

SubgoalXL: SUBGOAL-BASED EXPERT LEARNING FOR THEOREM PROVING

Anonymous authors

Paper under double-blind review

ABSTRACT

Formal theorem proving, a field at the intersection of mathematics and computer science, has seen renewed interest with advancements in large language models (LLMs). This paper introduces SubgoalXL, a novel approach that synergizes subgoal-based proofs with expert learning to enhance LLMs’ capabilities in formal theorem proving within the Isabelle environment. SubgoalXL addresses two critical challenges: the scarcity of specialized mathematics and theorem-proving data, and the need for improved multi-step reasoning abilities in LLMs. By optimizing data efficiency and employing subgoal-level supervision, SubgoalXL extracts richer information from limited human-generated proofs. The framework integrates subgoal-oriented proof strategies with an expert learning system, iteratively refining formal statement, proof, and subgoal generators. Leveraging the Isabelle environment’s advantages in subgoal-based proofs, SubgoalXL achieves a new state-of-the-art performance of 56.1% in Isabelle on the standard miniF2F dataset, marking an absolute improvement of 4.9%. Notably, SubgoalXL successfully solves 41 AMC12, 9 AIME, and 3 IMO problems from miniF2F. These results underscore the effectiveness of maximizing limited data utility and employing targeted guidance for complex reasoning in formal theorem proving, contributing to the ongoing advancement of AI reasoning capabilities.

1 INTRODUCTION

Formal theorem proving, a field at the intersection of mathematics and computer science, has flourished alongside the development of languages like Lean (de Moura et al., 2015) and Isabelle (Paulson, 1994). These two prominent communities have been instrumental in advancing the field’s core challenge: mechanizing mathematical reasoning and proof verification (Li et al., 2020). Through the creation of rigorously verified proofs, this discipline strengthens the foundations of mathematical certainty, potentially opening doors to new mathematical discoveries.

The field has recently garnered renewed attention, driven by advancements in large language models (LLMs). Despite their impressive capabilities, current LLMs often face limitations in performing complex reasoning tasks required for formal theorem proving, including the need for logically rigorous, multi-step proofs (Wu et al., 2022; Jiang et al., 2022a; Zhao et al., 2024; Xin et al., 2023; Lin et al., 2024). Conventional approaches struggle to align informal human intuition with the strict formalism required by theorem-proving languages, leading to inefficiencies in generating high-quality proofs. This highlights a pressing need to refine models that not only handle the depth of logical reasoning but also make more efficient use of available data while bridging the gap between informal and formal mathematical reasoning.

In this work, we introduce SubgoalXL (Figure 1), a novel approach that synergizes subgoal-based proofs with expert learning to enhance LLMs’ capabilities in formal theorem proving. SubgoalXL tackles the scarcity of specialized mathematics and theorem-proving data (Lin et al., 2024; Wu et al., 2024) by maximizing data utility through subgoal-level decomposition of proofs, allowing for more granular supervision and iterative refinement. This approach extracts deeper structural information from human-generated proofs by focusing on intermediate subgoals, effectively breaking down the reasoning process into smaller, manageable steps. Consequently, SubgoalXL enhances multi-step reasoning abilities, ensuring that each generated subgoal aligns with both the informal intuition and the formal proof structure. At its core, SubgoalXL integrates subgoal-oriented strategies

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

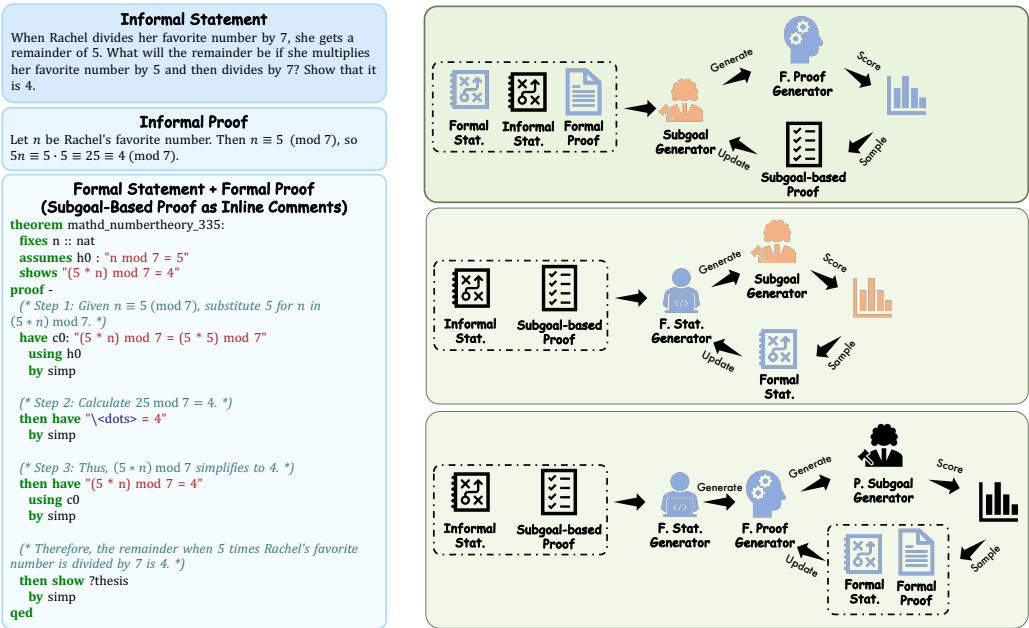


Figure 1: **Left:** Examples of informal statement, informal proof, formal statement, formal proof, and subgoal-based proof. **Right:** Overview of the subgoal-based expert learning framework. Abbreviations: “Stat.” for “Statement”, “F.” for “Formal”, and “P.” for “Posterior”. Each iteration samples subgoal-based proofs, formal statements, and formal proofs from their optimal distributions.

with an expert learning framework, refining the formal statement, proof, and subgoal generators through sampling from estimated optimal distributions, thereby improving the LLMs’ proficiency in navigating intricate logical structures and producing accurate formal proofs.

Leveraging the Isabelle environment’s advantages in subgoal-based proofs, Subgoal1XL significantly advances theorem-proving capabilities. It achieves a new state-of-the-art performance of 56.1% in Isabelle on the standard miniF2F dataset (Zheng et al., 2021), an absolute improvement of 4.9% over Zheng et al. (2023). Subgoal1XL successfully solves 41 AMC12, 9 AIME, and 3 IMO problems from miniF2F. The iterative expert learning process drives steady performance gains, underscoring Subgoal1XL’s robustness and effectiveness. These results highlight the critical role of maximizing limited data utility and employing effective guidance for complex reasoning, complementing large-scale data efforts (Wu et al., 2024; Xin et al., 2024a;b).

2 RELATED WORK

Formal theorem proving has advanced significantly through machine learning, focusing on enhancing proof search strategies and leveraging Large Language Models (LLMs) for autoformalization Polu & Sutskever (2020); Polu et al. (2022); Jiang et al. (2022a). Improvements in the proof search include self-supervised strategies in Expert Iteration (Polu et al., 2022) and PACT (Han et al., 2021), integrations of language models with automated provers in HyperTree Proof Search (HTPS)(Lample et al., 2022) and Thor(Jiang et al., 2022b), and transformer-based premise selection in Magnushammer (Mikuła et al., 2023). Despite these advancements, scalability remains a challenge due to the increasing complexity of theorems. The application of LLMs for autoformalization and proof generation has also been explored, with Wu et al. (2022) and Jiang et al. (2022a) demonstrating the conversion of mathematical problems into formal specifications. Baldur (First et al., 2023) further enhances proving capabilities by producing full proofs and incorporating a proof repair model. Additionally, LEGO-Prover (Xin et al., 2023) and Lyra Zheng et al. (2023) contribute uniquely to theorem proving by focusing on the incremental development of reusable theorems and integrating error messages from external verifiers for proof post-processing, respectively. DeepSeek-Prover (Xin

et al., 2024a) highlights the potential of large-scale synthetic data, while Lean-STaR (Lin et al., 2024) leverages informal information to boost theorem-proving capabilities by training language models to produce informal thoughts before each proof step. InternLM2-StepProver (Wu et al., 2024) addresses data scarcity by utilizing extensive formal data from Lean 4 repositories on GitHub. Zhao et al. (2024) introduce a subgoal-based demonstration learning framework that constructs and refines distinct subgoals for each example, significantly enhancing proof search efficiency in LLMs.

Nevertheless, challenges remain in addressing data scarcity and enhancing deep, multi-step reasoning in formal theorem proving. Building upon the insights from subgoal-based demonstration learning (Zhao et al., 2024), we introduce SubgoalXL, a novel framework that combines subgoal-based proofs with an expert learning system. This approach iteratively enhances formal statement, proof, and subgoal generation, aiming to improve data efficiency and achieve robust performance. SubgoalXL complements existing methods by focusing on maximizing the utility of limited data and refining complex reasoning strategies.

