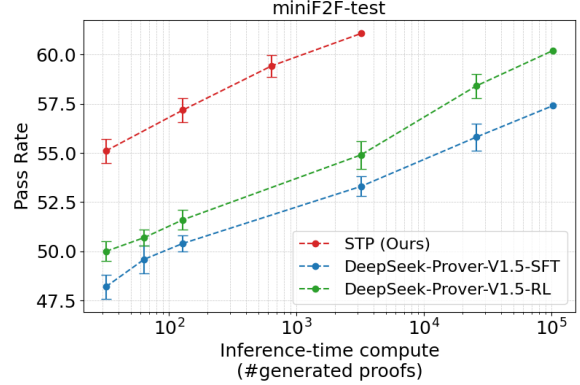


Figure 3: Comparison of pass rates on miniF2F-test (y-axis) with different numbers of inference-time samples (x-axis). The model trained with STP consistently outperforms the DeepSeek-Prover-V1.5 series.

For the Lean experiments, we aim for the best performance and therefore choose DeepSeek-Prover-V1.5-SFT [Xin et al., 2024b] as the base model for STP. As shown in Fig. 2, after a self-play training of roughly 120M generated proofs and 2M generated conjectures, we successfully prove 26.3% of the statements in the training dataset LeanWorkbook [Ying et al., 2024], doubling the previous best result of 13.2% [Wu et al., 2024] achieved by expert iteration. In Fig. 3, we compare the inference-time performance of existing models and the final model trained with STP by taking multiple independent samples on a common benchmark, miniF2F-test [Zheng et al., 2021]. Our model significantly outperforms the DeepSeek-Prover-V1.5 models across various sampling budgets. We also achieve state-of-the-art performance among whole-proof generation methods on miniF2F-test (61.1%, pass@3200), ProofNet-test (23.1%, pass@3200) [Azerbayev et al., 2023a] and PutnamBench (8/644, pass@64) [Tsoukalas et al., 2024], where pass@k represents the percentage of statements proved with  $k$  independently sampled proofs per statement.

In the Isabelle experiments, we study the scalability of STP by starting from a generic math-focused model Llemma-7b [Azerbayev et al., 2023b] and run STP for many more iterations (300M generated proofs in total). We compare the scaling of STP with expert iteration and parallel sampling, by taking several model checkpoints during the STP training run and then switching to the baseline methods. The results clearly demonstrate that STP achieves a better scaling behavior starting from various checkpoints with different capability (see Fig. 4 (Left) in Section 4.3). Ablation study also demonstrates that the main performance gain stems from the dense training signals given by the conjectures. Expert iteration wasted its compute on generating unsuccessful proofs to challenging theorems in the training dataset—at a checkpoint where the pass rate is around 11.4% on LeanWorkbook, only 131 out of 2.5M generated proofs of the unproved statements are correct, resulting in very limited training signals. In contrast, at least 47% of the generated conjectures in STP training are successfully proved because the conjecturer is trained to generate more approachable statements thanks to the design of its reward (see Fig. 4 (Right)).

## 2 Additional Related Works

We refer the readers to Bibel [2013], Loveland [2016] and the reference therein for classical automated theorem proving. Below, we discuss recent works on modern LLM-based theorem provers in addition to what has been discussed in the intro.

**Autoformalization.** A relatively efficient way to create formal proof data is autoformalization, that is, translating natural language math statements and/or proofs to formal language [Jiang et al., 2023, Lu et al., 2024]. A line of research focuses on generating proofs or reasoning steps in natural language and then formalizing the proofs [Jiang et al., 2022a, Zheng et al., 2023, Wang et al., 2023]. Most recently, AlphaProof [2024], Xin et al. [2024a,b] autoformalize statements and then train with expert iteration / RL to write proofs, achieving significant improvement over prior works thanks to the large-scale natural language datasets.

**Formal conjecturing.** Prior works also study how to generate new formal statements/conjectures by neural networks [Urban and Jakubův, 2020, Einarsdóttir et al., 2024, Johansson and Smallbone, 2023] or human-written generators [Polu et al., 2022, Trinh and Luong, 2024], and find that the synthetic statements are generally useful for training the provers [Wang and Deng, 2020, Wu et al., 2020]. Synthetic statements and proofs can also be extracted from an incorrect proof trajectory during RL with hindsight experience replay (HER) Andrychowicz et al. [2017] to speed up the training process [Aygün et al., 2022, Dong et al., 2024]. However, even though the training efficiency is improved, we argue that the final performance is still bounded by difficulty level of the existing dataset because synthetic statements are most likely easier than the given ones in the dataset.

**Self-play and automatic goal generation.** The closest related work to this paper is the concurrent work of Poesia et al. [2024] which also designs a self-play training that iterates between conjecturing and theorem proving. The key difference between this paper and Poesia et al. [2024] is that we start with a pre-trained model and work on practical formal languages like Lean and Isabelle with an infinite space of possible proof steps (which are actions in the RL algorithm), whereas Poesia et al. [2024] operates in a simplified and constrained setting with a finite action space and trains from scratch. As a result, Poesia et al. [2024] rely on constrained decoding to force the validity of generated conjectures, while we solely rely on the LLM itself to generate valid conjectures. Technically, since our training process is much longer (more than 50 iterations) than Poesia et al. [2024] (5 iterations), we must carefully design the conjecturing reward to maintain the diversity and relevance of the generated conjectures (see Section 3.2).

The idea of generating new tasks by the model is also explored in other domains such as alignment [Ye et al., 2024], programming puzzles [Haluptzok et al., 2022, Teodorescu et al., 2023, Pourcel et al., 2024b], video games [Zhang et al., 2023, Pourcel et al., 2024a], and classic RL environments [Parker-Holder et al., 2022, Colas et al., 2022].

More generally, self-play training has demonstrated its potential to achieve super-human performance on two-player games in a fixed environment like Go [Silver et al., 2016].

## 3 Method

On the high level, Self-play Theorem Prover (STP) involves three training stages: (1) model initialization by supervised finetuning, (2) self-play training (visualized in Fig. 1), and (3) final re-training. Unless otherwise stated, we use the term ‘statement’ to refer to the statements in given datasets, and ‘conjecture’ the generated conjectures.

### 3.1 Model initialization by supervised finetuning

In this stage, we initialize the model with two roles, conjecturer and prover, by finetuning a generic LLM (such as the Llama [Touvron et al., 2023]) on a SFT dataset constructed from existing proof libraries such as Mathlib [mathlib Community, 2020]. The proof libraries are organized into files containing human-written formal proofs of known mathematical theorems, and each file formalizes a relatively self-contained result, such as a chapter of a textbook. Our SFT data consists of the following two parts, for finetuning the prover and conjecturer, respectively. Also see concrete examples in Appendix A.1.

**Prover SFT dataset.** We construct a SFT dataset to teach the model to write formal proofs in the given format, where each example is the concatenation of a system prompt (to instruct the model to generate in formal language), a statement, and its corresponding proof. We only compute the next token prediction loss on the proof (which is

the expected output of the model), while the rest is treated as input. To build this dataset, we simply extract all the statement-proof pairs in the proof library files and add a system prompt.

**Conjecturer SFT dataset.** Generally, the conjecturer is to generate a new, related conjecture, given a seed statement with proof that provide the initial ideas. Technically, to further guide the generation of conjecturer, we also provide it a lemma used in the proof of the seed statement, which can be extracted from the verifier,<sup>1</sup> so that the generated conjectures are more likely to be related to the theorem through the lemma. Therefore, the input is a concatenation of the system prompt, a lemma, and a seed statement and its proof, separated by special formatting tokens, and the expected output is a conjecture on which we compute training loss. We also allow the model to generate conjectures with a fixed trivial lemma. To construct this dataset, we extract (lemma, theorem X, theorem Y) tuples from every proof library file such that (a) the lemma and two theorems appears in the file in this particular order, and (b) the lemma is used in the proof of both theorems. The lemma and theorem X will be part of the inputs, and theorem Y will be the output.

### 3.2 Self-play training

Our self-play training stage of STP is shown in Fig. 1. The main difference compared to expert iteration is the conjecturer in Steps 1 and 4, highlighted in a yellow background.

