LeanTree: Accelerating White-Box Proof Search with Factorized States in Lean 4

Matěj Kripner¹ Michal Šustr² Milan Straka¹

Abstract

Automated theorem proving (ATP) has been a classical problem in artificial intelligence since its inception, yet it remains challenging due to its vast state and action space. Large language models (LLMs) have recently emerged as a promising heuristic for ATP, but they lack correctness guarantees and thus require interaction with a proof verifier. Such interactions typically follow one of two approaches: black-box interaction, which does not utilize intermediate proof states, or white-box approaches, which allow for incremental proof construction and examination of intermediate states. While black-box approaches have directly benefited from recent LLM advances, white-box methods have comparatively lagged behind. In this paper, we address this gap by introducing LeanTree, which consists of (i) a tool built in the Lean 4 language that factorizes complex proof states into simpler, independent branches, and (ii) a dataset of these factorized intermediate states. Our white-box tooling offers several advantages over black-box approaches: it simplifies evaluation, reduces necessary context, generates richer training data, enables parallel search across multiple states, supports efficient reuse of states, and provides feedback in case of errors. Our preliminary results hint that white-box approaches outperform black-box alternatives in some settings.

1. Introduction

Automated theorem proving (ATP) has long been a foundational task in artificial intelligence (Gelernter, 1959), with applications ranging from the formalization of mathematics (van Doorn et al., 2023), physics (Tooby-Smith, 2024) and chemistry (Bobbin et al., 2024) to software verification (Avigad et al., 2025), cryptography (Doussot, 2024), SQL query verification (Chu et al., 2018), and general reasoning (Jiang et al., 2024).

Despite extensive research using logical formalisms, ATP continues to face significant challenges due to the combinatorial explosion of its underlying state and action spaces (Harrison, 2009; Baader, 2003; Barrett & Tinelli, 2018) and the infinite branching factor introduced by exogenous term creation (Trinh et al., 2024).

Recent advances in large language models (LLMs) capable of sophisticated reasoning in human languages (Guo et al., 2025) have sparked interest in their application to ATP (Ren et al., 2025), serving as a powerful heuristic in navigating the proof search. Conversely, ATP and formal verification present an avenue for improving reliability and interpretability of LLMs (Achim & Tenev, 2023). Open-source tooling and datasets are crucial in leveraging this synergy.

The existing approaches for ATP can be categorized into two distinct methodologies. In *white-box* generation, the prover interacts with a formal verifier iteratively, generating individual proof steps conditioned on the internal proof state, interleaving proof search with verification. In contrast, *black-box* provers only use the formal verifier to verify the final proof, generating proof steps with no feedback about the validity of preceding steps or the current internal proof state. Thus, in a black-box search, it falls upon the proof step generator to deduce the evolving proof state, as this is a necessary precondition to determining the validity of a subsequent proof step. This coupling of policy model and world model increases complexity of proof step generation, necessitating the use of more powerful models, while also requiring the proof generator to be retrained when modifying the formal verifier in case of upgrades or patches.

While black-box whole-proof generation methods directly benefit from recent advances in LLMs (Wang et al., 2025), white-box tree-search methods have seen comparatively little improvement since their introduction. Notably, the hypertree proof search approach introduced by Lample et al. (2022) has not been replicated in open source or built upon, except for follow-up work by the same research

¹Charles University, Faculty of Mathematics and Physics ²Czech Technical University, Faculty of Electrical Engineering. Correspondence to: Matěj Kripner <kripner@ufal.mff.cuni.cz>.

The second AI for MATH Workshop at the 42nd International Conference on Machine Learning, Vancouver, Canada.

group (Gloeckle et al., 2024). We attribute this primarily to insufficient tooling and the lack of suitable datasets. Specifically, machine learning-based white-box techniques require nontrivial preprocessing of existing proofs to extract intermediate proof states. Additionally, the size of proof states in verification languages like Lean tends to increase as the proof search progresses, exacerbating the distribution shift for LLMs pretrained on general text corpora. As such, it is desirable to simplify proof states by factorizing them into branches that can be solved independently. Although Lample et al. (2022) demonstrated the benefit of this white-box approach using factorized states, their tooling and dataset remain practically unusable for the ATP community (see Section 5).

To address these challenges, we introduce *LeanTree*, which is a) a tool for white-box interaction and data extraction in Lean 4 (Moura & Ullrich, 2021), and b) a dataset of preprocessed proofs from Mathlib (Community, 2020) and DeepSeekProver-V1 (Xin et al., 2024a) in a unified format. Our contributions are as follows:

- We build on top of Lean REPL (Community, 2025) to enable programmatic proof search over factorized proof states, meaning that independent goals can be solved individually (Section 3.2). To enable this, we detect dependence between goals caused by metavariable coupling (Section 3.3). Other improvements over Lean REPL include better error reporting, richer information about open goals, and a new incremental strategy for proof verification in Lean REPL using the Lean's verification kernel (Appendix B).
- We build a data extraction module by integrating Lean REPL, PaperProof (Karunus & Kovsharov, 2024), and a custom algorithm for proof tree building and tactic simplification (Section 4.1). The extracted dataset also contains information about the current Lean context for each theorem and the distance to the end of the proof for each step.
- We release the resulting suite of tools and datasets, dubbed *LeanTree*, in a format directly accessible by the community¹² with a user-friendly Python interface.
- We present a preliminary experiment hinting that supplying information about the intermediate proof state to a pretrained model during proof search outperforms black-box generation.
- We identify and prevent false-positive errors in proof search that result from the usage of library search tactics such as apply? (Appendix A). For example, such errors are present in the proofs presented by DeepSeek-Prover-V2 (Ren et al., 2025).