3 APPROACH

3.1 PROBLEM FORMALIZATION

Suppose we have an informal dataset $\mathcal{G} = \{(s_i^g, p_i^g)\}_{i=1}^{|\mathcal{G}|}$, where s_i^g is an informal statement and p_i^g is an informal proof. Similarly, we have a formal dataset $\mathcal{F} = \{(s_i^f, \mathcal{P}_i^f)\}_{i=1}^{|\mathcal{F}|}$, where s_i^f is a formal statement and \mathcal{P}_i^f is a formal proof. The goal is to train a language model $p_{\text{fpg}}(p, \mathcal{P} | s, \mathcal{S})$ using both \mathcal{G} and \mathcal{F} . Consequently, given a new informal statement s and its formal version \mathcal{S} , the model can generate both the informal proof p and the formal proof \mathcal{P} , following the distribution $p_{\text{fpg}}(p, \mathcal{P} | s, \mathcal{S})$. In this paper, we treat (p, \mathcal{P}) as a sequence of language tokens, with p representing the prefix and \mathcal{P} representing the suffix. In cases where an informal proof p is available, the model can directly generate the formal proof \mathcal{P} following $p_{\text{fpg}}(\mathcal{P} | s, \mathcal{S}, p)$.

The challenges mainly lie in (1) the limited effectiveness of informal proofs in \mathcal{G} due to discrepancies between human-written informal proofs and the established practices of formal proofs in theorem-proving languages; and (2) the difficulty in constructing the full training dataset, which requires aligned $(s, \mathcal{S}, p, \mathcal{P})$ quadruples. Inspired by Zhao et al. (2024), we use subgoal-based proofs (Figure 1a) to replace informal proofs in \mathcal{G} , achieving better consistency with the structure of formal proofs (see §3.2). Additionally, we develop an expert learning framework (Figure 1b) that samples $(s, \mathcal{S}, p, \mathcal{P})$ quadruples by estimating their optimal distributions through iterative refinement, leveraging probabilistic modeling and gradient estimation techniques (see §3.3).

3.2 SUBGOAL-BASED PROOF

To annotate subgoal-based proofs for the informal statements in \mathcal{G} , we begin by manually creating demonstration examples to serve as input for in-context learning (see Figure 1a). We select a subset of problems from the miniF2F validation set and manually construct the verified formal proof for each problem. Then, we prompt GPT-4o to generate subgoal-based proofs g , conditioned on the informal statement s , formal statement \mathcal{S} , and formal proof \mathcal{P} . This process ensures that: (1) the subgoal-based proofs are produced by autoregressive models; (2) they exhibit a consistent style, reducing the learning burden, as noted by Gu et al. (2018); and (3) each subgoal corresponds to a corresponding formal intermediate goal in Isabelle. These demonstrations are then used as in-context examples to annotate subgoal-based proofs for the informal statements in \mathcal{G} (see Appendix C for further details).

3.3 SUBGOAL-BASED EXPERT LEARNING

We introduce the SubgoalXL framework, which comprises a formal proof generator (p_{fpg}), a formal statement generator (p_{fsg}), and a subgoal generator (p_{sg}). Inspired by gradient estimation in probabilistic modeling (Schulman et al., 2015), this framework estimates optimal training data distributions

for each component and iteratively refines these components by fine-tuning on data sampled from the respective distributions.¹ The overall algorithm is presented in Algorithm 1.

Components. The core components include a formal statement generator, a formal proof generator, and a subgoal generator. The formal statement generator annotates formal statements for informal ones in \mathcal{I} following $p_{\text{fsg}}(\mathcal{S} \mid \mathcal{I})$. Subsequently, the formal proof generator produces formal proofs for the informal data in \mathcal{I} , based on $p_{\text{fpg}}(\mathcal{P} \mid \mathcal{I}, \mathcal{S}, g)$ and using subgoal-based proofs (see §3.2). The subgoal generator labels subgoal-based proofs for formal data in \mathcal{F} according to $p_{\text{sg}}(g \mid \mathcal{I}, \mathcal{S})$ after informal statements have been generated for each data point in \mathcal{F} (refer to Appendix D.2 for details).

Additionally, the formal proof generator assesses the performance of the subgoal generator by evaluating the likelihood of reconstructing formal proofs in \mathcal{F} . Conversely, the subgoal generator evaluates the formal statement generator by assessing the likelihood of reconstructing subgoal-based proofs in \mathcal{I} . We also introduce an auxiliary component, the posterior subgoal generator $p_{\text{psg}}(g \mid \mathcal{I}, \mathcal{S}, \mathcal{P})$, which evaluates the formal proof generator based on the likelihood of reconstructing subgoal-based proofs in \mathcal{I} . The formal statement generator, formal proof generator, and subgoal generator iteratively improve through expert learning, with only the formal proof generator used during testing.

Initialization. We begin by annotating the formal statements and proofs for the informal dataset \mathcal{I} using in-context learning, retaining only those verified by Isabelle. Next, we annotate the informal statements and proofs for the formal dataset \mathcal{F} in the same manner. The initial formal proof generator, denoted as $p_{\text{fpg}}^{(0)}$, is then trained on $\{(s_i^{\mathcal{I}}, \mathcal{S}_i^{\mathcal{I}}, g_i^{\mathcal{I}}, \mathcal{P}_i^{\mathcal{I}})\}_{i=1}^{|\mathcal{I}|} \cup \{(s_i^{\mathcal{F}}, \mathcal{S}_i^{\mathcal{F}}, p_i^{\mathcal{F}}, \mathcal{P}_i^{\mathcal{F}})\}_{i=1}^{|\mathcal{F}|}$. Similarly, the formal statement generator $p_{\text{fsg}}^{(0)}$ and subgoal generator $p_{\text{sg}}^{(0)}$ are trained on $\{(s_i^{\mathcal{I}}, \mathcal{S}_i^{\mathcal{I}})\}_{i=1}^{|\mathcal{I}|} \cup \{(s_i^{\mathcal{F}}, \mathcal{S}_i^{\mathcal{F}})\}_{i=1}^{|\mathcal{F}|}$ and $\{(s_i^{\mathcal{I}}, \mathcal{S}_i^{\mathcal{I}}, g_i^{\mathcal{I}})\}_{i=1}^{|\mathcal{I}|} \cup \{(s_i^{\mathcal{F}}, \mathcal{S}_i^{\mathcal{F}}, p_i^{\mathcal{F}})\}_{i=1}^{|\mathcal{F}|}$, respectively.

For training the posterior subgoal generator p_{psg} , we first obtain a version of the formal proof with all in-line comments removed, denoted as $\overline{\mathcal{P}}$. This component is trained on $\{(s_i^{\mathcal{I}}, \mathcal{S}_i^{\mathcal{I}}, \overline{\mathcal{P}}_i^{\mathcal{I}}, g_i^{\mathcal{I}})\}_{i=1}^{|\mathcal{I}|} \cup \{(s_i^{\mathcal{F}}, \mathcal{S}_i^{\mathcal{F}}, \overline{\mathcal{P}}_i^{\mathcal{F}}, p_i^{\mathcal{F}})\}_{i=1}^{|\mathcal{F}|}$. The posterior subgoal generator remains fixed during the expert learning process.

Expert Learning. Given the uncertainty in the quality of generated statements and proofs, we employ probabilistic modeling to compute the reward for each component. This allows us to derive the optimal distribution from which we sample statements and proofs in each iteration. For instance, in training the formal proof generator, the optimization objective in the k -th iteration is:

$$\max_p \mathbb{E}_{(\mathcal{S}, \mathcal{P}) \sim p} [\log p_{\text{psg}}(g \mid \mathcal{I}, \mathcal{S}, \overline{\mathcal{P}})] - \beta \mathbb{D}_{\text{KL}}[p(\mathcal{S}, \mathcal{P} \mid \mathcal{I}, g) \parallel p^{(k-1)}(\mathcal{S}, \mathcal{P} \mid \mathcal{I}, g)], \quad (1)$$

where $p(\mathcal{S}, \mathcal{P} \mid \mathcal{I}, g) = p_{\text{fpg}}(\mathcal{P} \mid \mathcal{I}, \mathcal{S}, g)p_{\text{fsg}}(\mathcal{S} \mid \mathcal{I})$ and $\log p_{\text{psg}}(g \mid \mathcal{I}, \mathcal{S}, \overline{\mathcal{P}})$ represents the reward which is derived using gradient estimators (Schulman et al., 2015). Intuitively, within the informal dataset, the formal statement and proof are treated as random variables, with optimal selections maximizing the likelihood of reconstructing the informal proof or subgoal-based proof. We also include KL-constraint terms to prevent overoptimization towards the reward. The optimal distribution of the formal proof is given by:

$$p^*(\mathcal{S}, \mathcal{P} \mid \mathcal{I}, g) = \frac{1}{Z(\mathcal{I}, g)} p^{(k-1)}(\mathcal{S}, \mathcal{P} \mid \mathcal{I}, g) \exp\left(\frac{1}{\beta} (\log p_{\text{psg}}(g \mid \mathcal{I}, \mathcal{S}, \overline{\mathcal{P}}))\right), \quad (2)$$

where $Z(\mathcal{I}, g) = \sum_{\mathcal{S}, \mathcal{P}} p^{(k-1)}(\mathcal{S}, \mathcal{P} \mid \mathcal{I}, g) \exp\left(\frac{1}{\beta} (\log p_{\text{psg}}(g \mid \mathcal{I}, \mathcal{S}, \overline{\mathcal{P}}))\right)$. The optimal distributions for formal statements $p_{\text{fsg}}^*(\mathcal{S} \mid \mathcal{I})$ and subgoal-based proofs $p_{\text{sg}}^*(g \mid \mathcal{I}, \mathcal{S})$ follow a similar pattern, as detailed in Appendix D.3.