**Generating conjectures and proofs (Steps 1 & 2).** The self-play training starts with collecting a list of the conjecturer’s inputs in the same format as in the conjecture SFT dataset (system prompt, lemma, and theorem), but from theorem-proof pairs where the theorems are from the given dataset without proofs and proofs are previously generated. We extract a seed lemma from the proof, using the verifier.<sup>2</sup> To prevent the model from only focusing on a few particular proof techniques, we de-duplicate the list based on the seed statement and lemma, and randomly drop some inputs whose lemma appears excessively. Then, the LLM generates conjectures from the inputs, and we randomly select a subset of the generated conjectures with size no larger than the number of remaining unproved statements in the given dataset, so that the prover’s compute budget is split equally between the conjectures and statements. (See the pseudo-code and details in Appendix A.2.) For the prover’s inputs, we combine the generated conjectures and the unproved statements in the existing dataset. Then, we independently sample  $K$  proofs per statement/conjecture in Step 2.

**Reward assignments (Step 4).** The major technical challenge of STP is to design the reward function for the conjecturer (in other words, construct the conjecturer dataset in Step 4). The ultimate goal is to incentivize conjecturer to generate diverse, relevant, approachable yet challenging conjectures to provide enough training signals to the prover.

In Step 4, we first organize all generated conjectures and proofs into a list of examples  $\mathcal{D} = \{(t_i, p_i^t, l_i, c_i, p_i^c)\}_{i=1}^n$  where  $t_i$  and  $p_i^t$  represents a seed statement and its proof,  $l_i$  is a lemma used in the proof  $p_i^t$ , and  $c_i, p_i^c$  are the generated conjectures and the generated proof. We will filter  $\mathcal{D}$  as described below and then use  $(t_i, p_i^t, l_i)$  as the input to the conjecturer and  $c_i$  as the output, and  $p_i^c$  as the output of the prover w.r.t. the input  $c_i$ .

To decide whether a conjecture  $c$  is challenging, we use the (empirical) pass rate of the prover estimated by the  $K$  independently generated proofs:

$$\hat{P}(c) \triangleq (\#\{i : c_i = c, p_i^c \text{ is correct}\}) / (\#\{i : c_i = c\}).$$

Then, we select the examples in  $\mathcal{D}$  where (a) lemma  $l_i$  is used in the proof of conjecture  $p_i^c$ , and (b) the pass rate of the conjecture,  $\hat{P}(c_i)$ , is between  $(0, 1/4]$ :

$$\begin{aligned} \overline{\mathcal{D}} \leftarrow \{ & (t_i, p_i^t, l_i, c_i) \mid (t_i, p_i^t, l_i, c_i, p_i^c) \in \mathcal{D}, \\ & \hat{P}(c_i) \in (0, 1/4], p_i^c \text{ is correct}, l_i \text{ is used in } p_i^c \}. \end{aligned}$$

<sup>1</sup>There is no fundamental difference between lemmas and theorems in formal proofs — the naming is purely for better exposition.

<sup>2</sup>In our implementation, lemmas are extracted together with proof verification in Step 3 by configuring the verifiers accordingly.

Here we discard the proofs (of the conjecture)  $p_i^c$  since they are not needed to train the conjecturer, and we remove the duplicated conjectures (that have multiple proofs).

Then, we apply a heuristic elegance filter to discourage the model from generating artificially hard conjectures with complicated goals — we remove conjectures whose minimum proof length divided by the length of the conjecture is in the lowest 20% of remaining examples.

Finally, we re-weight the selected conjectures to preserve the diversity of the conjecturer — the reward for conjecturer cannot only depend on the generated conjectures individually because otherwise the conjecturer’s optimal policy may degenerate to a singular distribution, whereas in reality, the given dataset typically has multiple modes because the statements focus on different topics like algebra, number theory, and calculus. Therefore, our idea is to push the *distribution* of the selected conjectures toward the unproved statements in the existing dataset to maintain the balance between multiple modes. To this end, we compute a distribution  $P$  supported on the selected conjectures that minimizes the Wasserstein distance to the uniform distribution over unproved theorems, denoted by  $Q$ . The matching cost or similarity metric between a conjecture and a statement, used for computing the Wasserstein distance between  $P$  and  $Q$ , is defined as the negative cosine similarity between their embeddings (given by the current model). Finally, we use the distribution  $P$  as the training set for the conjecturer. Pseudo-code of this step is in Appendix A.3, and an efficient implementation is in Appendix A.5.

For the prover dataset, we only select correct generated proofs where the empirical pass rate of the corresponding statement/conjecture is below 1/2. (We consider other correct proofs trivial). We de-duplicate the prover dataset by exact match. Then, the prover is trained on a replay buffer containing the selected proofs from the last three iterations.

**LLM training (Step 5).** We use weighted cross entropy loss computed on the conjectures or proofs (but not the inputs of the model). For the proof dataset, we weight the examples reciprocally to the number of verified proofs to the corresponding statement/conjecture. We also use a length penalization of the form  $\gamma^L$  to reward simpler proofs, where  $\gamma < 1$  is the discount factor and  $L$  is the length of the proof.

### 3.3 Final re-training

To avoid training instability caused by the changing data distribution during self-play, we re-train the final model checkpoint from the base model (before the SFT stage) on a combination of the SFT dataset and all the correct proofs generated during the self-play training whose corresponding statement/conjecture has an empirical pass rate no larger than 1/4. For every statement/conjecture, we randomly keep at most 16 distinct proofs to speedup the training.

## 4 Experiments

This section presents our implementation details of STP, the results of Isabelle and Lean experiments, and the ablation studies, followed by examples of generated conjectures.

### 4.1 Implementation details

**Training datasets.** We use the de-duplicated LeanWorkbook [Ying et al., 2024], which contains around 89K Lean4 statements, as the given dataset without proofs (see Appendix A.4 for details). For the Isabelle experiments, we translate the Lean4 statements to Isabelle using the DeepSeek V2.5 with few-shot prompting.

The SFT dataset for the Isabelle experiments is extracted from AFP<sup>3</sup> and Isabelle built-in files such as HOL. For the Lean experiments, we first sample 32 proofs per statement in LeanWorkbook since our base model, DeepSeek-Prover-V1.5-SFT, is already trained on it, and combine the correct proofs with examples extracted from the proof library Mathlib4 [mathlib Community, 2020] as the SFT dataset.

---

<sup>3</sup><https://www.isa-afp.org/>



**Periodic refreshing.** With a limited replay buffer, the model may forget some proof skills learned in the SFT stage after many iterations. Therefore, during our STP training, we periodically re-train the model from the base model on all previously generated correct proofs, following a procedure similar to the final re-training in Section 3.3. After refreshing, we reset the replay buffer and restart the self-play training using the re-trained model checkpoint.

**Verifiers’ setup.** To study the scalability of STP with limited compute, in the Isabelle experiments, we disable the advanced proof tactics `sledgehammer`, `mason`, `smt`, `metis`, `sos`, which require huge CPU compute, to allow more training iterations, sacrificing verification strength and overall performance. We use PISA [Jiang et al., 2021] to interact with Isabelle, and enforce a 10s timeout for any proof step and 360s timeout for entire proofs. For Lean, we follow Xin et al. [2024b], which allows all proof tactics, and set a 200s timeout and a 15GB memory limit for each proof.

**Hyperparameters.** For inference, we cap the number of generated tokens to 1024, and set the sampling temperature to 0.7 for Llemma-7b and 1.0 for DeepSeek-Prover, following Dong et al. [2024], Xin et al. [2024b], respectively. For training, we use batch size 2048 and Adam [Kingma and Ba, 2014] with a constant learning rate of  $5e-5$  in STP, and  $1e-4$  in SFT and final re-training. The discount factor is  $\gamma = \exp(-0.001)$ .

In each iteration of STP, we sample  $K = 32$  proofs per conjecture/statement. For the expert iteration and parallel sampling, we use  $K = 64$ . Since we maintain the number of generated conjectures per iteration to be at most the number of unproved statements in the given dataset, STP has the same sample budget as the baseline methods per iteration.

## 4.2 Results with Lean

For the Lean experiments, we choose DeepSeek-Prover-V1.5-SFT as our base model, which was trained on proofs collected by expert iteration on a combination of public, such as LeanWorkbook, miniF2F-valid [Zheng et al., 2021], and ProofNet-valid [Azerbayev et al., 2023a], and other proprietary datasets. We ran 24 iterations of STP and generated 2M conjectures, 120M proofs, and 19.8B tokens in total. We use the cumulative pass rate, defined by the fraction of statements proved during the entire training, as the main metric for training progress. Fig. 2 plots the cumulative pass rate of STP and two major baselines, expert iteration, and parallel sampling, on the training dataset LeanWorkbook [Ying et al., 2024]. Expert iteration alternates between generating proofs to the statements in the given dataset and finetuning the model on correct proofs. (See discussions and comparison about variants of expert iteration in Appendix A.6.) Parallel sampling simply generates proofs with the given model. Fig. 2 shows that STP achieves significantly better scaling than expert iteration, which simulates the performance of DeepSeek’s model as if it were trained for more iterations.