2. Related work

We give a brief overview of Lean, black-box and white-box ATP approaches, and existing datasets and benchmarks.

Lean 4 (Moura & Ullrich, 2021) is a leading formal language for theorem proving with applications both in the formalization of mathematics (van Doorn et al., 2023) and in numerous other domains (Tooby-Smith, 2024; Bobbin et al., 2024; Doussot, 2024; Jiang et al., 2024). One of the pleasant attributes of Lean is its extensibility, stemming from the fact that its theorem proving component is written in Lean itself, and the user can modify both its syntax and its semantics. Such an approach is possible because Lean 4 is also a fully-fledged functional programming language. This is in contrast to previous versions of Lean, where some of the existing Lean interaction tools had to bind directly to the Lean kernel written in C++.

Black-box ATP. Automated theorem proving can be formulated as a text-to-text problem so that any existing approach based on a large language model (LLM) is directly applicable. This approach to theorem proving was pioneered by Baldur (First et al., 2023). Recently, a LLM with chain-of-thought refined using reinforcement learning was utilized by DeepSeek-Prover-V2 (Ren et al., 2025), Kimina-Prover (Wang et al., 2025), and InternLM-Math (Ying et al., 2024b), securing top places on relevant benchmarks.

White-box ATP. In white-box ATP, the prover interacts with the verification system iteratively during a proof search, generating proof steps conditioned on the internal proof state, with consequences discussed in the introduction of this paper. In neural theorem proving, this approach was pioneered by Holophrasm (Whalen, 2016) in Metamath (Megill & Wheeler, 2019), using a UCT-based (Kocsis & Szepesvári, 2006) AND-OR tree search with a RNN sequence-to-sequence model for tactic enumeration and payoff prediction. This approach was extended by GPT-f (Polu & Sutskever, 2020) and Evariste (Lample et al., 2022), with the latter additionally targeting Lean 3. ABEL (Gloeckle et al., 2024) extended Evariste with incremental improvements and Lean 4 compatibility, but did not release their code or dataset.

Datasets. The standard library of Lean, Mathlib (Community, 2020), offers more than 107k formal definitions and 220k formal proofs. Together with numerous formalization projects created by the community, this offers a sizable human-written dataset for black-box theorem proving.

In contrast, white-box ATP approaches require data in a more preprocessed format that reveals intermediate proof states. First, such a dataset was provided by LeanStep (Han et al., 2021) in Lean 3 by extracting tactic invocation data from Mathlib. LeanDojo (Yang et al., 2023) extends this for Lean 4, additionally containing information about the

¹https://github.com/Kripner/leantree
²https://huggingface.co/datasets/ufal/
leantree

premises used in each tactic application, suitable for training retrieval-augmented approaches. Similar datasets are also offered by Pantograph (Aniva et al., 2024) and lean-trainingdata (Morrison, 2023). Additionally, Lean-Workbook (Ying et al., 2024a) and DeepSeek-Prover-V1 (Xin et al., 2024a) provide large collections of automatically formalized proofs of various quality. However, none of these datasets offers tactic proofs in a structured, simplified way, which we describe in this paper.

Benchmarks. ATP benchmarks typically consist of a set of formalized theorem statements, optionally accompanied by their informal counterparts in natural language, where the objective is to generate a valid proof. MiniF2F (Zheng et al., 2021) offers 488 problems drawn from mathematical olympiad and undergraduate mathematics courses, formalized in parallel in Lean 3, Metamath, Isabelle, and HOL Light. The benchmark was later ported to Lean 4 by Yang et al. (2023), with some formalization errors corrected by Wu et al. (2024).

Similarly, ProofNet (Azerbayev et al., 2023a) provides 371 problems formalized in Lean 3, together with their natural language statement and natural language proof, making it suitable also for evaluating autoformalization approaches. This benchmark was ported to Lean 4 by Xin et al. (2024b).

PutnamBench (Tsoukalas et al., 2024) comprises of 1 099 problems formalized in Lean 4, Isabelle, and Coq. This benchmark is comparatively difficult, with DeepSeek-Prover-V2 (Ren et al., 2025), the current state-of-the-art, attaining only 7.2 % success rate with pass@1024.³

Lean Tooling. The comparison of LeanTree with existing tools is given in Section 5.

3. Lean programmatic interaction

We build on top of Lean REPL (Community, 2025) to enable incremental proof execution. In this section, we describe the LeanTree interaction module.