Let $\hat{\mathcal{S}}^{(k)}$ and $(\tilde{\mathcal{S}}^{(k)}, \tilde{\mathcal{P}}^{(k)})$ be drawn from the distributions $p_{\text{fsg}}^*(\mathcal{S} \mid \mathcal{I})$ and $p^*(\mathcal{S}, \mathcal{P} \mid \mathcal{I}, g)$, respectively, for the informal dataset. Similarly, let $g^{(k)}$ be drawn from the distribution $p_{\text{sg}}^*(g \mid \mathcal{I}, \mathcal{S})$ for the formal

¹We do not include an iterative bootstrapping process for an informal statement generator, as generating informal statements from formal statements is significantly less challenging than the other three tasks. Instead, we annotate the informal statements for each formal statement in the formal dataset using in-context learning. The prompt template can be found in Appendix D.2.

dataset. In the k -th iteration, the formal statement generator updates with samples from $\{(\mathcal{J}, \hat{\mathcal{S}}^{(k)})\}$, while the formal proof generator is trained on $\{(\mathcal{J}, \hat{\mathcal{S}}^{(k)}, \mathcal{g}, \tilde{\mathcal{P}}^{(k)})\} \cup \{(\mathcal{J}, \mathcal{S}, g^{(k)}, \mathcal{P})\}$. Simultaneously, the subgoal generator refines its parameters using $\{(\mathcal{J}, \hat{\mathcal{S}}, g^{(k)})\}$. These updates are augmented by the corresponding training data generated during the initialization phase, contributing to increased data diversity and model robustness throughout training.

Diversity Efforts. We employ various strategies to enhance the diversity of model outputs, thereby improving the efficiency of the search process. (1) During the initialization phase, we train four distinct language models for the formal proof generator. These models are derived from combinations of two prompt templates (see Appendix D.1) and two proof lengths. For the proof lengths, one model retains the entire dataset, while the other selectively excludes shorter proofs based on indicators drawn from Bernoulli distributions.² (2) In each iteration of the expert learning phase, we reinitialize the components from the Llama-3-8B rather than from the previous iteration’s checkpoints.

Algorithm 1 Subgoal-based Expert Learning

Requires: \mathcal{G} : informal dataset ($\mathcal{G} = \{\mathcal{J}_i^g, p_i^g\}_{i=1}^{|\mathcal{G}|}$).
 \mathcal{F} : formal dataset ($\mathcal{F} = \{\mathcal{S}_i^f, \mathcal{P}_i^f\}_{i=1}^{|\mathcal{F}|}$).
 K_{\max} : maximum iterations for expert learning.
 m : sample size in expert learning.

1: $\mathcal{D}_{\text{fsg}}^{(0)}, \mathcal{D}_{\text{fpg}}^{(0)}, \mathcal{D}_{\text{sg}}^{(0)}, \mathcal{D}_{\text{psg}} \leftarrow \emptyset$ ▷ Initialize datasets for training all components
2: **for** $i = 1$ to $|\mathcal{G}|$ **do**
3: Annotate subgoal-based proof g_i^g , formal statement \mathcal{S}_i^g , and formal proof \mathcal{P}_i^g for (\mathcal{J}_i^g, p_i^g) .
4: Remove inline comments in \mathcal{P}_i^g to obtain $\overline{\mathcal{P}}_i^g$.
5: Update $\mathcal{D}_{\text{fsg}}^{(0)} \leftarrow \mathcal{D}_{\text{fsg}}^{(0)} \cup \{(\mathcal{J}_i^g, \mathcal{S}_i^g)\}$ and $\mathcal{D}_{\text{fpg}}^{(0)} \leftarrow \mathcal{D}_{\text{fpg}}^{(0)} \cup \{(\mathcal{J}_i^g, \mathcal{S}_i^g, g_i^g, \mathcal{P}_i^g)\}$.
6: Update $\mathcal{D}_{\text{sg}}^{(0)} \leftarrow \mathcal{D}_{\text{sg}}^{(0)} \cup \{(\mathcal{J}_i^g, \mathcal{S}_i^g, g_i^g)\}$ and $\mathcal{D}_{\text{psg}} \leftarrow \mathcal{D}_{\text{psg}} \cup \{(\mathcal{J}_i^g, \mathcal{S}_i^g, \overline{\mathcal{P}}_i^g, g_i^g)\}$.
7: **end for**
8: **for** $i = 1$ to $|\mathcal{F}|$ **do**
9: Annotate informal statement \mathcal{J}_i^f and informal proof p_i^f for $(\mathcal{S}_i^f, \mathcal{P}_i^f)$.
10: Update $\mathcal{D}_{\text{fsg}}^{(0)}, \mathcal{D}_{\text{fpg}}^{(0)}, \mathcal{D}_{\text{sg}}^{(0)}$, and \mathcal{D}_{psg} accordingly.
11: **end for**
12: Fine-tune models to obtain $p_{\text{fsg}}^{(0)}, p_{\text{fpg}}^{(0)}, p_{\text{sg}}^{(0)}$, and p_{psg} using $\mathcal{D}_{\text{fsg}}^{(0)}, \mathcal{D}_{\text{fpg}}^{(0)}, \mathcal{D}_{\text{sg}}^{(0)}$, and \mathcal{D}_{psg} respectively.
13: **for** $k = 1$ to K_{\max} **do** ▷ Begin expert learning iterations
14: $\mathcal{D}_{\text{fsg}}^{(k)} \leftarrow \mathcal{D}_{\text{fsg}}^{(0)}, \mathcal{D}_{\text{fpg}}^{(k)} \leftarrow \mathcal{D}_{\text{fpg}}^{(0)}, \mathcal{D}_{\text{sg}}^{(k)} \leftarrow \mathcal{D}_{\text{sg}}^{(0)}$.
15: **for** $i = 1$ to $|\mathcal{G}|$ **do**
16: **for** $j = 1$ to m **do**
17: Sample $\mathcal{S}_{i,j}^{(k)}$ according to Eq.3, then update $\mathcal{D}_{\text{fsg}}^{(k)} \leftarrow \mathcal{D}_{\text{fsg}}^{(k)} \cup \{(\mathcal{J}_i^g, \mathcal{S}_{i,j}^{(k)})\}$.
18: Sample $(\mathcal{S}_{i,j}^{(k)}, \mathcal{P}_{i,j}^{(k)})$ according to Eq.2, then update $\mathcal{D}_{\text{fpg}}^{(k)} \leftarrow \mathcal{D}_{\text{fpg}}^{(k)} \cup \{(\mathcal{J}_i^g, \mathcal{S}_{i,j}^{(k)}, g_{i,j}^g, \mathcal{P}_{i,j}^{(k)})\}$.
19: **end for**
20: **end for**
21: **for** $i = 1$ to $|\mathcal{F}|$ **do**
22: **for** $j = 1$ to m **do**
23: Sample $g_{i,j}^{(k)}$ according to Eq.4.
24: Update $\mathcal{D}_{\text{sg}}^{(k)} \leftarrow \mathcal{D}_{\text{sg}}^{(k)} \cup \{(\mathcal{J}_i^g, \mathcal{S}_i^g, g_{i,j}^{(k)})\}$ and $\mathcal{D}_{\text{fpg}}^{(k)} \leftarrow \mathcal{D}_{\text{fpg}}^{(k)} \cup \{(\mathcal{J}_i^g, \mathcal{S}_i^g, g_{i,j}^{(k)}, \mathcal{P}_i^f)\}$.
25: **end for**
26: **end for**
27: Fine-tune models to obtain $p_{\text{fsg}}^{(k)}, p_{\text{fpg}}^{(k)}$, and $p_{\text{sg}}^{(k)}$ using $\mathcal{D}_{\text{fsg}}^{(k)}, \mathcal{D}_{\text{fpg}}^{(k)}$, and $\mathcal{D}_{\text{sg}}^{(k)}$ respectively.
28: **end for**

4 EXPERIMENTS

4.1 DATASET AND EVALUATION

Dataset. We evaluate our approach using the miniF2F dataset (Zheng et al., 2021), which includes 488 formal mathematical problems from high-school competitions, expressed in three formal lan-

²For formal proofs with lengths 1, 2, and 3, the drop rates are 0.8, 0.6, and 0.4, respectively.