To achieve the best performance on common benchmarks, we further train the model using statements from LeanWorkbook, miniF2F-valid, and ProofNet-valid for another 8 iterations. Table 1 compares our models with prior works on miniF2F-test and ProofNet-test, which contain formal statements of high-school level and college level math questions, respectively. Among the whole-proof generation methods, STP significantly outperforms DeepSeek-Prover-V1.5-RL (which is continuously trained with RL on top of their SFT model) and achieves SoTA performance across various inference-time sample budgets.

Table 1 also compares STP with tree search methods such as InternLM2.5-StepProver [Wu et al., 2024], which use LLMs to generate single proof steps conditioned on the current verifier’s proof state and then find a complete proof by best first search or MCTS. The sample budget of these methods are not directly comparable with whole-proof generation methods because (a) the number of steps in a generated proof varies significantly, (b) LLMs in tree search methods need to process additional tokens related to the verifier’s proof state, and (c) methods like InternLM2.5-StepProver [Wu et al., 2024] require an additional LLM as the value function. Moreover, it’s conceivable that tree search methods can also be used with STP, so essentially these are orthogonal methods. Nonetheless, we compute the total number of proof steps per statement generated by STP as an proxy for the total number of LLM output tokens for STP and tree search methods, ignoring the additional compute required by tree search methods to process the verifier’s proof state and query the value function. Results in Table 1 indicate that STP also outperforms prior tree search methods with similar (estimated) inference-time budgets.

Table 1: Pass rate on miniF2F [Zheng et al., 2021] and ProofNet [Azerbaiyev et al., 2023a] with different inference-time sample budgets. Our method, STP, achieves state-of-the-art performance among whole-proof generation methods across various sample budgets. For reference, we also include tree search methods, even though they are orthogonal to our main contribution. The sample budgets of tree search methods are not fully comparable to that of the whole proof generation because they also use the LLM to process the verifier’s internal proof state.

Method	Sample budget (#Proofs)	Sample budget (#Steps)	MiniF2F-test	ProofNet-test
<i>Whole-Proof Generation Methods</i>				
TheoremLlama [Wang et al., 2024]	128	-	33.6%	-
DSP [Jiang et al., 2022a]	100	-	39.3%	-
DeepSeek-Prover-V1.5-SFT	128	-	50.4% $\pm$ 0.4%	15.9% $\pm$ 0.6%
[Xin et al., 2024b]	3200	-	53.3% $\pm$ 0.5%	21.0% $\pm$ 0.9%
DeepSeek-Prover-V1.5-RL	128	-	51.6% $\pm$ 0.5%	18.2% $\pm$ 0.5%
[Xin et al., 2024b]	3200	-	54.9% $\pm$ 0.7%	22.0% $\pm$ 0.5%
	25,600	-	58.4% $\pm$ 0.6%	23.7%
	102,400	-	60.2%	-
STP	128	1.1K	57.2% $\pm$ 0.6%	18.0% $\pm$ 0.7%
	3200	28.2K	61.1%	23.1%
STP	128	1.1K	<b>57.7% <math>\pm</math> 0.6%</b>	<b>18.5% <math>\pm</math> 0.8%</b>
(w/ miniF2F-valid, ProofNet-valid) <sup>4</sup>	3200	28.0K	<b>61.7% <math>\pm</math> 0.6%</b>	<b>23.1% <math>\pm</math> 0.5%</b>
	25,600	224K	<b>63.5%</b>	<b>25.8%</b>
<i>Tree Search Methods<sup>5</sup></i>				
ReProver [Yang et al., 2024]	-	-	26.5%	-
PACT [Zheng et al., 2021]	-	$8 \times 16 \times 512 = 66K$	29.2%	-
GPT-f [Polu et al., 2022]	-	$64 \times 8 \times 512 = 262K$	36.6%	-
HTPS [Lample et al., 2022]	-	$64 \times 5000 = 320K$	41.0%	-
Lean-STaR [Lin et al., 2024]	-	$64 \times 1 \times 50 = 3.2K$	46.3%	-
DeepSeek-Prover-V1.5-RL + RMaxTS <sup>6</sup>	3200	-	55.0% $\pm$ 0.7%	21.5% $\pm$ 0.8%
[Xin et al., 2024b]	25,600	-	59.6% $\pm$ 0.6%	25.3%
	204,800	-	63.5%	-
InternLM2.5-StepProver	-	$4 \times 32 \times 600 = 77K$	58.5% $\pm$ 0.9%	-
[Wu et al., 2024]	-	$16 \times 32 \times 600 = 307K$	62.5% $\pm$ 0.5%	-
	-	$256 \times 32 \times 600 = 4.9M$	65.9%	-

As shown in Table 3, on PutnamBench [Tsoukalas et al., 2024] which consists of undergraduate-level mathematics competition questions, STP solves 8 out of 644 problems with 64 samples per problem, and STP (w/ miniF2F valid, ProofNet valid) solves 9 problems with 3200 samples per problem, outperforming the best result of 6 problems in prior works achieved by Wu et al. [2024].

### 4.3 Results with Isabelle

For Isabelle experiments, we start with the Llemma-7b [Azerbaiyev et al., 2023b], math-focused model, and run 58 iterations of STP to study its scalability. We take several checkpoints during STP training and then switch to the expert iteration and parallel sampling baselines to study the scalability of the algorithm from checkpoints with various

<sup>4</sup>Our base model, DeepSeek-Prover-V1.5-SFT, is trained on miniF2F-valid and ProofNet-valid, although we only include these statements in the last 8 iterations of this experiment.

<sup>5</sup>The #Steps for tree search methods is typically calculated by #Independent runs  $\times$  #Tactics generated per search step  $\times$  #Search steps, or #Independent runs  $\times$  #Search steps.

<sup>6</sup>DeepSeek-Prover-V1.5-RL + RMaxTS is a tree search method that uses the LLMs to generate complete proofs during the search instead of single proof steps. Therefore, we treat their sample budget as the number of generated proofs instead of steps.



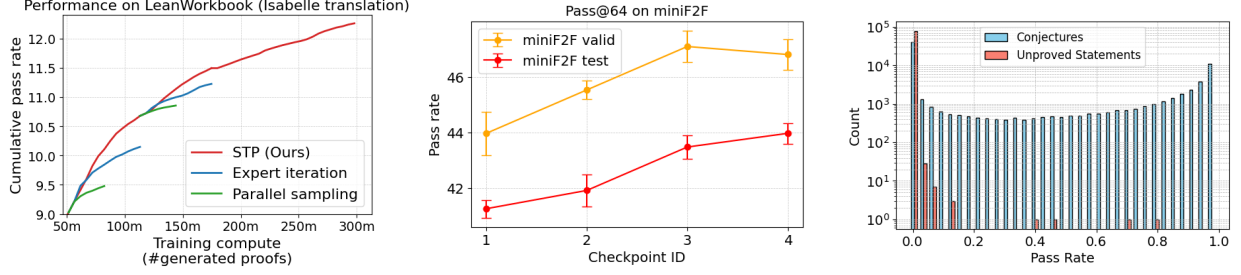


Figure 4: **Left:** Cumulative pass rate on LeanWorkbook (translated into Isabelle) of STP, expert iteration, and parallel sampling, started from two checkpoints in STP training. STP achieves better scaling starting from both checkpoints. For better visualization, the x-axis starts with 50m in this figure, and we defer the full plot to Fig. 5 (Right) in Appendix B.2. **Middle:** The performance of our model on miniF2F gradually improves during the training process. Note that our model is not trained on miniF2F valid and we disallow advanced tactics such as `sos`. The checkpoints are taken roughly per 68M generated proofs. **Right:** Histogram of empirical pass rates of generated conjectures and unproved statements in the training dataset at a checkpoint where the cumulative pass rate on LeanWorkbook (Isabelle translation) is 11.4%. The generated conjectures are significantly more likely to be proved (i.e., has a positive pass rate) than the unproved statements in the dataset, and therefore provide denser training signal. Note that the y-axis is in log scale.

capability. Fig. 4 (Left) compares their cumulative pass rates on LeanWorkbook (Isabelle translation), showing that STP consistently achieves a better scaling across the training process. The model also gradually improves on miniF2F over the training process, as shown in Fig. 4 (Middle).