Preliminaries. In Lean, theorem proving consists of producing a term whose type is equal to the theorem type. For example, consider the following theorem:

theorem sub_zero (a b :
$$\mathbb{N}$$
) (h : b = 0)
: a - b = a (1)

Proving Theorem 1 consists of supplying a term of type a - b = a using available free variables a, b of type \mathbb{N} and h of type b = 0. While such a term can be written down explicitly, it is common to utilize Lean's *tactic mode* to structure the proof in a format more aligned with human-

³https://trishullab.github.io/

written proofs.

Upon entering tactic mode using the by keyword, Lean creates a *metavariable* with the desired type a - b = a and no value assigned yet. Metavariables can be thought of as holes in a proof that have to be filled in before the proof is considered finished. This gives us an intermediate state in the proof construction, where the proof term is type-correct but still contains holes to be filled in later. The still unassigned metavariables are called *open goals* so that each goal has a target type and a set of available free variables called *hypotheses*. Because goals are metavariables, the two terms are often used interchangeably.

We then proceed by applying a *tactic* to the list of open goals. A tactic is a procedure that assigns values to one or more open goals, with each of the values possibly containing new unassigned metavariables. In this way, a tactic invocation transforms the list of open goals into a new list of open goals. Typically, a tactic only affects the *main goal* (defined to be the first goal), either solving it or reducing it to one or more simpler sub-goals. The proof is finished once there are no more open goals, i.e., every metavariable has been assigned.

3.1. LeanTree interaction module

LeanTree offers a Python interface, where the user can initiate a proof search either from an individual theorem statement or from a Lean file. After applying a Lean tactic to a proof state, LeanTree returns a list of resulting sub-states. This interface is suitable for tree search and backtracking as the interaction can be started from any intermediate proof state. Additionally, given the high CPU requirements of Lean 4, LeanTree implements a dynamic pool of environments suitable for parallel execution.

3.2. Factorized proof states

We follow Lample et al. (2022) in factorizing proof states into independent goals, which can be explored and proven individually. Note that in this section, we disregard metavariable coupling, which is then addressed in Section 3.3.

The utility of proof state factorization has been identified by Lample et al. (2022), citing the "re-use and parallelization in the proof search algorithm" as their primary motivation. An example of potential state reuse is shown in Figure 1, where the same node is reached using two different search paths. To illustrate proof search parallelization, consider the following theorem:

theorem mul_eq_zero_iff (n m :
$$\mathbb{N}$$
)
: n * m = 0 \leftrightarrow n = 0 \vee m = 0 (2)

A possible strategy to prove Theorem 2 is to branch out based on the values of n and m. We achieve this by execut-

PutnamBench/leaderboard.html (retrieved Jun 2025)

ing the tactic cases n <;> cases m. Since each of the two numbers can be either zero or a successor of another number following the Peano axioms, the resulting proof state consists of four goals. For example, one of the states is $(n'+1) \cdot 0 = 0 \Leftrightarrow n'+1 = 0 \lor 0 = 0$, corresponding to the case where n is a successor and m is zero.

To complete the proof, all four goals must be closed. In a non-factorized setting, the proof state consisting of four goals is assessed by the policy and critic models as a whole. Since tactics typically affect only the main goal, the last goal can only be worked on once all others have been closed. In contrast, factorized states allow the search algorithm to focus on the goals individually, in any order the algorithm deems fit.

As a final advantage, simplifying states reduces the difficulty of the policy task, deciding which tactic to apply next, and of the critic task, estimating the value of the current goal. Note that proof states can grow arbitrarily complex during proof search since each state corresponds to a list of goals.

AND-OR Search. Search over factorized states naturally leads to an AND-OR tree, where a prover selects a tactic in an OR node and subsequently has to prove all goals in an AND node. Figure 1 shows an example where choosing either the rw tactic or the cases tactic in the root leads to a proof. For a more in-depth introduction to proof search in Lean, we recommend (Limperg & From, 2023).

3.3. Metavariable coupling

Contrary to the simplified view presented in Section 3.2, it is not always possible to fully factorize the proof state into individual goals that can be proven independently. Subsequently, nodes in the proof tree are allowed to contain more than one goal. Specifically, this occurs when two goals share the same metavariable. For example, applying the transitivity lemma Nat.le_trans to the goal $2 \le 5$ produces goals $2 \le ?m$ and $?m \le 5$ sharing a metavariable ?m whose value has not yet been decided.⁴ Subsequently, choosing a value for ?m in one of the goals affects the provability of the other goal, and therefore the two goals cannot be explored independently. For a more in-depth explanation of this phenomenon, we recommend Aniva et al. (2024).

LeanTree detects all dependencies between goals caused by metavariables and offers strategies to deal with them during proof search. By default, goals are factorized to the maximum extend allowed by metavariable coupling.

3.4. Proof verification

Originally, the Lean REPL did not utilize Lean verification kernel to type-check assignments introduced by tactic exe-

cutions, leading to incorrect proofs being accepted in some cases. Our prover managed to exploit these inconsistencies, finding 12 incorrect proofs on MiniF2F that passed verification in the REPL. We describe these incorrect proofs and recent approaches to fix Lean REPL in the appendix. Note that all proofs found by our prover are also independently verified directly using Lean, which guarantees their correctness.