270 guages: Lean, HOL-Light, and Isabelle. The dataset is split into a validation set and a test set, each
 271 containing 244 problems. These problems come from three different sources: 260 problems are from
 272 the MATH dataset (Hendrycks et al., 2021), 160 problems are from real high-school mathematical
 273 competitions (AMC, AIME, and IMO), and 68 problems are designed to match the difficulty level of
 274 these competitions.

275
 276 **Evaluation.** The task involves generating formal sketches for problems in the miniF2F dataset.
 277 The validity of a formal sketch must meet two criteria: it should not contain “cheating” keywords
 278 like “sorry” and “oops” that end a proof prematurely, and it must be verifiable by the interactive
 279 theorem prover Isabelle. To facilitate working with Isabelle, we use the Portal-to-Isabelle API
 280 introduced by Jiang et al. (2022a). We use the pass rate to measure our results, reporting it for both
 281 the miniF2F-valid set and the miniF2F-test set. Further details about the formal environments are
 282 provided in Appendix A.

283 4.2 BASELINES

284 To assess the performance of our approach, we compare it against several established baselines.

285
 286 **Symbolic Automated Provers.** We first apply Sledgehammer, a proof automation tool extensively
 287 used within the Isabelle environment. Sledgehammer incorporates a 120-second timeout and utilizes
 288 five automated theorem provers (Z3, CVC4, SPASS, Vampire, E). Following Jiang et al. (2022a), we
 289 enhance Sledgehammer with a set of 11 common tactics (e.g., auto, simp, blast, fastforce, force, eval,
 290 presburger, sos, arith, linarith, auto simp: field_simps). If these tactics fail or take longer than 10
 291 seconds, the system defaults to the basic Sledgehammer configuration.

292
 293 **Search-based Approaches.** We also employ search-based methods, particularly Monte-Carlo tree
 294 search (Silver et al., 2016), to explore proof possibilities. This includes Thor (Jiang et al., 2022b)
 295 and a version enhanced with expert iteration on autoformalized data (Thor+expert iteration (Wu
 296 et al., 2022)). Thor integrates language models with automated theorem provers to efficiently
 297 select premises from large libraries, while Thor+expert iteration further refines this by training on
 298 autoformalized theorems.

299
 300 **LLM-based Approaches.** In the LLM-based category, we evaluate several frameworks: Draft,
 301 Sketch, and Prove (DSP) (Jiang et al., 2022a), LEGO-Prover (Xin et al., 2023), Lyra (Zheng et al.,
 302 2023), and Subgoal-Prover (Zhao et al., 2024). DSP uses the 540B Minerva model (Lewkowycz
 303 et al., 2022) to generate formal sketches from informal proofs. LEGO-Prover incrementally develops
 304 reusable theorems to enhance proof efficiency, while Lyra integrates feedback from external verifiers
 305 to optimize the verification process. Subgoal-Prover improves LLM performance in formal theorem
 306 proving by replacing informal proofs with subgoal-based proofs and using diffusion models to
 307 organize demonstrations optimally. Notably, all these methods employ Sledgehammer for consistency
 308 across evaluations.

309
 310 Comparisons with theorem proving methods based on Lean (de Moura et al., 2015), a system utilizing
 311 distinct tactics and automation mechanisms that are not directly comparable to Isabelle, are deferred
 312 to Appendix B for thorough analysis.

313 4.3 IMPLEMENTATION DETAILS

314
 315 We collected past AMC8, AMC10, AMC12, AIME, and IMO problems from the AOPS website ³
 316 and combined them with training data from GSM8K (Cobbe et al., 2021) and MATH (Hendrycks
 317 et al., 2021) to build the informal dataset. The formal dataset was constructed using the AFP-
 318 2021 ⁴ library and the HOL library from Isabelle 2021 ⁵. This resulted in a total of 18k
 319 ⟨Informal Statement, Informal Proof⟩ pairs and 195k ⟨Formal Statement, Formal Proof⟩ pairs. To
 320

321 ³<https://artofproblemsolving.com/community>

322 ⁴<https://www.isa-afp.org/release/afp-2021-10-22.tar.gz>

323 ⁵https://isabelle.in.tum.de/website-Isabelle2021/dist/Isabelle2021_linux.tar.gz

Table 1: Performance on the miniF2F dataset. Methods marked with [†] incorporate human-written informal proofs either fully or partially during the proof search process. Bold numbers denote the highest performance achieved.

Model	Base	miniF2F-valid	miniF2F-test
Sledgehammer	-	9.9%	10.4%
Sledgehammer+heuristic	-	18.0%	20.9%
Thor (Jiang et al., 2022b)	-	28.3%	29.9%
Thor + expert iteration (Wu et al., 2022)	-	37.3%	35.2%
DSP (Jiang et al., 2022a) [†]	Codex	42.6%	39.3%
Subgoal-Prover (Zhao et al., 2024)	GPT-3.5-Turbo	48.0%	45.5%
LEGO-Prover (Xin et al., 2023) [†]	GPT-3.5-Turbo	55.3%	50.0%
Lyra (Zheng et al., 2023) [†]	GPT-4	55.3%	51.2%
SubgoalXL (ours) [†]	Llama-3-8B	61.9%	56.1%

prevent data leakage, we filtered out problems that had more than 10% 3-gram overlap with those from the miniF2F dataset.

For the initialization phase, we employed a mixture of deepseek-math-base and Llama-3-8B, with a maximum generation length of 2048 tokens and temperature settings of 0.6 and 0.8. This yielded 27k quadruples for the informal dataset and 174k quadruples for the formal dataset. The training of Llama-3-8B was performed with a learning rate of 1e-5 over 3 epochs, utilizing a sequence length of 8192 tokens. These hyperparameters were also applied during the expert learning phase. All training was performed on a single SN20 node.

For the expert learning phase, we retained 11k problems from the informal dataset (after excluding GSM8K problems) and 10k problems from the formal dataset (after selecting 10k problems from the HOL library). The maximum number of expert learning iterations, K_{\max} , was set to 3, with a sample size m of 2. At each iteration, we trained 4 formal proof generators using combinations of two prompt templates and two proof lengths, leading to a total of 16 models after 3 iterations. The number of verified proofs generated during each iteration was 3156, 3592, and 4117, respectively. Adding the 27k verified proofs obtained during the initialization phase, a total of 38k verified proofs were generated.

For inference, each model generated 512 samples with and without human-written informal proofs, resulting in a total of 16384 proof attempts across all iterations for the miniF2F dataset. This includes 8192 attempts with human-written informal proofs and 8192 attempts without them. The inference process was executed across 4 SN40 nodes.

Verification was carried out using both Isabelle 2021 and Isabelle 2022. A formal proof was deemed correct if it passed verification in either version of Isabelle. The verification process was conducted on 2048 CPUs.

4.4 MAIN RESULTS

Our main experimental results, as shown in Table 1, highlight several important findings: (1) SubgoalXL achieves the best performance, setting a new state-of-the-art with 56.1% on the miniF2F-test dataset, surpassing previous methods by an absolute improvement of up to 4.9%. (2) The success of both SubgoalXL and Subgoal-Prover emphasizes the effectiveness of subgoal-based proofs in enhancing the capabilities of large language models in formal theorem proving. (3) The benefits of expert iteration are evident, as demonstrated by the performance gains of Thor + expert iteration and SubgoalXL, reinforcing the value of iterative refinement in boosting theorem proving accuracy.

Table 2: Ablation study results on the miniF2F dataset.

Model	miniF2F-valid	miniF2F-test
SubgoalXL	46.3%	39.3%
-subgoal	34.8%	36.5%

5 ANALYSIS

5.1 ABLATION STUDY

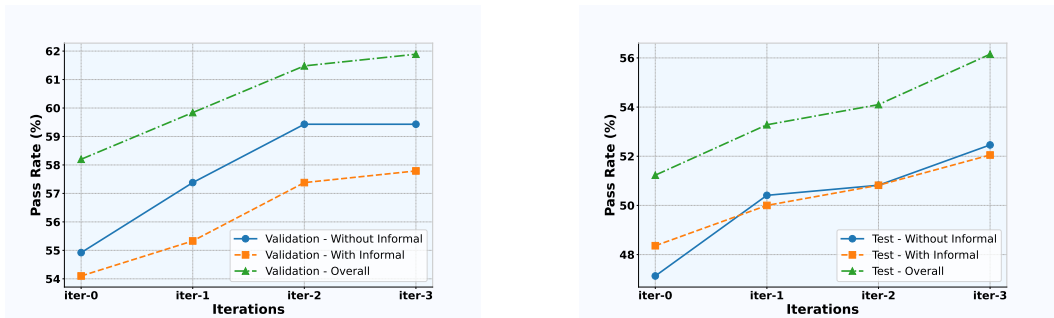
In our study, we conducted ablation experiments on our proposed model using a search budget of 64 to assess the impact of the subgoal-based framework. We evaluated two configurations: the complete model and a variant without subgoal-based proofs (-subgoal). Results in Table 2 demonstrate the importance of the subgoal-based component, as removing it (-subgoal) led to a significant decrease in performance. Specifically, the full model achieved 46.3% on the miniF2F-valid and 39.3% on the miniF2F-test, whereas the -subgoal variant saw a reduction to 34.8% on miniF2F-valid and 36.5% on miniF2F-test.