#### 4.4 Ablation study

**Generated conjectures provide denser training signals.** Fig. 4 (Right) shows the histogram of empirical pass rates of the generated conjectures and the unproved statements in LeanWorkbook using a checkpoint in the Isabelle experiment. Only 131 out of 2.5M generated proofs for the 79K unproved statements are correct. As a result, finetuning the model on correct proofs has almost no effect, and thus expert iteration plateaus. In contrast, generated conjectures by STP offer has higher pass rates and thus more training signals, leading to better scaling.

**Re-training with generated conjectures still helps downstream performance.** One may hypothesis that the self-play algorithm and generated conjectures only help improve the pass rate on LeanWorkBook. It turns out that in the final re-training stage, it is still beneficial to re-train with the generated conjectures in addition to the successfully proved statements in LeanWorkBook even for performance on miniF2F-test and ProofNet-test—it leads to about 1% performance gain (for pass@128) than re-training only on the latter (See Appendix B.1).

#### 4.5 Examples of generated conjectures

In this section, we list three manually selected examples of the generated conjectures at the last iteration of the Lean experiment to demonstrate the quality of generated conjectures.

**Example 1.** The generated conjecture says  $(1 + x)^{2n} \geq 1 + x^n$  when  $n \geq 1$  is an integer and  $x \in [0, 1]$ :

```
theorem lean_workbook_36081' (x : ℝ) (hx : 0 ≤ x ∧ x ≤ 1) : ∀ n : ℕ, n ≥ 1 → (1 + x) ^ (2 * n) ≥ 1 + x ^ n
```

The seed statement says  $1 + x^2 \leq (1 + x)^2$  when  $x \in [0, 1]$ :

```
theorem lean_workbook_36081 (x : ℝ) (hx : 0 ≤ x ∧ x ≤ 1) : 1 + x ^ 2 ≤ (1 + x) ^ 2
```

In this case, the conjecture is harder than the original statement but is proved with similar techniques — expanding the powers of a binomial and then using the fact that  $x \geq 0$ .

**Example 2.** The generated conjecture says  $(x^n - 1) \bmod (x - 1) \leq 1$  if  $x, n$  are integers:

```
theorem lean_workbook_54038' (x : ℕ) (n : ℕ) (hn : 1 < n) : (x^n - 1) %
```

The seed statement says  $n - 1$  divides  $n^k - 1$ :

```
theorem lean_workbook_54038 (n : ℕ) (k : ℕ) (hn : 2 ≤ n) : n - 1 | n^k - 1
```

In this case, our model generates a variant of the original statement by realizing that  $b \bmod a$  equals zero if  $a$  divides  $b$ . This conjecture may help the model connect its proof technique in algebra and number theory. However, the conjecture itself is somewhat unusual and the inequality is not tight. Therefore it is unlikely to be included in any datasets.

**Example 3.** The generated conjecture says  $\sum_{i \geq 0} ((1/4)^i \cdot a) = \frac{a}{1-1/4}$  if  $0 < a \leq 1$ .

```
theorem lean_workbook_plus_46203' (a : ℝ) (ha : 0 < a ∧ a ≤ 1) : Σ' (i : ℕ), 1 / 4  
^ i * a = a / (1 - 1 / 4)
```

The seed statement is a special case where  $a = \sqrt{5}/3$ :

```
theorem lean_workbook_plus_46203 :  
Σ' k : ℕ, (1 / 4)^k * (Real.sqrt 5 / 4) = (Real.sqrt 5 / 3)
```

In this case, the conjecture generalizes the given statement by replacing `Real.sqrt 5 / 4` with a variable  $a$ .

## 5 Conclusion

This paper designs Self-play Theorem Prover (STP) that simultaneously has two roles, conjecturer and prover. By providing training signals to each other, STP goes beyond the statements in the given dataset and its performance continuously improves. Our final model significantly outperforms Deepseek-Prover-V1.5 series and achieves state-of-the-art performance among whole-proof generation methods on common formal proof benchmarks.

## Acknowledgment

The authors would like to thank Yinuo Ren, Zhizhou Ren, Woosuk Kwon, David Hall, Huajian Xin and Kaiyue Wen for their helpful discussions. The authors would also like to thank the support from NSF RI 2211780, and NSF CIF 2212263, and the Google TPU Research Cloud for the computing resources that enabled these experiments.

## References

- AlphaProof. Ai achieves silver-medal standard solving international mathematical olympiad problems. 2024. URL <https://deepmind.google/discover/blog/ai-solves-imo-problems-at-silver-medal-level/>.
- Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. *Advances in neural information processing systems*, 30, 2017.
- Thomas Anthony, Zheng Tian, and David Barber. Thinking fast and slow with deep learning and tree search. *Advances in neural information processing systems*, 30, 2017.

- Eser Aygün, Ankit Anand, Laurent Orseau, Xavier Glorot, Stephen M McAleer, Vlad Firoiu, Lei M Zhang, Doina Precup, and Shihab Mourad. Proving theorems using incremental learning and hindsight experience replay. In *International Conference on Machine Learning*, pages 1198–1210. PMLR, 2022.
- Zhangir Azerbayev, Bartosz Piotrowski, Hailey Schoelkopf, Edward W Ayers, Dragomir Radev, and Jeremy Avigad. Proofnet: Autoformalizing and formally proving undergraduate-level mathematics. *arXiv preprint arXiv:2302.12433*, 2023a.
- Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster, Marco Dos Santos, Stephen McAleer, Albert Jiang, Jia Deng, Stella Biderman, and Sean Welleck. Llemma: An open language model for mathematics. In *The 3rd Workshop on Mathematical Reasoning and AI at NeurIPS’23*, 2023b.
- Wolfgang Bibel. *Automated theorem proving*. Springer Science & Business Media, 2013.
- Cédric Colas, Tristan Karch, Olivier Sigaud, and Pierre-Yves Oudeyer. Autotelic agents with intrinsically motivated goal-conditioned reinforcement learning: a short survey. *Journal of Artificial Intelligence Research*, 74:1159–1199, 2022.
- Kefan Dong, Arvind Mahankali, and Tengyu Ma. Formal theorem proving by rewarding llms to decompose proofs hierarchically. *arXiv preprint arXiv:2411.01829*, 2024.
- Sólrún Halla Einarsdóttir, Yousef Alhessi, Emily First, and Moa Johansson. On lemma conjecturing using neural, symbolic and neuro-symbolic approaches. 2024.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Patrick Haluptzok, Matthew Bowers, and Adam Tauman Kalai. Language models can teach themselves to program better. *arXiv preprint arXiv:2207.14502*, 2022.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- Albert Q Jiang, Sean Welleck, Jin Peng Zhou, Wenda Li, Jiacheng Liu, Mateja Jamnik, Timothée Lacroix, Yuhuai Wu, and Guillaume Lample. Draft, sketch, and prove: Guiding formal theorem provers with informal proofs. *arXiv preprint arXiv:2210.12283*, 2022a.
- Albert Q Jiang, Wenda Li, and Mateja Jamnik. Multilingual mathematical autoformalization. *arXiv preprint arXiv:2311.03755*, 2023.
- Albert Qiaochu Jiang, Wenda Li, Jesse Michael Han, and Yuhuai Wu. Lisa: Language models of Isabelle proofs. In *6th Conference on Artificial Intelligence and Theorem Proving*, pages 378–392, 2021.
- Albert Qiaochu Jiang, Wenda Li, Szymon Tworkowski, Konrad Czechowski, Tomasz Odrzygóźdź, Piotr Miłoś, Yuhuai Wu, and Mateja Jamnik. Thor: Wielding hammers to integrate language models and automated theorem provers. *Advances in Neural Information Processing Systems*, 35:8360–8373, 2022b.
- Moa Johansson and Nicholas Smallbone. Exploring mathematical conjecturing with large language models. 2023.
- Cezary Kaliszyk, Josef Urban, Henryk Michalewski, and Miroslav Olšák. Reinforcement learning of theorem proving. *Advances in Neural Information Processing Systems*, 31, 2018.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pages 611–626, 2023.
- Guillaume Lample, Timothee Lacroix, Marie-Anne Lachaux, Aurelien Rodriguez, Amaury Hayat, Thibaut Lavril, Gabriel Ebner, and Xavier Martinet. Hypertree proof search for neural theorem proving. *Advances in neural information processing systems*, 35:26337–26349, 2022.
- Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. Starcoder: may the source be with you! *arXiv preprint arXiv:2305.06161*, 2023.
- Haohan Lin, Zhiqing Sun, Yiming Yang, and Sean Welleck. Lean-star: Learning to interleave thinking and proving. *arXiv preprint arXiv:2407.10040*, 2024.
- Donald W Loveland. *Automated theorem proving: A logical basis*. Elsevier, 2016.
- Jianqiao Lu, Yingjia Wan, Zhengying Liu, Yinya Huang, Jing Xiong, Chengwu Liu, Jianhao Shen, Hui Jin, Jipeng Zhang, Haiming Wang, et al. Process-driven autoformalization in lean 4. *arXiv preprint arXiv:2406.01940*, 2024.
- The mathlib Community. The lean mathematical library. In *Proceedings of the 9th ACM SIGPLAN International Conference on Certified Programs and Proofs*, CPP 2020, page 367–381, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450370974. doi: 10.1145/3372885.3373824. URL <https://doi.org/10.1145/3372885.3373824>.
- Leonardo de Moura and Sebastian Ullrich. The lean 4 theorem prover and programming language. In *Automated Deduction–CADE 28: 28th International Conference on Automated Deduction, Virtual Event, July 12–15, 2021, Proceedings 28*, pages 625–635. Springer, 2021.
- Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. Codegen: An open large language model for code with multi-turn program synthesis. *arXiv preprint arXiv:2203.13474*, 2022.
- Tobias Nipkow, Markus Wenzel, and Lawrence C Paulson. *Isabelle/HOL: a proof assistant for higher-order logic*. Springer, 2002.
- Jack Parker-Holder, Minqi Jiang, Michael Dennis, Mikayel Samvelyan, Jakob Foerster, Edward Grefenstette, and Tim Rocktäschel. Evolving curricula with regret-based environment design. In *International Conference on Machine Learning*, pages 17473–17498. PMLR, 2022.
- Aske Plaat, Annie Wong, Suzan Verberne, Joost Broekens, Niki van Stein, and Thomas Back. Reasoning with large language models, a survey. *arXiv preprint arXiv:2407.11511*, 2024.
- Gabriel Poesia, David Broman, Nick Haber, and Noah D Goodman. Learning formal mathematics from intrinsic motivation. *arXiv preprint arXiv:2407.00695*, 2024.
- Stanislas Polu, Jesse Michael Han, Kunhao Zheng, Mantas Baksys, Igor Babuschkin, and Ilya Sutskever. Formal mathematics statement curriculum learning. *arXiv preprint arXiv:2202.01344*, 2022.
- Rémy Portelas, Cédric Colas, Lilian Weng, Katja Hofmann, and Pierre-Yves Oudeyer. Automatic curriculum learning for deep rl: A short survey. *arXiv preprint arXiv:2003.04664*, 2020.
- Guillaume Pourcel, Thomas Carta, Grgur Kovač, and Pierre-Yves Oudeyer. Autotelic llm-based exploration for goal-conditioned rl. In *Intrinsically Motivated Open-ended Learning Workshop at NeurIPS 2024*, 2024a.
- Julien Pourcel, Cédric Colas, Gaia Molinaro, Pierre-Yves Oudeyer, and Laetitia Teodorescu. Aces: generating diverse programming puzzles with autotelic language models and semantic descriptors. *Neurips*, 2024b.

- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, YK Li, Y Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- Noah Shinn, Federico Cassano, Beck Labash, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning.(2023). *arXiv preprint cs.AI/2303.11366*, 2023.
- David Silver, Aja Huang, Christopher J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7676): 484–503, 2016.
- Laetitia Teodorescu, Cédric Colas, Matthew Bowers, Thomas Carta, and Pierre-Yves Oudeyer. Codeplay: Autotelic learning through collaborative self-play in programming environments. In *IMOL 2023-Intrinsically Motivated Open-ended Learning workshop at NeurIPS 2023*, 2023.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Trieu Trinh and Thang Luong. Alphageometry: An olympiad-level ai system for geometry. *Google DeepMind*, 17, 2024.
- George Tsoukalas, Jasper Lee, John Jennings, Jimmy Xin, Michelle Ding, Michael Jennings, Amitayush Thakur, and Swarat Chaudhuri. Putnambench: Evaluating neural theorem-provers on the putnam mathematical competition. *arXiv preprint arXiv:2407.11214*, 2024.
- Josef Urban and Jan Jakubův. First neural conjecturing datasets and experiments. In *Intelligent Computer Mathematics: 13th International Conference, CICM 2020, Bertinoro, Italy, July 26–31, 2020, Proceedings 13*, pages 315–323. Springer, 2020.
- Haiming Wang, Huajian Xin, Chuanyang Zheng, Lin Li, Zhengying Liu, Qingxing Cao, Yinya Huang, Jing Xiong, Han Shi, Enze Xie, et al. Lego-prover: Neural theorem proving with growing libraries. *arXiv preprint arXiv:2310.00656*, 2023.
- Mingzhe Wang and Jia Deng. Learning to prove theorems by learning to generate theorems. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, pages 18146–18157, 2020.
- Ruida Wang, Jipeng Zhang, Yizhen Jia, Rui Pan, Shizhe Diao, Renjie Pi, and Tong Zhang. Theoremllama: Transforming general-purpose llms into lean4 experts. *arXiv preprint arXiv:2407.03203*, 2024.
- Minchao Wu, Michael Norrish, Christian Walder, and Amir Dezfouli. Tacticzero: Learning to prove theorems from scratch with deep reinforcement learning. *Advances in Neural Information Processing Systems*, 34:9330–9342, 2021.
- Yuhuai Wu, Albert Qiao Chu Jiang, Jimmy Ba, and Roger Grosse. Int: An inequality benchmark for evaluating generalization in theorem proving. *arXiv preprint arXiv:2007.02924*, 2020.
- Zijian Wu, Suozhi Huang, Zhejian Zhou, Huaiyuan Ying, Jiayu Wang, Dahua Lin, and Kai Chen. Internlm2. 5-stepprover: Advancing automated theorem proving via expert iteration on large-scale lean problems. *arXiv preprint arXiv:2410.15700*, 2024.
- Huajian Xin, Daya Guo, Zhihong Shao, Zhizhou Ren, Qihao Zhu, Bo Liu, Chong Ruan, Wenda Li, and Xiaodan Liang. Deepseek-prover: Advancing theorem proving in llms through large-scale synthetic data. *arXiv preprint arXiv:2405.14333*, 2024a.

- Huajian Xin, ZZ Ren, Junxiao Song, Zhihong Shao, Wanjia Zhao, Haocheng Wang, Bo Liu, Liyue Zhang, Xuan Lu, Qiushi Du, et al. Deepseek-prover-v1.5: Harnessing proof assistant feedback for reinforcement learning and monte-carlo tree search. *arXiv preprint arXiv:2408.08152*, 2024b.
- Kaiyu Yang, Aidan Swope, Alex Gu, Rahul Chalamala, Peiyang Song, Shixing Yu, Saad Godil, Ryan J Prenger, and Animashree Anandkumar. Leandjo: Theorem proving with retrieval-augmented language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.
- Ziyu Ye, Rishabh Agarwal, Tianqi Liu, Rishabh Joshi, Sarmishta Velury, Quoc V Le, Qijun Tan, and Yuan Liu. Evolving alignment via asymmetric self-play. *arXiv preprint arXiv:2411.00062*, 2024.
- Huaiyuan Ying, Zijian Wu, Yihan Geng, Jiayu Wang, Dahua Lin, and Kai Chen. Lean workbook: A large-scale lean problem set formalized from natural language math problems. *arXiv preprint arXiv:2406.03847*, 2024.
- Jenny Zhang, Joel Lehman, Kenneth Stanley, and Jeff Clune. Omni: Open-endedness via models of human notions of interestingness. *arXiv preprint arXiv:2306.01711*, 2023.
- Chuanyang Zheng, Haiming Wang, Enze Xie, Zhengying Liu, Jiankai Sun, Huajian Xin, Jianhao Shen, Zhenguo Li, and Yu Li. Lyra: Orchestrating dual correction in automated theorem proving. *arXiv preprint arXiv:2309.15806*, 2023.
- Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. minif2f: a cross-system benchmark for formal olympiad-level mathematics. In *International Conference on Learning Representations*, 2021.



## A Additional Implementation Details

In this section, we list the missing implementation details.