4. Supervised data extraction

To enable machine learning-based approaches to imitate proofs without interfacing with Lean, LeanTree offers a data extraction module to create supervised data from existing Lean proofs. The data extraction module is fully compatible with the interaction module in the sense that the extracted proof trees can be directly executed and verified.

4.1. Proof tree building

To achieve compatibility with the interaction module, the proof states extracted from supervised data must be factorized into independent branches as described in Section 3.2. This process yields a proof tree where each node contains a list of goals (typically a singleton list), and each edge corresponds to a tactic application.

To illustrate some challenges in building a proof tree, consider the Lean proof and its corresponding proof tree shown in Figure 2, where we prove n < 2n under the condition n > 0. For the purposes of this illustration, we do not utilize powerful tactics like linarith, which would solve this trivial theorem nearly automatically. Instead, we illustrate branching, metavariable coupling, and the introduction of a new hypothesis.

Note that in the original Lean proof in Figure 2, tactics are nontrivially structured via branching using cases, focusing a goal via the center-dot operator, and merging multiple rw commands into one. Additionally, each Lean tactic consumes the set of currently open goals – the unfactorized proof state – and produces a new set of open goals – a new proof state. In contrast, each edge of the constructed proof tree corresponds to a simple tactic applied to a minimal list of goals that cannot be split further due to metavariable coupling.

To illustrate metavariable coupling, observe that applying Nat.le_trans yields three goals that are not split due to a shared metavariable ?m. However, once the metavariable is assigned using exact, the remaining two goals are now independent and can be split.

Finally, note that in the proof tree, the tactic case succ n' and its parent node are synthetic in the sense that they do not have a counterpart in the original proof. This synthetic

⁴A third goal of type \mathbb{N} is created for technical reasons.



Figure 1: Example of an AND-OR proof tree for Theorem 1. Nodes connected with a dashed line could potentially be merged and correspond to transpositions in board games. This is possible thanks to factorization (see Section 3.2).

tactic is necessary to assign the name n' to the otherwise inaccessible variable – denoted n+ before the assignment – so that it can be referenced by later tactics.

Other Lean features not shown in Figure 2 that further complicate proof tree construction include tactic combinators such as try, iterate, and any_goals, where the set of affected goals is determined in runtime, constructs such as all_goals and <; > that go against the tree formulation by operating on multiple goals, or tactics like switch and rotate_left that change the order of open goals.

We implement the transformation illustrated in Figure 2 in three stages described in the following three subsections.

4.1.1. SINGLETON TREES

In the first stage of proof tree building, a Lean proof is transformed into a *singleton tree*, where each node contains a single goal, by matching tactic applications with the goals they affect. This is nontrivial, because, in general, a tactic application transforms a list of open goals into an entirely new list of empty goals without any guarantees about the transformation.

To build the singleton tree, we integrate Lean REPL with the BetterParser module of PaperProof (Karunus & Kovsharov, 2024), which solves the analogous problem for the purposes of visualization.

4.1.2. TACTIC SIMPLIFICATION

In the second stage, the singleton tree is modified so that complex tactics are broken down into simpler ones and each tactic is a valid self-contained proof step. Here we describe some of these modifications.

Nested Tactics. Tactics parametrized by a term can themselves contain nested tactic blocks. For example, consider the following tactic from the Complex Analysis module in Mathlib.

```
exact (by
rw [aeval_algHom_apply, hw, map_zero], (3)
rfl)
```

In Listing 3, the exact tactic is parametrized by a term which contains a tactic block introduced using the by keyword. We break down this complex proof step by masking the nested tactic block and instead spawning a child goal, yielding the following series of simpler proof steps.

Note that the rfl tactic is not masked since it is used as a term instead of as a tactic.

Structuring tactics. The cases and induction tactics introduce syntactic branching by partitioning a proof based on the value of an inductive type. The respective tactic optionally contains a solving branch for each constructor of such type. For example, the cases tactic in Listing 2 provides separate proofs for the zero and succ constructors of variable n. This cases tactic spans 10 lines, which is contrary to our effort of breaking down the proof into simple steps. For this reason, we move the individual proof branches from the cases tactic into their separate proof steps.

However, branches corresponding to parametrized constructors can optionally assign names to some of the constructors' parameters. For example, in Listing 2, the succ constructor is parametrized by a single natural number which is given the name n'. When the cases tactic is broken down, the name n' has to be assigned manually, since Lean disallows referring to variables with no explicitly assigned name. We achieve this by adding a synthetic tactic case which has no direct tactic counterpart in the original proof. Note that the same effect could be achieved using the rename tactic, which is, however, rare in human-written proofs.

Merged rw tactics. In Lean, subsequent rw tactics can be merged into a single tactic, which is commonplace in human-writen proofs. We partition such composite tactics back into individual rw tactics. This is also the case for the rwa tactic, where we additionally add the assumption tactic to preserve semantics.



Figure 2: Proof tree builder in LeanTrees transforms a Lean proof with complex structure (left) into a proof tree (right) with simple states (lists of goals that cannot be split further) and transitions (tactics). Synthetic elements in the proof tree are marked using dashed border.