Table 3: Impact of human-written informal proofs on the performance of SubgoalXL on the miniF2F dataset. The miniF2F benchmark includes human-written informal proofs for each problem, provided by the benchmark’s publishers.

Model	miniF2F-valid	miniF2F-test
SubgoalXL (w/o informal proof)	59.4%	52.5%
SubgoalXL (with informal proof)	57.8%	52.1%

5.2 IMPACT OF HUMAN-WRITTEN INFORMAL PROOFS

We investigated the effect of human-written informal proofs on the performance of our model by conducting experiments with and without these proofs, using a search budget of 8192. Table 3 presents the results on the miniF2F-valid and miniF2F-test datasets. Our model without informal proofs achieved 59.4% on miniF2F-valid and 52.5% on miniF2F-test, while the version incorporating informal proofs reached 57.8% on miniF2F-valid and 52.1% on miniF2F-test. These results suggest that the inclusion of human-written informal proofs does not significantly enhance the model’s performance. Our model’s generation of subgoal-based proofs appears to be more effective than utilizing informal proofs in certain scenarios (refer to §5.6 for detailed examples).



(a) Validation pass-rate over iterations

(b) Test pass-rate over iterations

Figure 2: Pass-rate comparisons across different iterations on the miniF2F dataset.

5.3 ITERATIVE PERFORMANCE ANALYSIS

To evaluate our model’s iterative improvement, we conducted experiments with and without human-written informal proofs, tracking validation and test pass rates over several iterations in the expert learning process. Figures 2a and 2b present these pass rates across four iterations. In the miniF2F-valid split (Figure 2a), the model without informal proofs began at 54.92% in iteration 0 and plateaued at 59.43% by iteration 2, maintaining this performance in iteration 3. The model with informal proofs started at 54.10%, peaking at 57.79% in iteration 3. Overall validation performance increased consistently from 58.20% in iteration 0 to 61.89% in iteration 3. In the miniF2F-test split (Figure 2b), the model without informal proofs improved from 47.13% in iteration 0 to 52.46% in iteration 3, while the model with informal proofs started at 48.36% and reached 52.05% by iteration 3. Overall test performance increased from 51.23% in iteration 0 to 56.15% in iteration 3. These results indicate that our subgoal-based framework drives iterative performance improvements, with the exclusion of informal proofs often yielding better results.

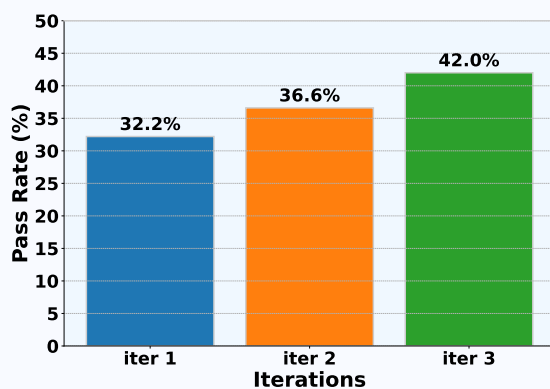


Figure 3: Synthetic proof pass-rate over iterations.

5.4 SYNTHETIC PROOF PASS RATE ANALYSIS

We analyzed the pass rates of synthetic proofs over three iterations to evaluate the iterative learning process. The results, depicted in Figure 3, show a steady increase in performance. In iteration 1, the pass rate was 32.18%. This improved to 36.63% in iteration 2 and further to 41.98% in iteration 3. These results indicate a consistent improvement in the generation of synthetic proofs as the iterations progress, highlighting the effectiveness of the iterative learning framework in enhancing the model’s proof generation capabilities.

5.5 ERROR ANALYSIS IN PROOF GENERATION

To gain insights into the errors encountered during proof generation, we categorized and quantified various error types. The results, depicted in Figure 4, reveal the frequency of each error category. The most prevalent error was “Outer syntax error” with 1, 510, 737 occurrences, followed by “Failed to finish proof” (127, 082), and “Undefined fact” (124, 611). Other notable errors included “Type unification failed” (90, 664), “Timeout” (74, 459), and “Failed to apply initial proof method“ (58, 659). This detailed error analysis highlights common failure points in the proof generation process, providing a clear direction for targeted improvements.

5.6 CASE STUDY

We evaluated the effectiveness of subgoal-based proofs versus informal proofs using a specific theorem. As shown in Figure 5, the leftmost example represents a successful proof using subgoal-based methods, while the other examples depict failed attempts using informal proofs. The subgoal-based proof demonstrated robustness and effectiveness, whereas the informal proof attempts failed to sufficiently establish the necessary conditions, leading to incomplete proofs.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