### A.1 Examples of inputs and outputs of our model

Here we present some concrete examples to demonstrate the input and output formats of our model.

**Examples of the conjecturer.** In the following, we show examples of the conjecturer’s inputs and outputs. Note that `<lemma>`, `<easy theorem>`, `<hard theorem>`, `</hard theorem>` are the formatting tokens, and the system prompt is the first three lines in the input examples.

*Input:*

```
Complete the following Lean 4 code:

```lean4
<lemma>
lemma sq_nonneg (a :  $\alpha$ ) : 0 ≤ a ^ 2
<easy theorem>
theorem lean_workbook_9742 (a b c :  $\mathbb{R}$ ) (ha : a ≥ 0) (hb : b ≥ 0) (hc : c ≥ 0) (hab :
  a + b + c = 1) : a^2 + b^2 + c^2 ≤ 1 / 4 + a^3 + b^3 + c^3 := by
nlinarith [sq_nonneg (a - 1 / 2), sq_nonneg (b - 1 / 2), sq_nonneg (c - 1 / 2),
mul_nonneg ha hb, mul_nonneg hb hc, mul_nonneg ha hc, hab]
<hard theorem>
theorem
```

*Output:*

```
lean_workbook_9229 (a b c :  $\mathbb{R}$ ) (ha : 0 < a ∧ 0 < b ∧ 0 < c) (hab : a + b + c = 1) :
  a^3 / (a * b) + b^3 / (b * c) + c^3 / (c * a) ≥ 1 - a^2 - b^2 - c^2</hard theorem>
```
```

In the following example, the lemma is replaced by a trivial statement, and the model can generate conjectures without focusing on any particular direction (see Section 4.1).

*Input:*

```
Complete the following Lean 4 code:

```lean4
<lemma>
theorem true: True
<easy theorem>
theorem lean_workbook_980 (a b :  $\mathbb{R}$ ) (ha : 0 < a) (hb : 0 < b) (hab : 1/a + 1/b = 1) :
  1/(a + 1) + 2/(2 * b + 1) ≤ 3/4 := by
field_simp [ha.ne', hb.ne', ha.ne', hb.ne'] at hab
field_simp [ha.ne', hb.ne', ha.ne', hb.ne']
rw [div_le_div_iff (by positivity) (by positivity)]
nlinarith [sq_nonneg (a - b), sq_nonneg (a - 2 * b), sq_nonneg (2 * a - b), hab]
<hard theorem>
theorem
```

*Output:*

```
lean_workbook_982 (a b c :  $\mathbb{R}$ ) (ha : 0 < a) (hb : 0 < b) (hc : 0 < c) (hab : a + b + c
= 1) :
a / a ^ 2 + b / b ^ 2 + c / c ^ 2 ≥ 3 / 4 ∧ 1 / (a + b + c) ^ 2 + 1 / (b + c + a) ^ 2
+ 1 / (c + a + b) ^ 2 ≥ 3 / 4</hard theorem>
```

```
'''
```

**Examples of the prover.** In the following, we show an example of the prover’s inputs and outputs.

*Input:*

Complete the following Lean 4 code:

```
'''lean4
theorem lean_workbook_52370 (p q : ℕ) (hp : Nat.Prime p) (hq : Nat.Prime q) (hpq: p
  >= q) (hq5 : q >= 5) : 24 ∣ (p^2 - q^2) := by
```

*Output:*

```
have : p^2 %
rw [Nat.pow_mod]
have : q^2 %
rw [Nat.pow_mod]
rw [← Nat.mod_add_div p 24, ← Nat.mod_add_div q 24]
ring_nf
norm_num
have : p %
have := hp.eq_one_or_self_of_dvd 2
have := hp.eq_one_or_self_of_dvd 3
have := hp.eq_one_or_self_of_dvd 4
have := hp.eq_one_or_self_of_dvd 5
omega
have : q %
have := hq.eq_one_or_self_of_dvd 2
have := hq.eq_one_or_self_of_dvd 3
have := hq.eq_one_or_self_of_dvd 4
have := hq.eq_one_or_self_of_dvd 5
omega
rcases this with (h | h | h | h | h | h | h | h) <|> rcases this with (h' | h' | h' | h' |
  h' | h' | h' | h') <|>
simp [h, h', pow_two, Nat.mul_mod, Nat.add_mod, Nat.mod_mod, Nat.mod_self, zero_add,
Nat.zero_mul, Nat.mod_one, Nat.one_mul, Nat.mod_add_mod]
<|> omega
```

## A.2 Pseudo-code for selecting the conjecturer’s inputs

In the following, we present the pseudo-code for selecting the conjecturing imports. Recall that the input for the conjecturer consists of a statement, its proof, and a lemma used in the proof (c.f., Section 3.1). In Step 1, we construct the prompts by taking the correct proofs to statements in the given dataset, and extract the lemmas used in the proof by the formal verifiers. We also allow the model to propose conjectures without focusing on any particular lemma, which is implemented by replacing the lemma statement with a fixed trivial statement in the prompt (see Appendix A.1 for concrete examples). Finally, we de-duplicate the inputs by the (statement, lemma) pair. After generating the conjectures, we randomly select a subset whose size does not exceed the number of remaining unproved statements in the given dataset, so that the prover’s sample budget is distributed equally between the conjectures and the statements.

We run two heuristic methods to ensure the diversity of the inputs. First, we make sure that each lemma  $l$  appears at most  $0.1n$  times in the inputs because we observe that some lemmas (e.g., `sq_nonneg`, `mul_self_nonneg`) are much more likely to be included. Second, we make sure that every statement-lemma pair only appear at most once in the prompt, even if there are multiple correct proofs.

Alg. 1 shows the complete pseudo-code for selecting the conjecturer’s inputs.

---

**Algorithm 1** Prepare inputs for the conjecturer.

---

- 1: **Input:** a list of statements and proofs  $L = \{(t_i, p_i)\}_{i=1}^n$ .
  - 2: Initialize prompt list  $P = []$ .
  - 3: **for**  $(t, p) \in L$  **do**
  - 4:     Parse the proof and get the set of used lemmas  $S$ .
  - 5:     With probability 0.5, add the trivial lemma to  $S$ .
  - 6:     For every lemma  $l \in S$ , add  $(t, p, l)$  to the prompt list  $P$ .
  - 7: **for**  $l \in S$  **do**
  - 8:     **if**  $\sum_{(t', p', l') \in P} \mathbb{I}[l = l'] > 0.1n$  **then**
  - 9:         Randomly keep at most  $0.1n$  prompts with lemma  $l$  in  $P$ .
  - 10: De-duplicate  $P$  randomly so that every (statement, lemma) pair  $(t, l)$  appears at most once.
  - 11: **Return:** de-duplicated list of prompts  $P$ .
- 

### A.3 Pseudo-code for preparing the conjecturer dataset.

The pseudo-code for preparing the conjecturer dataset is shown in Alg. 2. The motivations and explanations of each step in Alg. 2 can be found in Section 3.2.

---

**Algorithm 2** Prepare the conjecturer dataset.

---

- 1: **Input:** a list of (seed statement, proof of the seed statement, lemma, generated conjecture, generated proof of the conjecture) tuples  $\mathcal{D} = \{(t_i, p_i^t, l_i, c_i, p_i^c)\}_{i=1, \dots, n}$ , and unproved statements  $Q = \{t_j\}_{j=1, \dots, m}$ .
- 2: For each conjecture  $c$ , compute its empirical pass rate

$$\hat{P}(c) \triangleq (\#\{i : c_i = c, p_i^c \text{ is correct}\}) / (\#\{i : c_i = c\}).$$

- 3: Select conjecturing examples that (a) have low but positive pass rates, and (b) the lemma  $l$  is used in the proof  $p^c$ :

$$\overline{\mathcal{D}} = \{(t, p^t, l, c) \mid (t, p^t, l, c, p^c) \in \mathcal{D}, \hat{P}(c) \in (0, 1/4], \\ p^c \text{ is correct, } l \text{ is used in } p^c\}.$$

- 4: De-duplicate  $\overline{\mathcal{D}}$  based on the conjecture  $c$ .
- 5: Compute the elegance score

$$E(c) \triangleq \frac{\min\{\text{len}(p_i^c) : 1 \leq i \leq n, p_i^c \text{ is correct, } c_i = c\}}{\text{len}(c)}$$

- 6: Let  $\kappa$  be the 20%-quantile of  $E(c)$  for conjectures in  $\overline{\mathcal{D}}$ .
  - 7: Apply elegance filter:  $\tilde{\mathcal{D}} = \{(t, p^t, l, c) \in \overline{\mathcal{D}} \mid E(c) \geq \kappa\}$ .
  - 8: Find a distribution  $P$  supported on the conjectures in  $\tilde{\mathcal{D}}$  that minimizes the Wasserstein distance  $W(P, Q)$  (Alg. 4).
  - 9: **Return:**  $\tilde{\mathcal{D}}$  re-weighted by the density of  $P$ .
- 

### A.4 Pre-processing LeanWorkbook

LeanWorkbook is a dataset that contains statements translated from natural language math statements (a.k.a., auto-formalization). The original dataset contains 140K (natural language statement, formal statement) pairs.

We de-duplicate the LeanWorkbook dataset by keeping only one formal statements per natural language statement. After de-duplication, we get 89,221 formal Lean4 statements as our existing dataset w/o proofs for Lean experiments.

For the Isabelle experiments, we translate the Lean4 statements to Isabelle using DeepSeek-V2.5 API with few-shot examples. The prompt to the model is listed below.

Please translate the following lean statement into Isabelle. Please make sure that

1. All the variables are well-typed.
2. All the functions are correctly translated into the corresponding Isabelle functions.
3. All the symbols are correctly translated into corresponding Isabelle symbols.
4. Please directly output the translation without explanation.

Here are some hints for the translation:

1. In Isabelle, the second operand of the operator  $\wedge$  should be integer. For real numbers, please use `powr` instead.
2. Please define the types of numerals.
3. `'Real.logb x y'` should be translated to `'log x y'`.
4. `'Real.sqrt x'` should be translated to `'sqrt x'`.
5. Variables with subscripts should be disallowed. For any variable names of form `a_b`, translate it to `ab`.
6. Please translate superscripts to the corresponding exponential form. For example,  $x^{-1}$  should be translated to `(x powr -1)`.
7. `'a | b'` should be translated to `'a dvd b'`.
8. `'x  $\equiv$  y [ZMOD p]'` should be translated to `'x mod p = y mod p'`.
9. `'x  $\in$  zmod p'` should represent that `x` is nat and `x < p`.