4.1.3. MERGING METAVARIABLE-DEPENDENT GOALS

To obtain the final proof tree, sibling nodes in a singleton tree are merged if they share a metavariable. This corresponds to the fact that metavariable-dependent goals cannot be explored independently.

During this merging process, the entire proof tree is executed in Lean REPL starting from the root state and following each tree edge, verifying that each proof branch ends in a proven state. In this way, we ensure the correctness of the whole proof tree building process. Proof trees that fail this verification are not included in the final dataset (cf. Section 4.2).

4.2. LeanTree dataset

We release the extracted proof tree data in an unified format from two sources: 1) a recent version of Mathlib 4,⁵ the standard library of human-written proofs in Lean, and 2) a collection of 27.5K proofs autoformalized by DeepSeek-Prover-V1.⁶ Importantly, each sample in the LeanTree dataset corresponds to a Lean file rather than just an individual theorem. This is necessary to capture the structure of a real-world Lean project like Mathlib where a proof can depend on any definition located above it in the source file.

⁵leanprover/lean4:v4.19.0, under Apache-2.0 license

⁶https://huggingface.co/datasets/

deepseek-ai/DeepSeek-Prover-V1, MIT License

Each file in the LeanTree dataset contains a list of theorems, and each theorem contains a list of all tactic proofs in its proof term. Note that there can be more than one tactic proof for a theorem if its proof contains more than one non-nested by-blocks. For each tactic proof, LeanTree then contains a proof tree with nodes corresponding to factorized proof states and edges corresponding to tactic applications.

To demonstrate a possible use case for proof trees, the dataset also contains the size and depth for each proof tree node. These can serve as objectives for a critic model in various proof search algorithm.

Additionally, the LeanTree dataset contains information about the surrounding context, namely the list of imported modules for each Lean file and the list of open namespaces for each theorem. The correspondence between samples in the dataset and the underlying Lean repository is given by character offsets specifying the span of each theorem, proof, and tactic execution.

Overall, LeanTree contains 74 706 factorized tactic proofs from Mathlib and 26 201 from DeepSeek-Prover-V1. Since Lean was not designed to enable factorized proof tree search out-of-the-box, there are a large number of small technical challenges to overcome during the proof tree building. While we are continually working on perfecting this process, not all tactic proofs can currently be converted. Specifically, 23.0% of tactic proofs in Mathlib and 4.7% in DeepSeek-Prover-V1 were not converted. We note that for Mathlib, 28.6% of the issues stem from the usage of calc and conv tactics that fundamentally change the structure of a proof.

5. Comparison to existing tools

Several tools have recently emerged both for programmatic interaction with Lean and machine learning data extraction from human-written Lean proofs. In this section, we compare them to LeanTree. We argue that LeanTree is unique in its data extraction capabilities.

Evariste (Lample et al., 2022) corresponds most closely to LeanTree. Evariste enables factorizing proof state into individual goals, which can be explored independently. In addition, Evariste detects metavariable coupling and does not split metavariable-dependent goals, similar to LeanTree.

However, the Evariste dataset has not been released and the tool itself is available only in a limited, nonexecutable format.⁷ This makes it challenging to build on top of Evariste's contributions.

Importantly, Evariste was only implemented for Lean 3,

which is now deprecated. Since its core logic is implemented using C++, directly binding to Lean 3 internals, its migration to Lean 4 would be difficult. In contrast, the core logic of LeanTree is implemented in Lean 4 itself, which is a common Lean 4 paradigm (Moura & Ullrich, 2021). Furthermore, LeanTree builds on top of the semi-official Lean REPL project, and it is therefore set to benefit from future improvements by the community.

LeanDojo (Yang et al., 2023) offers both an interaction module and a supervised data extraction module. The extracted data also contains information about the premises used in each tactic application, enabling training of retrieval-augmented generation models. The extracted datasets are available for both Lean 3 and Lean 4.

However, LeanDojo does not offer proof state factorization, tactic simplification, or proof tree building. First, this means that a tactic can only be applied to a whole proof state consisting of a list of goals, yielding a new set of goals. This prevents parallel search over independent goals.

Second, tactics in the LeanDojo dataset are left in the form in which they appear in the abstract syntax tree of a Lean source code. This includes complex or nested tactics that span multiple lines. For example, the cases tactic shown in Figure 2 is left unaltered. The user of the dataset can then choose to filter out such samples, losing valuable training data, or keep them in the training set, but increase the difficulty of the modeling task and reduce the granularity of the inference-time proof search.

Third, LeanDojo does not contain information about the subproof size and sub-proof depth for each node, and provides tactics in an unstructured format that is not conducive to proof tree building.

Pantograph (Aniva et al., 2024) In contrast to LeanDojo, Pantograph's interaction and data extraction modules offer proof state factorization, meaning that independent goals can be explored independently. Pantograph also provides information about which goals are coupled via a shared metavariable.

However, tactic simplification and proof tree building are not supported by the data extraction module, leading to the same drawbacks as described for the case of LeanDojo. Information about sub-proof size and sub-proof depth is also not provided.

In addition, the Pantograph API is not well suited for gradually building the current Lean context by interleaving definitions and theorems that can depend on any previously declared ones. Instead, theorems are treated as independent units.