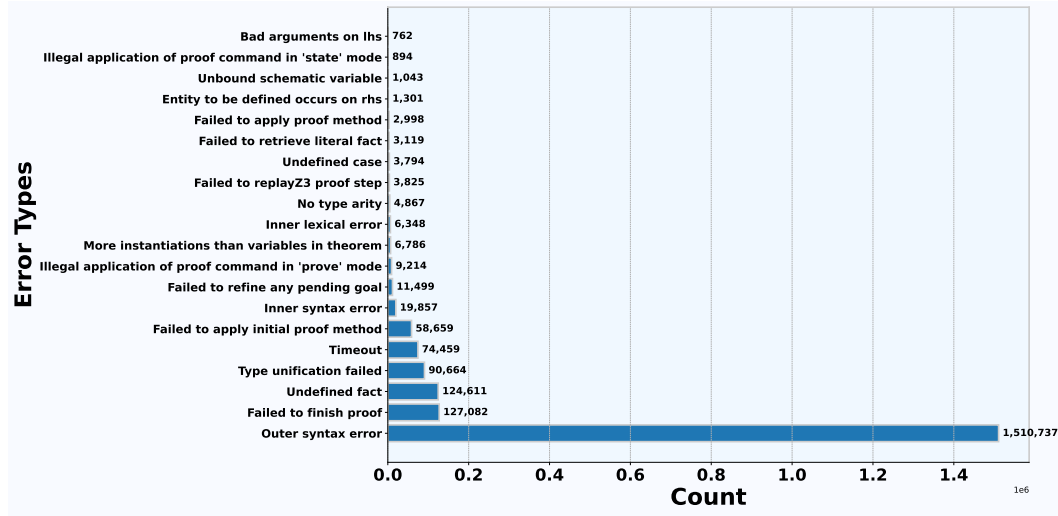


Figure 4: Counts of Different Error Types

Success (w/ Subgoal-Based Proof)	Failure 1 (w/ Informal Proof)	Failure 2 (w/ Informal Proof)
<pre> theorem math.numbertheory_435: fixes k :: nat assumes "0 < k" and "\forall r. r < k. gcd (6 * n + k) (6 * n + 3) = 1" and "\forall r. r < k. gcd (6 * n + k) (6 * n + 2) = 1" and "\forall r. r < k. gcd (6 * n + k) (6 * n + 1) = 1" shows "\exists l. l < k" proof - (* We will show that for any positive integer n, 6n + k is relatively prime to each of 6n + 3, 6n + 2, and 6n + 1 (This is true by assumption). *) have "\forall r. r < k. gcd (6 * n + k) (6 * n + 3) = 1" using assms(2) sledgehammer moreover have "\forall r. r < k. gcd (6 * n + k) (6 * n + 2) = 1" using assms(3) sledgehammer moreover have "\forall r. r < k. gcd (6 * n + k) (6 * n + 1) = 1" using assms(4) sledgehammer (* We will show that for any positive integer n, 6n + k is relatively prime to each of 6n + 3, 6n + 2, and 6n + 1 (This is true by assumption). *) ultimately have "\forall r. r < k. coprime (6 * n + k) (6 * n + 3) \and coprime (6 * n + k) (6 * n + 2) \and coprime (6 * n + k) (6 * n + 1)" sledgehammer (* Now we will show that k is at least 5. Suppose that k < 5. Then k = 1, 2, 3, 4. We will show that for each of these values of k, there exists a positive integer n such that 6n + k is not relatively prime to one of 6n + 3, 6n + 2, and 6n + 1. *) have "\exists n. n < 5" proof (rule ccontr) assume "\forall n. n < 5" then have "k < 5" sledgehammer with assms(1) consider "k = 1" "k = 2" "k = 3" "k = 4" sledgehammer then show False proof cases case 1 then have "gcd (6 * 2 + k) (6 * 2 + 1) < 6" sledgehammer then show False sledgehammer next case 2 then have "gcd (6 * 5 + k) (6 * 5 + 3) < 6" sledgehammer then show False sledgehammer next case 3 then have "gcd (6 * 4 + k) (6 * 4 + 2) < 6" sledgehammer then show False sledgehammer next case 4 then have "gcd (6 * 2 + k) (6 * 2 + 3) < 6" sledgehammer then show False sledgehammer qed qed (* Thus k is at least 5. *) then show thesis sledgehammer qed </pre>	<pre> theorem math.numbertheory_435: fixes k :: nat assumes "0 < k" and "\forall r. r < k. gcd (6 * n + k) (6 * n + 3) = 1" and "\forall r. r < k. gcd (6 * n + k) (6 * n + 2) = 1" and "\forall r. r < k. gcd (6 * n + k) (6 * n + 1) = 1" shows "\exists l. l < k" proof - (* Obviously, we have that k > 3, because otherwise two of the integers would be identical and not be relatively prime. *) have "k > 3" using assms sledgehammer (* Start by testing k = 4. 6n + 4 and 6n + 3 are relatively prime because they are consecutive integers, but 6n + 4 and 6n + 2 are both even and are therefore not relatively prime. *) have "\exists n. n < 4" using assms sledgehammer (* The next candidate to test is k = 5. Firstly, we have that gcd(6n + 5, 6n + 3) = gcd(6n + 3, (6n + 5) - (6n + 3)) = gcd(6n + 3, 2) Since 6n + 3 is always odd, the two integers 6n + 5 and 6n + 3 are relatively prime. Secondly, gcd(6n + 5, 6n + 2) = gcd(6n + 2, (6n + 5) - (6n + 2)) = gcd(6n + 2, 3) Note that 6n + 3 is always divisible by 3, so 6n + 2 is never divisible by 3. As a result, we have that 6n + 5 and 6n + 2 are relatively prime. Finally, gcd(6n + 5, 6n + 1) = gcd(6n + 1, (6n + 5) - (6n + 1)) = gcd(6n + 1, 4) Note that 6n + 1 is always odd, so 6n + 5 and 6n + 1 are also relatively prime. Therefore, the smallest positive integer k that permits 6n + k to be relatively prime with each of 6n + 3, 6n + 2, and 6n + 1 is k = 5. *) have "k > 5" using assms sledgehammer then show thesis sledgehammer qed </pre>	<pre> theorem math.numbertheory_435: fixes k :: nat assumes "0 < k" and "\forall r. r < k. gcd (6 * n + k) (6 * n + 3) = 1" and "\forall r. r < k. gcd (6 * n + k) (6 * n + 2) = 1" and "\forall r. r < k. gcd (6 * n + k) (6 * n + 1) = 1" shows "\exists l. l < k" proof - (* Obviously, we have that k > 3, because otherwise two of the integers would be identical and not be relatively prime. *) have "k > 3" using assms sledgehammer (* Start by testing k = 4. 6n + 4 and 6n + 3 are relatively prime because they are consecutive integers, but 6n + 4 and 6n + 2 are both even and are therefore not relatively prime. The next candidate to test is k = 5. *) have "\exists n. n < 5" using assms sledgehammer (* Firstly, we have that gcd(6n + 5, 6n + 3) = gcd(6n + 3, (6n + 5) - (6n + 3)) = gcd(6n + 3, 2) Since 6n + 3 is always odd, the two integers 6n + 5 and 6n + 3 are relatively prime. Secondly, gcd(6n + 5, 6n + 2) = gcd(6n + 2, (6n + 5) - (6n + 2)) = gcd(6n + 2, 3) Note that 6n + 3 is always divisible by 3, so 6n + 2 is never divisible by 3. As a result, we have that 6n + 5 and 6n + 2 are relatively prime. Finally, gcd(6n + 5, 6n + 1) = gcd(6n + 1, (6n + 5) - (6n + 1)) = gcd(6n + 1, 4) Note that 6n + 1 is always odd, so 6n + 5 and 6n + 1 are also relatively prime. Therefore, the smallest positive integer k that permits 6n + k to be relatively prime with each of 6n + 3, 6n + 2, and 6n + 1 is k = 5. *) then show thesis sledgehammer qed </pre>

Figure 5: Case study comparing subgoal-based and informal proofs. The left example shows a successful attempt using subgoal-based proofs, while the right examples depict failed attempts with informal proofs.

6 CONCLUSION

In conclusion, SubgoalXL marks a significant step forward in AI-powered theorem proving within the Isabelle environment. By addressing the challenges of complex multi-step reasoning, SubgoalXL demonstrates the efficacy of integrating subgoal-based proofs with an expert learning framework. This method iteratively refines three key components: a formal statement generator, a formal proof generator, and a subgoal generator, leading to improved performance on theorem-proving tasks. The empirical results confirm the effectiveness of SubgoalXL, achieving state-of-the-art performance on the standard miniF2F dataset with a score of 56.1%. This work paves the way for further innovations in applying AI to tackle advanced mathematical challenges in formal theorem proving.

REFERENCES

- 540
541
542 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
543 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve
544 math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 545 Leonardo de Moura, Soonho Kong, Jeremy Avigad, Floris Van Doorn, and Jakob von Raumer. The
546 lean theorem prover (system description). In *Automated Deduction-CADE-25: 25th International
547 Conference on Automated Deduction, Berlin, Germany, August 1-7, 2015, Proceedings 25*, pp.
548 378–388. Springer, 2015.
- 549 Emily First, Markus N Rabe, Talia Ringer, and Yuriy Brun. Baldur: Whole-proof generation and
550 repair with large language models. *arXiv preprint arXiv:2303.04910*, 2023.
- 551 Jiatao Gu, James Bradbury, Caiming Xiong, Victor O.K. Li, and Richard Socher. Non-autoregressive
552 neural machine translation. In *International Conference on Learning Representations*, 2018. URL
553 <https://openreview.net/forum?id=B118Bt1Cb>.
- 554 Jesse Michael Han, Jason Rute, Yuhuai Wu, Edward W Ayers, and Stanislas Polu. Proof artifact
555 co-training for theorem proving with language models. *arXiv preprint arXiv:2102.06203*, 2021.
- 556 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
557 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv
558 preprint arXiv:2103.03874*, 2021.
- 559 Albert Q Jiang, Sean Welleck, Jin Peng Zhou, Wenda Li, Jiacheng Liu, Mateja Jamnik, Timothée
560 Lacroix, Yuhuai Wu, and Guillaume Lample. Draft, sketch, and prove: Guiding formal theorem
561 provers with informal proofs. *arXiv preprint arXiv:2210.12283*, 2022a.
- 562 Albert Qiaochu Jiang, Wenda Li, Szymon Tworowski, Konrad Czechowski, Tomasz Odrzygóźdź,
563 Piotr Miłoś, Yuhuai Wu, and Mateja Jamnik. Thor: Welding hammers to integrate language
564 models and automated theorem provers. *Advances in Neural Information Processing Systems*, 35:
565 8360–8373, 2022b.
- 566 Guillaume Lample, Timothee Lacroix, Marie-Anne Lachaux, Aurelien Rodriguez, Amaury Hayat,
567 Thibaut Lavril, Gabriel Ebner, and Xavier Martinet. Hypertree proof search for neural theorem
568 proving. *Advances in Neural Information Processing Systems*, 35:26337–26349, 2022.
- 569 Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ra-
570 masesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative
571 reasoning problems with language models. *arXiv preprint arXiv:2206.14858*, 2022.
- 572 Wenda Li, Lei Yu, Yuhuai Wu, and Lawrence C Paulson. Modelling high-level mathematical
573 reasoning in mechanised declarative proofs. *arXiv preprint arXiv:2006.09265*, 2020.
- 574 Haohan Lin, Zhiqing Sun, Yiming Yang, and Sean Welleck. Lean-star: Learning to interleave thinking
575 and proving. *arXiv preprint arXiv:2407.10040*, 2024.
- 576 Maciej Miłoś, Szymon Antoniak, Szymon Tworowski, Albert Qiaochu Jiang, Jin Peng Zhou,
577 Christian Szegedy, Łukasz Kuciński, Piotr Miłoś, and Yuhuai Wu. Magnushammer: A transformer-
578 based approach to premise selection. *arXiv preprint arXiv:2303.04488*, 2023.
- 579 Lawrence C Paulson. *Isabelle: A generic theorem prover*. Springer, 1994.
- 580 Lawrence C Paulsson and Jasmin C Blanchette. Three years of experience with sledgehammer,
581 a practical link between automatic and interactive theorem provers. In *Proceedings of the 8th
582 International Workshop on the Implementation of Logics (IWIL-2010)*, Yogyakarta, Indonesia.
583 *EPiC*, volume 2, 2012.
- 584 Stanislas Polu and Ilya Sutskever. Generative language modeling for automated theorem proving.
585 *arXiv preprint arXiv:2009.03393*, 2020.
- 586 Stanislas Polu, Jesse Michael Han, Kunhao Zheng, Mantas Baksys, Igor Babuschkin, and Ilya
587 Sutskever. Formal mathematics statement curriculum learning. *arXiv preprint arXiv:2202.01344*,
588 2022.