```
## Input:
'''lean
theorem lean_workbook_50 (a b c : ℝ)
(ha : a ≥ 0 ∧ b ≥ 0 ∧ c ≥ 0)
(hab : a + b + c = 3)
: a^3 + b^3 + c^3 + 216 * (a * b + b * c + c * a) / (24 + a * b + b * c + c * a) ≤ 27
:= by sorry
'''

## Output:
'''Isabelle
theorem lean_workbook_50:
fixes a b c :: real
assumes "a ≥ 0 ∧ b ≥ 0 ∧ c ≥ 0"
assumes "a + b + c = 3"
shows "a^3 + b^3 + c^3 + 216 * (a * b + b * c + c * a) / (24 + a * b + b * c + c * a)
≤ 27"
sorry
'''

## Input:
'''lean
{}
'''

## Output:
```

## A.5 Re-weighting the conjecturing dataset

In this section, we describe the motivations and implementation details of the re-weighting method for the conjecturing dataset.

**Motivation.** In our early experiments, we observe that the generated conjectures tend to have mode collapse issue after several iterations of self-play training. For example, the generated conjectures are mostly about algebraic manipulations even when the seed statements contain questions about, for example, number theory. This is partly because the LeanWorkbook dataset contains a significant portion of inequality questions.

Therefore, in addition to the particular conjecturing format where we require that the proof of the conjecture must use the lemma given in the input, we also re-weight the conjecturing examples at every iteration. Intuitively, if there is a distance function that can separate statements of different topics, the Wasserstein projection of the conjectures will have a similar distribution of topics, and therefore alleviates the mode collapsing issue.

**Cost function.** We compute the cost  $d(x, y)$  of matching conjecture  $x$  to a statement  $y$  by the negative of the cosine similarity between their embeddings, and the embedding is computed by the last hidden layer of the current model averaged over the sequence dimension. Since our model is trained to generate proofs of conjectures and statements, we expect that statements with similar proof techniques tend to have similar embeddings, and therefore smaller cost for the matching.

**Algorithm.** On the high level, our method computes a re-weighting of the generated conjectures that minimizes its Wasserstein distance to the unproved statements in the given dataset. Abstractly speaking, let  $\mathcal{X}$  be the set of generated conjectures, and  $\mathcal{Q}$  the set of unproved statements. Let  $d(x, y)$  be the distance between a conjecture  $x$  and a statement  $y$ . Then, the optimization problem can be written as

$$\operatorname{argmin}_{P: P \text{ is a valid distribution, } \operatorname{supp}(P) \subseteq \mathcal{X}} W(P, Q), \quad (1)$$

where  $W(P, Q)$  is the Wasserstein distance between  $P$  and  $Q$  (with little abuse of notation, we use  $Q$  to represent the uniform distribution over the unproved statements). The Wasserstein distance  $W(P, Q)$  is defined by the following optimal transportation problem where  $\mu$  is a matching between the distribution  $P$  and  $Q$ :

$$W(P, Q) = \min_{\mu} \sum_{x \in \operatorname{supp}(P), y \in \operatorname{supp}(Q)} \mu(x, y) d(x, y) \quad (2)$$

$$\text{s.t.} \quad \sum_{y \in \operatorname{supp}(Q)} \mu(x, y) = P(x), \quad (3)$$

$$\sum_{x \in \operatorname{supp}(P)} \mu(x, y) = Q(y), \quad (4)$$

$$\mu(x, y) \geq 0, \quad \forall x, y. \quad (5)$$

Combining the equations above, the re-weighting distribution  $P$  can be computed by

$$\operatorname{argmin}_{P: \operatorname{supp}(P) \subseteq \mathcal{X}} \min_{\mu} \sum_{x \in \operatorname{supp}(P), y \in \operatorname{supp}(Q)} \mu(x, y) d(x, y) \quad (6)$$

$$\text{s.t.} \quad \sum_{y \in \operatorname{supp}(Q)} \mu(x, y) = P(x), \quad (7)$$

$$\sum_{x \in \operatorname{supp}(P)} \mu(x, y) = Q(y), \quad (8)$$

$$\mu(x, y) \geq 0, \quad \forall x, y, \quad (9)$$

$$P(x) \geq 0, \quad \forall x, \quad (10)$$

$$\sum_{x \in \mathcal{X}} P(x) = 1, \quad (11)$$

where the last two constraint ensures that  $P$  is a valid distribution. Equivalently, we get the following program,

$$\operatorname{argmin}_{P: \operatorname{supp}(P) \subseteq \mathcal{X}} \min_{\mu} \sum_{x \in \operatorname{supp}(P), y \in \operatorname{supp}(Q)} \mu(x, y) d(x, y) \quad (12)$$

$$\text{s.t.} \quad \sum_{x \in \text{supp}(P)} \mu(x, y) = Q(y), \quad (13)$$

$$\sum_{x \in \mathcal{X}, y \in \text{supp}(Q)} \mu(x, y) = 1, \quad (14)$$

$$\mu(x, y) \geq 0, \quad \forall x, y, \quad (15)$$

$$P(x) = \sum_{y \in \text{supp}(Q)} \mu(x, y). \quad (16)$$

Since  $Q(y)$  is given, we can optimize  $\mu(x, y)$  for every fixed  $y$  separately, and then compute the final  $P(x)$  using Eq. (16). As a result, the program above has a closed-form solution  $\mu^*(x, y) = Q(y) \mathbb{I}[x = \arg\min_{x' \in \mathcal{X}} d(x', y)]$  and  $P(x) = \sum_{y \in \text{supp}(Q)} \mu^*(x, y)$ . In other words, the optimal matching  $\mu(x, y)$  for any given  $y$  is only supported at the  $x$  that minimizes the distance  $d(x, y)$ . Therefore, the (theoretical) algorithm that computes the optimal re-weighting is given in Alg. 3. Note that the last line in Alg. 3 is to make sure that the sum of the weights equals the number of generated conjectures (i.e., the sum of weights before re-weighting).

---

**Algorithm 3** Computing the optimal re-weighting (theory).

---

- 1: **Input:** generated conjectures  $\mathcal{X} = \{x_1, \dots, x_n\}$  of size  $n$ , unproved statements  $Q$  with size  $m$ , and a distance function  $d(x, y)$ .
  - 2: Initialize the optimal re-weighting  $P = [0, 0, \dots, 0]$ .
  - 3: **for**  $y \in Q$  **do**
  - 4:     Compute  $x^* = \arg\min_{x \in \mathcal{X}} d(x, y)$ .
  - 5:      $P(x^*) \leftarrow P(x^*) + 1/m$ .
  - 6: **Return:** the optimal re-weighting is  $[P(x_1) * n, P(x_2) * n, \dots, P(x_n) * n]$ .
- 

In the actual implementation, we additionally requires that the weighting  $P$  for every conjecture  $x$  cannot be too big because otherwise it might cause instability of the LLM training with weighted cross entropy loss. Therefore, our practical implementation is shown in Alg. 4.

---

**Algorithm 4** Computing the optimal re-weighting.

---

- 1: **Input:** generated conjectures  $\mathcal{X} = \{x_1, \dots, x_n\}$  of size  $n$ , unproved statements  $Q$  with size  $m$ , and a distance function  $d(x, y)$ .
  - 2: Initialize the optimal re-weighting  $P = [0, 0, \dots, 0]$ .
  - 3: Initialize the masks  $M(x) = 1, \forall x \in \mathcal{X}$ .
  - 4: **for**  $y \in Q$  **do**
  - 5:     Compute  $x^* = \arg\min_{x \in \mathcal{X}} d(x, y)M(x)$ .
  - 6:      $P(x^*) \leftarrow P(x^*) + 1/m$ .
  - 7:     **if**  $P(x^*) * n > 3$  **then**
  - 8:          $M(x^*) \leftarrow 0$ .
  - 9: **Return:** the optimal re-weighting is  $[P(x_1) * n, P(x_2) * n, \dots, P(x_n) * n]$ .
- 

## A.6 Implementation details for expert iteration.

In this section, we describe two different implementations of expert iteration and compare their performance.

**Vanilla expert iteration.** For vanilla expert iteration, we only sample proofs to the *unproved* statements in the given dataset. The LLM training dataset consists of all the correct proofs generated in this and previous iterations, and in each iteration, the model is trained from the base model.



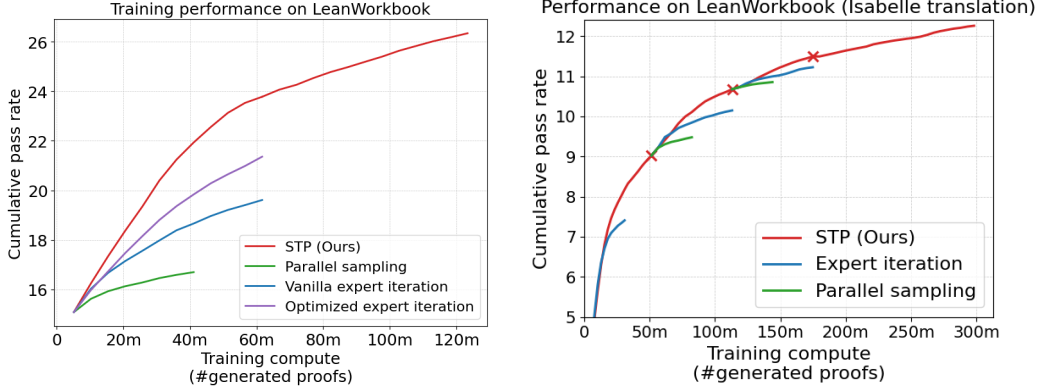


Figure 5: **Left:** Comparison of pass rates between STP, two implementations of expert iteration, and parallel sampling methods on LeanWorkbook. **Right:** Comparison of pass rates between STP and baseline methods on LeanWorkbook, starting with the SFT model. The red crosses shows the points where we refresh the self-play training as described in Section 4.1.

**Optimized expert iteration.** The most significant issue of vanilla expert iteration is the limited correct proofs generated in each iteration. As a result, even though the model is re-trained at every iteration, the difference between two models in consecutive iterations are limited.

Therefore, in our optimized implementation of expert iteration, we generate proofs to all statements in the given dataset, regardless of whether they are previously proved or not. Then, to construct the LLM training dataset, we randomly choose at most 16 proofs per statement (so that the model does not overfit to the easy problems with many correct proofs). Note that this implementation requires slightly more sample budget per iteration. However, since the pass rate on the given dataset is low (less than 30% even for our best model), this difference is not significant.

In Fig. 5 (Left), we plot the cumulative pass rate of two implementations of expert iteration, STP and parallel sampling. STP outperforms both implementations of expert iteration, and the optimized implementation of expert iteration is better than the vanilla implementation.

For the figures of Isabelle experiments, we always use the optimized implementation of expert iteration. For Fig. 2, we use the vanilla implementation.

## A.7 Additional details for interacting with the Isabelle verifier

For the Isabelle experiments, we have an additional filter for the conjectures — if the generated conjecture is equivalent to the statement in the prompt (tested by `solve_direct` in Isabelle), we consider it invalid.

We disallow the tactics `sledgehammer`, `mason`, `smt`, `metis`, `sos` by invalidating proofs that contain any of these sub-strings. However, following the implementation of Jiang et al. [2022b], we still use the keyword ‘sledgehammer’ to replace the following simple tactics