⁷Specifically, the codebase states: "code does NOT run out of the box as we removed references to Meta-internal systems" https://github.com/facebookresearch/ Evariste



Figure 3: Performance on MiniF2F using linear rollouts with the Llemma-7B model. Error bars show standard deviation over 5 runs.

6. Experiments

We investigate the utility of white-box approach in proof search using the Llemma-7B model (Azerbayev et al., 2023b) (MIT license), reporting the success rate on MiniF2F test set (Zheng et al., 2021) modified for Lean 4 by Wu et al. (2024). For each input theorem, we run N = 10 independent linear rollouts, starting from an empty proof and sequentially sampling tactics proposed by the model for a maximum of M = 25 steps. The search is successful if the concatenated tactics from a rollout are a valid proof of the input theorem. If no proof is found during the N rollouts, the search is unsuccessful.

In Figure 3, we report the results for the black-box and white-box rollouts averaged over 5 runs with error bars of one standard deviation. We additionally report the success rate of black-box whole proof generation, where the model is prompted to generate the entire proof at once. All our measurements are averaged over 5 runs. Lastly, we include the success rate of proof generation using best-first search (Polu & Sutskever, 2020) reported by Azerbayev et al. (2023b), which performs white-box rollouts in a tree fashion and is therefore a natural upper bound on the white-box rollout performance.

Our results show that in our specific setting, supplying the prover with internal proof states improves proof search performance. Interestingly, black-box linear rollouts achieve worse performance than black-box whole-proof generation, indicating that constraining the model to generate only one tactic at a time hurts its performance. However, rollouts with access to the internal state consistently outperform both black-box methods.

We run all our experiments on a single node equipped with 8 AMD MI210 accelerators and 192 CPU nodes, parallelizing both model inference and Lean interaction. In total, our experiments necessitated 59 GPU-hours and 450 CPU-hours.

7. Conclusion

In this work, we introduced *LeanTree*, a novel white-box tool designed to accelerate white-box approaches to automated theorem proving in Lean 4 via factorized proof states and structured proof trees. Our approach addresses limitations of existing black-box methods by leveraging intermediate proof states, enabling parallelized and efficient proof search, and providing structured feedback during theorem proving. LeanTree's integration with Lean REPL ensures that it will continue to benefit from patches and fixes by the Lean community.

Additionally, we released a unified dataset of factorized intermediate proof states derived from Mathlib and DeepSeek-Prover-V1. A limitation of LeanTree is that currently not all tactic proofs can be converted; however, we are working on expanding the coverage.

Finally, our preliminary experiments hint that white-box approaches, facilitated by LeanTree, can outperform blackbox methods in automated theorem proving. We hope that LeanTree serves as a stepping stone toward more efficient, parallelizable, and reliable proof search and encourages further research into white-box automated theorem proving.

Acknowledgements

Matěj Kripner was supported by the grant no. 25-18031S of the Czech Science Foundation (GAČR). This research was partially supported by SVV project number 260 821. We acknowledge VSB – Technical University of Ostrava, IT4Innovations National Supercomputing Center, Czech Republic, for awarding this project access to the LUMI supercomputer, owned by the EuroHPC Joint Undertaking, hosted by CSC (Finland) and the LUMI consortium through the Ministry of Education, Youth and Sports of the Czech Republic through the e-INFRA CZ (grant ID: 90254). The authors thank Martin Suda, Matej Straka, Tomáš Čížek, Radovan Haluška, Martin Schmid, and the Lean community for helpful suggestions and corrections.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Achim, T. and Tenev, V. Harmonic: Building mathematical superintelligence. https://harmonic.fun/, 2023. Accessed: 2025-05-14.
- Aniva, L., Sun, C., Miranda, B., Barrett, C., and Koyejo, S. Pantograph: A machine-to-machine interaction interface for advanced theorem proving, high level reasoning, and data extraction in lean 4. arXiv preprint arXiv:2410.16429, 2024.
- Avigad, J., Goldberg, L., Levit, D., Seginer, Y., and Titelman, A. A proof-producing compiler for blockchain applications. J. Autom. Reason., 69(2), April 2025. ISSN 0168-7433. doi: 10.1007/s10817-025-09723-y. URL https: //doi.org/10.1007/s10817-025-09723-y.
- Azerbayev, Z., Piotrowski, B., Schoelkopf, H., Ayers, E. W., Radev, D., and Avigad, J. ProofNet: Autoformalizing and formally proving undergraduate-level mathematics (2023). URL https://arxiv. org/abs/2302.12433, 2023a.
- Azerbayev, Z., Schoelkopf, H., Paster, K., Santos, M. D., McAleer, S., Jiang, A. Q., Deng, J., Biderman, S., and Welleck, S. Llemma: An open language model for mathematics. arXiv preprint arXiv:2310.10631, 2023b.
- Baader, F. The description logic handbook: Theory, implementation and applications. Cambridge university press, 2003.
- Barrett, C. and Tinelli, C. Satisfiability modulo theories. *Handbook of model checking*, pp. 305–343, 2018.
- Bobbin, M. P., Sharlin, S., Feyzishendi, P., Dang, A. H., Wraback, C. M., and Josephson, T. R. Formalizing chemical physics using the lean theorem prover. *Digital Discovery*, 3(2):264–280, 2024.
- Chu, S., Murphy, B., Roesch, J., Cheung, A., and Suciu, D. Axiomatic foundations and algorithms for deciding semantic equivalences of sql queries. *Proc. VLDB Endow.*, 11(11):1482–1495, July 2018. ISSN 2150-8097. doi: 10. 14778/3236187.3236200. URL https://doi.org/ 10.14778/3236187.3236200.