- 594 John Schulman, Nicolas Heess, Theophane Weber, and Pieter Abbeel. Gradient estimation using
595 stochastic computation graphs. *Advances in neural information processing systems*, 28, 2015.
596
- 597 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche,
598 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering
599 the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- 600 Yuhuai Wu, Albert Qiaochu Jiang, Wenda Li, Markus Rabe, Charles Staats, Mateja Jamnik, and
601 Christian Szegedy. Autoformalization with large language models. *Advances in Neural Information
602 Processing Systems*, 35:32353–32368, 2022.
603
- 604 Zijian Wu, Jiayu Wang, Dahua Lin, and Kai Chen. Lean-github: Compiling github lean repositories
605 for a versatile lean prover. *arXiv preprint arXiv:2407.17227*, 2024.
- 606 Huajian Xin, Haiming Wang, Chuanyang Zheng, Lin Li, Zhengying Liu, Qingxing Cao, Yinya Huang,
607 Jing Xiong, Han Shi, Enze Xie, et al. Lego-prover: Neural theorem proving with growing libraries.
608 *arXiv preprint arXiv:2310.00656*, 2023.
609
- 610 Huajian Xin, Daya Guo, Zhihong Shao, Zhizhou Ren, Qihao Zhu, Bo Liu, Chong Ruan, Wenda Li,
611 and Xiaodan Liang. Deepseek-prover: Advancing theorem proving in llms through large-scale
612 synthetic data. *arXiv preprint arXiv:2405.14333*, 2024a.
- 613 Huajian Xin, ZZ Ren, Junxiao Song, Zhihong Shao, Wanxia Zhao, Haocheng Wang, Bo Liu, Liyue
614 Zhang, Xuan Lu, Qiushi Du, et al. Deepseek-prover-v1. 5: Harnessing proof assistant feedback for
615 reinforcement learning and monte-carlo tree search. *arXiv preprint arXiv:2408.08152*, 2024b.
616
- 617 Xueliang Zhao, Wenda Li, and Lingpeng Kong. Subgoal-based demonstration learning for formal
618 theorem proving. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria
619 Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International
620 Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp.
621 60832–60865. PMLR, 21–27 Jul 2024. URL [https://proceedings.mlr.press/v235/
zhao24h.html](https://proceedings.mlr.press/v235/zhao24h.html).
622
- 623 Chuanyang Zheng, Haiming Wang, Enze Xie, Zhengying Liu, Jiankai Sun, Huajian Xin, Jianhao
624 Shen, Zhenguo Li, and Yu Li. Lyra: Orchestrating dual correction in automated theorem proving.
625 *arXiv preprint arXiv:2309.15806*, 2023.
- 626 Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. Minif2f: a cross-system benchmark for
627 formal olympiad-level mathematics. *arXiv preprint arXiv:2109.00110*, 2021.
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

A FORMAL ENVIRONMENT AND PROOF TOOLS

Interactive Theorem Provers. Interactive Theorem Provers (ITPs), such as Isabelle (Paulson, 1994), are essential tools in modern mathematical verification. They help incorporate mathematical definitions and theorems into a consistent logical framework, such as Higher-Order Logic or Dependent Type Theory, which their kernels implement. The kernel is vital in the verification process, carefully checking each theorem to ensure it is correctly recognized by the ITP, thus maintaining system integrity. Using an ITP involves expressing the theorem in its programming language and then breaking it down into simpler goals or subgoals. A theorem is proven once it is reduced to known facts. We chose Isabelle for our paper because it has an intuitive interface, supports various logical frameworks, and offers a comprehensive library of formalized mathematics.

Sledgehammer. Sledgehammer (Paulsson & Blanchette, 2012) is a powerful tool for automating reasoning within the interactive theorem prover Isabelle. It works by translating the goals expressed in Isabelle/HOL’s higher-order logic into other types of logic, such as first-order logic. These translated goals are then sent to automated theorem provers like E, CVC4, Z3, Vampire, and SPASS. If any of these automated provers succeed in finding proof, Sledgehammer reconstructs the proof within the Isabelle/HOL framework using certified provers, such as metis, meson, and smt. This reconstructed proof, being more understandable to humans, greatly improves the system’s usability and enhances the efficiency and effectiveness of interactive theorem proving.

B COMPARATIVE ANALYSIS WITH LEAN-BASED METHODS

Table 4: Performance on the miniF2F dataset. Methods marked with [†] integrate human-written informal proofs, either in full or partially, during the proof search. Methods denoted by [‡] employ a tree search strategy, where the search budget is structured as $N \times M$, representing N search trees with M as the budget per tree.

Model	Base	Synthetic Size	Search Budget	miniF2F-valid	miniF2F-test
Lean-based					
HTPS (Lample et al., 2022) [†]	-	-	-	58.6%	41.0%
Lean-STaR (Lin et al., 2024) [‡]	InternLM2-Math-plus	50k Theorems	64×50	-	46.3%
DeepSeek-Prover (Xin et al., 2024a)	DeepSeekMath-Base	712k Theorems	65536	-	50.0%
InternLM2-StepProver (Wu et al., 2024) [‡]	InternLM2-Math-plus	1.155B Tokens	64×3200	63.9%	54.5%
DeepSeek-Prover-V1.5 (Xin et al., 2024b) [‡]	DeepSeekMath-Base	9B Tokens	32×6400	-	63.5%
Isabelle-based					
Sledgehammer	-	-	-	9.9%	10.4%
Sledgehammer+heuristic	-	-	-	18.0%	20.9%
Thor (Jiang et al., 2022b) [‡]	-	-	300	28.3%	29.9%
Thor + expert iteration (Wu et al., 2022) [‡]	-	-	-	37.3%	35.2%
DSP (Jiang et al., 2022a) [†]	Codex	-	100	42.6%	39.3%
Subgoal-Prover (Zhao et al., 2024)	GPT-3.5-Turbo	-	100	48.0%	45.5%
LEGO-Prover (Xin et al., 2023) [†]	GPT-3.5-Turbo	-	100	55.3%	50.0%
Lyra (Zheng et al., 2023) [†]	GPT-4	-	200	55.3%	51.2%
SubgoalXL (ours) [†]	Llama-3-8B	38k Theorems	16384	61.9%	56.1%

In this section, we provide a detailed analysis and discussion of the performance of our approach in relation to Lean-based theorem proving methods, including HTPS (Lample et al., 2022), Lean-STaR (Lin et al., 2024), DeepSeek-Prover (Xin et al., 2024a), InternLM2-StepProver (Wu et al., 2024) and DeepSeek-Prover-V1.5 (Xin et al., 2024b). Lean employs a unique set of tactics and automation techniques that differ significantly from those used in the Isabelle proof assistant, making direct comparisons within the main body of our evaluation less meaningful.

Table 4 demonstrates the significant performance variation across Lean-based and Isabelle-based approaches on the miniF2F dataset. Lean-based methods generally operate on a much larger scale, utilizing between $15 \times$ and $120 \times$ more synthesized proofs than Isabelle-based systems. For example, DeepSeek-Prover-V1.5 employs 9 billion tokens of synthesized data, while our approach uses only 38k theorems. Furthermore, these Lean methods typically operate with a search budget that is approximately $12.5 \times$ larger than ours. InternLM2-StepProver, for instance, operates with a search budget of 64×3200 , compared to our search budget of 16,384.

Despite this, SubgoalXL outperforms most Lean-based methods, achieving 56.1% on the miniF2F-test dataset compared to 54.5% for InternLM2-StepProver and 50.0% for DeepSeek-Prover. On the miniF2F-valid dataset, SubgoalXL also exhibits strong performance, further demonstrating the efficiency of our subgoal-based approach, even with significantly smaller synthetic data and a more moderate search budget. These results suggest that our approach is highly effective in optimizing formal theorem-proving tasks. A potential direction for future work is to scale up the size of synthesized proofs and the search budget, which could further enhance performance and competitiveness, particularly in comparison to resource-intensive Lean-based methods.

C MORE DETAILS ABOUT SUBGOAL-BASED PROOF

After creating 26 demonstration examples, as detailed in §3.2, we used the prompt shown in Figure 6 to annotate subgoal-based proofs for the problems in the informal dataset.

```

Prompt to generate subgoal-based proofs

System:
To generate a subgoal-based proof from an informal statement and informal proof,
systematically decompose the informal proof into smaller logical steps (subgoals), then
prove each subgoal to logically build towards proving the main theorem.

User:
text \<open>
### Informal Statement:
When Rachel divides her favorite number by 7, she gets a remainder of 5. What will the
remainder be if she multiplies her favorite number by 5 and then divides by 7? Show that it
is 4.

### Informal Proof:
Let  $n$  be Rachel's favorite number. Then  $n \equiv 5 \pmod{7}$ , so  $5n \equiv 5 \cdot 5 \equiv 25 \equiv 4 \pmod{7}$ .
\<close>

Assistant:
Step 1: Given  $n \equiv 5 \pmod{7}$ , substitute 5 for  $n$  in  $(5 \cdot n) \pmod{7}$ .
Step 2: Calculate  $25 \pmod{7} = 4$ .
Step 3: Thus,  $(5 \cdot n) \pmod{7}$  simplifies to 4.
Therefore, the remainder when 5 times Rachel's favorite number is divided by 7 is 4.

User:
text \<open>
### Informal Statement:
{informal_statement}

### Informal Proof:
{informal_proof}
\<close>

Assistant:

```

Figure 6: Prompt to generate subgoal-based proofs.