```
[by auto, by simp, by blast, by fastforce, by force, by eval, by presburger,
by arith, by linarith, by (auto simp: field_simps)].
```

During proof verification, we try these tactics sequentially to replace the keyword ‘sledgehammer’. If any of the tactics succeed, we proceed to the remaining proof steps. Otherwise we flag the proof incorrect.

## A.8 Additional details for interacting with the Lean4 verifier

During the self-play training, we use the same imports as the miniF2F Lean4 project<sup>7</sup> instead importing the entire Mathlib to optimize the memory efficiency. This is because we do not have access to an additional CPU cluster for proof verification, and the available CPU memory in TPU-v4 VMs is limited.

<sup>7</sup><https://github.com/yangky11/miniF2F-lean4/tree/main/MiniF2F>

## A.9 Compute resources

Our experiments are primarily done on TPU-v4 VMs with 32 nodes. Each node contains 4 TPU chips (8 TPU cores), 240 CPU cores, and 400G memory. We use vLLM [Kwon et al., 2023] to generate LLM outputs, and Levanter<sup>8</sup> to train the LLM. In both STP and expert iteration, since the generated proofs are heavily filtered (based on the correctness, elegance, trivialness, etc.) when constructing the training dataset, the LLM training only takes less than 25% of the wall-clock time for TPU compute, and generating proofs takes the rest 75%. In the Lean experiments, running the verifier requires 2x wall-clock time than the TPU compute. Therefore, even though we parallelize the CPU and TPU compute, the TPU is still idle for more than 50% of the wall-clock time.

## B Additional Experiment Results

In this section we show the additional experiment results with both Lean and Isabelle formal verifier.

### B.1 Additional results with Lean

In Table 3, we compare the performance of our method with prior works on PutnamBench. Note that DSP [Jiang et al., 2022a] uses Isabelle verifier where PutnamBench only has 640 statements. Our model STP achieves state-of-the-art performance by solving 9 out of 644 problems.

Table 2 compares the model obtained by final re-training with and without the proofs of generated conjectures, as discussed in the ablation study section (Section 4.4). The results show that it is still beneficial to re-train with the generated conjectures in addition to the successfully proved statements in LeanWorkbook even for performance on miniF2F-test and ProofNet-test, which it leads to about 1% performance gain (for pass@128).

Table 2: Pass rate on miniF2F and ProofNet.

Method	Sample budget	MiniF2F-test	ProofNet-test
STP (w/o conjectures)	128	56.1% $\pm$ 0.5%	16.9% $\pm$ 0.8%
STP	128	57.2% $\pm$ 0.6%	18.0% $\pm$ 0.7%

### B.2 Additional results with Isabelle

In Fig. 5 (Right), we plot the pass rates of STP and baseline methods on LeanWorkbook starting from iteration 0. The red crosses shows the points where we refresh the training process as described in Section 4.1.

Table 3: Results on PutnamBench [Tsoukalas et al., 2024].

Method	Sample budget (#Proofs / #Steps)	Result
<i>Whole-Proof Generation Methods</i>		
DSP (GPT-4o) [Jiang et al., 2022a]	10	4/640
STP	64	8/644
	3200	8/644
STP	64	7/644
w/ miniF2F-valid, ProofNet-valid	3200	<b>9/644</b>
<i>Tree Search Methods</i>		
InternLM2.5-StepProver-BF+CG	$2 \times 32 \times 600$	6/644

<sup>8</sup><https://github.com/stanford-crfm/levanter>