- Community, L. Lean 4 repl. https://github.com/ leanprover-community/repl, 2025. Accessed: 2025-04-28.
- Community, M. The lean mathematical library. In Proceedings of the 9th ACM SIGPLAN International Conference on Certified Programs and Proofs, CPP 2020, pp. 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.
- Doussot, G. Cryptography experiments in lean 4: SHA-3 implementation. Cryptology ePrint Archive, Paper 2024/1880, 2024. URL https://eprint.iacr. org/2024/1880.
- First, E., Rabe, M. N., Ringer, T., and Brun, Y. Baldur: Whole-proof generation and repair with large language models. In Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, pp. 1229–1241, 2023.
- Gelernter, H. L. Realization of a geometry theorem-proving machine. In *Proceedings of the First International Conference on Information Processing (IFIP)*, pp. 273–281, 1959.
- Gloeckle, F., Limperg, J., Synnaeve, G., and Hayat, A. ABEL: Sample efficient online reinforcement learning for neural theorem proving. In *The 4th Workshop on Mathematical Reasoning and AI at NeurIPS'24*, 2024.
- Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. arXiv preprint arXiv:2501.12948, 2025.
- Han, J. M., Rute, J., Wu, Y., Ayers, E. W., and Polu, S. Proof artifact co-training for theorem proving with language models. arXiv preprint arXiv:2102.06203, 2021.
- Harrison, J. Handbook of practical logic and automated reasoning. Cambridge University Press, 2009.
- Jiang, D., Fonseca, M., and Cohen, S. B. Leanreasoner: Boosting complex logical reasoning with lean. North American Chapter of the Association for Computational Linguistics, 2024. doi: 10.48550/arxiv.2403.13312.
- Karunus, E. and Kovsharov, A. Paperproof: A new proof interface for lean 4. https://github.com/ Paper-Proof/paperproof, 2024. URL https: //github.com/Paper-Proof/paperproof. Accessed: 2025-04-30.

- Kocsis, L. and Szepesvári, C. Bandit based monte-carlo planning. In *European conference on machine learning*, pp. 282–293. Springer, 2006.
- Lample, G., Lacroix, T., Lachaux, M.-A., Rodriguez, A., Hayat, A., Lavril, T., Ebner, G., and Martinet, X. Hypertree proof search for neural theorem proving. *Advances in neural information processing systems*, 35:26337–26349, 2022.
- Limperg, J. and From, A. H. Aesop: White-box bestfirst proof search for lean. In *Proceedings of the 12th ACM SIGPLAN International Conference on Certified Programs and Proofs*, CPP 2023, pp. 253–266, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400700262. doi: 10.1145/ 3573105.3575671. URL https://doi.org/10. 1145/3573105.3575671.
- Megill, N. and Wheeler, D. A. *Metamath: a computer language for mathematical proofs.* Lulu. com, 2019.
- Morrison, K. lean-training-data: Tools for extracting training-data from lean libraries. https://github.com/kim-em/lean-training-data, 2023. Accessed: 2025-07-01.
- Moura, L. d. and Ullrich, S. 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, pp. 625–635. Springer, 2021.
- Polu, S. and Sutskever, I. Generative language modeling for automated theorem proving. *arXiv preprint arXiv:2009.03393*, 2020.
- Ren, Z., Shao, Z., Song, J., Xin, H., Wang, H., Zhao, W., Zhang, L., Fu, Z., Zhu, Q., Yang, D., et al. Deepseekprover-v2: Advancing formal mathematical reasoning via reinforcement learning for subgoal decomposition. arXiv preprint arXiv:2504.21801, 2025.
- Tooby-Smith, J. Formalization of physics index notation in lean 4. ArXiv, abs/2411.07667, 2024. URL https://api.semanticscholar.org/ CorpusID:273970152.
- Trinh, T. H., Wu, Y., Le, Q. V., He, H., and Luong, T. Solving olympiad geometry without human demonstrations. *Nature*, 625(7995):476–482, 2024.
- Tsoukalas, G., Lee, J., Jennings, J., Xin, J., Ding, M., Jennings, M., Thakur, A., and Chaudhuri, S. Putnambench: Evaluating neural theorem-provers on the putnam mathematical competition. *arXiv preprint arXiv:2407.11214*, 2024.