D MORE DETAILS ABOUT EXPERT LEARNING

D.1 IMPLEMENTATION DETAILS OF EACH COMPONENT

The prompt templates for the formal statement generator, subgoal generator, and posterior subgoal generator are shown in Figures 7-9, respectively. To improve the diversity of outputs generated by the formal proof generator, we employ two distinct prompt templates, as illustrated in Figures 10 and 11. In this study, all components are initialized with Llama-3-8B.

756
757
758
759
760
761
762
763
764
765

```

Prompt for formal statement generator

Translate the informal statement into a formal statement by defining variables and
assumptions explicitly, and then stating the main claim clearly using precise mathematical
notation.

### Informal Statement
{informal_statement}

### Formal Statement
{formal_statement}

```

766
767

Figure 7: Prompt for formal statement generator.

768
769
770

```

Prompt for subgoal generator

Generate a subgoal-based proof by identifying and breaking down the critical steps needed
to achieve the overall proof, explaining each subgoal with clear mathematical reasoning
and ensuring logical progression from one subgoal to the next until the final proof is
achieved.

### Informal Statement
{informal_statement}

### Formal Statement
{formal_statement}

### Subgoal-based Proof
{subgoal_based_proof}

```

771
772
773
774
775
776
777
778
779
780
781

782
783

Figure 8: Prompt for subgoal generator.

784
785
786

```

Prompt for posterior subgoal generator

Generate a subgoal-based proof by breaking down the formal proof into critical steps,
providing clear mathematical reasoning for each subgoal, and ensuring logical progression
from one subgoal to the next until the final proof is achieved.

### Informal Statement
{informal_statement}

### Formal Statement
{formal_statement}

### Formal Proof
{formal_proof}

### Subgoal-based Proof
{subgoal_based_proof}

```

787
788
789
790
791
792
793
794
795
796
797
798
799

800
801
802

Figure 9: Prompt for posterior subgoal generator.

803
804

D.2 ANNOTATION OF FORMAL AND INFORMAL DATA

805
806
807
808
809

For the problems within the informal dataset, we employed a mixture of deepSeek-math-base and Llama-3-8B using the prompt templates illustrated in Figures 12 and 13 to generate their corresponding formal statements and proofs. For the problems in the formal dataset, we use the prompt template shown in Figure 14 to annotate their informal statements and proofs.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

```

Prompt for formal proof generator (Template 1)
### Problem:
{informal_statement}

### Proof:
{formal_statement}

{formal_proof}

```

Figure 10: Prompt for formal proof generator (template 1).

```

Prompt for formal proof generator (Template 2)
(For problems from the formal dataset)
Develop formal proofs using Isabelle, leveraging auxiliary tools such as Sledgehammer to
enhance the proving process.

### Input
(* Informal Statement:
{informal_statement} *)
{formal_statement}

### Output
{formal_proof}

(For problems from the informal dataset)
Utilize Isabelle for the systematic verification of theorem proofs, employing Sledgehammer
as a supplementary tool. Approach the problem in a step-by-step manner.

### Problem
{informal_statement}

### Isabelle Proof
{formal_statement}

{formal_proof}

```

Figure 11: Prompt for formal proof generator (template 2).

D.3 OPTIMAL DISTRIBUTIONS FOR FORMAL STATEMENTS AND SUBGOAL-BASED PROOFS

The optimal distribution for the formal statement at the k -th iteration is given by:

$$p_{\text{fsg}}^*(\mathcal{S} \mid \mathcal{J}) = \frac{1}{Z(\mathcal{J})} p_{\text{fsg}}^{(k-1)}(\mathcal{S} \mid \mathcal{J}) \exp\left(\frac{1}{\beta} (\log p_{\text{sg}}(g \mid \mathcal{J}, \mathcal{S}))\right), \quad (3)$$

where $Z(\mathcal{J}) = \sum_{\mathcal{S}} p_{\text{fsg}}^{(k-1)}(\mathcal{S} \mid \mathcal{J}) \exp\left(\frac{1}{\beta} (\log p_{\text{sg}}(g \mid \mathcal{J}, \mathcal{S}))\right)$.

Similarly, the optimal distribution for the subgoal-based proof at the k -th iteration is determined by:

$$p_{\text{sg}}^*(g \mid \mathcal{J}, \mathcal{S}) = \frac{1}{Z(\mathcal{J}, \mathcal{S})} p_{\text{sg}}^{(k-1)}(g \mid \mathcal{J}, \mathcal{S}) \exp\left(\frac{1}{\beta} (\log p_{\text{fpg}}(\mathcal{P} \mid \mathcal{J}, \mathcal{S}, g))\right), \quad (4)$$

where $Z(\mathcal{J}, \mathcal{S}) = \sum_g p_{\text{sg}}^{(k-1)}(g \mid \mathcal{J}, \mathcal{S}) \exp\left(\frac{1}{\beta} (\log p_{\text{fpg}}(\mathcal{P} \mid \mathcal{J}, \mathcal{S}, g))\right)$.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

```
Prompt to generate formal statements  
System:  
Use Isabelle to formalize informal mathematical problems by defining the necessary  
variables and assumptions, followed by constructing the precise theorem statements.  
User:  
text \<open>  
When Rachel divides her favorite number by 7, she gets a remainder of 5. What will the  
remainder be if she multiplies her favorite number by 5 and then divides by 7? Show that it  
is 4.  
\<close>  
Assistant:  
theorem mathd_numbertheory_335:  
  fixes n :: nat  
  assumes h0 : "n mod 7 = 5"  
  shows "(5 * n) mod 7 = 4"  
User:  
text \<open>  
{informal_statement}  
\<close>  
Assistant:
```

Figure 12: Prompt to generate formal statements.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

```

Prompt to generate formal proofs

System:
Use Isabelle to systematically prove theorem statements. Use tools like sledgehammer to
assist in proving.

User:
text \<open>
When Rachel divides her favorite number by 7, she gets a remainder of 5. What will the
remainder be if she multiplies her favorite number by 5 and then divides by 7? Show that it
is 4.
\<close>

theorem mathd_numbertheory_335:
  fixes n :: nat
  assumes h0 : "n mod 7 = 5"
  shows "(5 * n) mod 7 = 4"

Assistant:
proof -
  (* Step 1: Given  $n \equiv 5 \pmod{7}$ , substitute 5 for n in  $(5 * n) \pmod{7}$ . *)
  have c0: "(5 * n) mod 7 = (5 * 5) mod 7"
    using h0
    by simp

  (* Step 2: Calculate  $25 \pmod{7} = 4$ . *)
  then have "\<dots> = 4"
    by simp

  (* Step 3: Thus,  $(5 * n) \pmod{7}$  simplifies to 4. *)
  then have "(5 * n) mod 7 = 4"
    using c0
    by simp

  (* Therefore, the remainder when 5 times Rachel's favorite number is divided by 7 is 4. *)
  then show ?thesis
    by simp
qed

User:
text \<open>
{informal_statement}
\<close>

{formal_statement}

Assistant:

```

Figure 13: Prompt to generate formal proofs.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Prompt to generate informal statements and proofs

System:

Clearly state the problem, simplify and explain each step in plain language, summarize the result, and use LaTeX for mathematical expressions.

User:

theorem mathd_numbertheory_335:

fixes $n :: \text{nat}$

assumes $h0 : "n \bmod 7 = 5"$

shows $"(5 * n) \bmod 7 = 4"$

proof -

(Step 1: Given $n \equiv 5 \pmod{7}$, substitute 5 for n in $(5 * n) \bmod 7$. *)*

have $c0 : "(5 * n) \bmod 7 = (5 * 5) \bmod 7"$

using $h0$

by simp

(Step 2: Calculate $25 \bmod 7 = 4$. *)*

then have $"\langle \dots \rangle = 4"$

by simp

(Step 3: Thus, $(5 * n) \bmod 7$ simplifies to 4. *)*

then have $"(5 * n) \bmod 7 = 4"$

using $c0$

by simp

(Therefore, the remainder when 5 times Rachel's favorite number is divided by 7 is 4. *)*

then show ?thesis

by simp

qed

Assistant:

Problem

When Rachel divides her favorite number by 7, she gets a remainder of 5. What will the remainder be if she multiplies her favorite number by 5 and then divides by 7? Show that it is 4.

Solution

Let n be Rachel's favorite number. Then $n \equiv 5 \pmod{7}$, so $5n \equiv 5 \cdot 5 \equiv 25 \equiv 4 \pmod{7}$.

User:

{formal_statement}

{formal_proof}

Assistant:

Figure 14: Prompt to generate informal statements and proofs.