- van Doorn, F., Massot, P., and Nash, O. Formalising the h-principle and sphere eversion. In *Proceedings of the 12th ACM SIGPLAN International Conference on Certified Programs and Proofs*, CPP 2023, pp. 121–134, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400700262. doi: 10.1145/3573105.3575688. URL https://doi. org/10.1145/3573105.3575688.
- Wang, H., Unsal, M., Lin, X., Baksys, M., Liu, J., Santos, M. D., Sung, F., Vinyes, M., Ying, Z., Zhu, Z., et al. Kimina-prover preview: Towards large formal reasoning models with reinforcement learning. *arXiv preprint arXiv:2504.11354*, 2025.
- Whalen, D. Holophrasm: a neural automated theorem prover for higher-order logic. *arXiv preprint arXiv:1608.02644*, 2016.
- Wu, Z., Huang, S., Zhou, Z., Ying, H., Wang, J., Lin, D., and Chen, K. Internlm2. 5-stepprover: Advancing automated theorem proving via expert iteration on large-scale lean problems. *arXiv preprint arXiv:2410.15700*, 2024.
- Xin, H., Guo, D., Shao, Z., Ren, Z., Zhu, Q., Liu, B., Ruan, C., Li, W., and Liang, X. Deepseek-prover: Advancing theorem proving in llms through large-scale synthetic data. arXiv preprint arXiv:2405.14333, 2024a.
- Xin, H., Ren, Z., Song, J., Shao, Z., Zhao, W., Wang, H., Liu, B., Zhang, L., Lu, X., Du, Q., et al. Deepseek-prover-v1.5: Harnessing proof assistant feedback for reinforcement learning and monte-carlo tree search. arXiv preprint arXiv:2408.08152, 2024b.
- Yang, K., Swope, A., Gu, A., Chalamala, R., Song, P., Yu, S., Godil, S., Prenger, R. J., and Anandkumar, A. LeanDojo: Theorem proving with retrieval-augmented language models. *Advances in Neural Information Processing Systems*, 36:21573–21612, 2023.
- Ying, H., Wu, Z., Geng, Y., Wang, J., Lin, D., and Chen, K. Lean workbook: A large-scale lean problem set formalized from natural language math problems. *arXiv* preprint arXiv:2406.03847, 2024a.
- Ying, H., Zhang, S., Li, L., Zhou, Z., Shao, Y., Fei, Z., Ma, Y., Hong, J., Liu, K., Wang, Z., et al. Internlmmath: Open math large language models toward verifiable reasoning. *arXiv preprint arXiv:2402.06332*, 2024b.
- Zheng, K., Han, J. M., and Polu, S. Minif2f: a cross-system benchmark for formal olympiad-level mathematics. *arXiv preprint arXiv:2109.00110*, 2021.

A. False-positives with library search tactics

The apply? tactic⁸ is intended for searching the Lean library for lemmas applicable in the current context during interactive proof construction, similar to exact? and rw?. However, the apply? tactic must not be left in a final proof, as it has the same semantics as the sorry keyword, leaving such a proof incomplete.

For example, the proofs of putnam_2005_a4 and putnam_2007_b4 presented by DeepSeek-Prover-V2 (Ren et al., 2025) utilize the apply? tactic and are thus incorrect. This highlights the potential of reinforcement learning to exploit unchecked ambiguities in a proving environment. It also hints as to why a smaller DeepSeek-Prover-V2-7B model beat DeepSeek-Prover-V2-671B in some contexts related to cardinal numbers, as the former might have learned to exploit the apply? tactic while the latter did not (cf. Section *Skill Discovery by Reinforcement Learning* in Ren et al. (2025)).

In LeanTree, we prevent such errors by banning the usage of apply? by default, leaving it as opt-in for users who explicitly want to integrate it into their proving system.

B. Incremental proof verification

Originally, the Lean REPL did not utilize the Lean verification kernel to type-check assignments introduced by tactic applications, leading to incorrect proofs being accepted in some cases. Our prover managed to exploit these inconsistencies, finding 12 incorrect proofs in the MiniF2F validation set that passed verification in the REPL when using the tactic mode. One of these incorrect proofs is listed below.

```
theorem mathd_algebra_422 (x : \mathbb{R}) (\sigma : Equiv \mathbb{R} \mathbb{R}) (h<sub>0</sub> : \forall x, \sigma.1 x = 5 * x - 12)
(h<sub>1</sub> : \sigma.1 (x + 1) = \sigma.2 x) : x = 47 / 24 := by
rw [h<sub>0</sub>] at h<sub>1</sub>
by_contra! h
apply h.symm
simpa using h _
```

Lean correctly rejects the above proof with the error: "*don't know how to synthesize placeholder*". However, when executing the proof in the tactic mode using the old version of Lean REPL, the proof passes as correct.

In a recent update,⁹ Lean REPL was extended to type-check the full proof term after each tactic application, mitigating all false positives in the verification process, at the cost of introducing some drawbacks detailed below. We propose a different strategy: Instead of verifying the type-correctness of the whole proof term, type-check only the new assignments introduced by a tactic application.

To illustrate this difference, assume that we are searching for a proof of type T and have already constructed a term t : T which still contains unassigned metavariables of types G_1, \ldots, G_n . Next, assume that we apply a tactic assigning terms $g_1 : G_1, \ldots, g_n : G_n$ to the metavariables in t. The current strategy of Lean REPL is to type-check the whole assignment t : T, which contains the new terms g_1, \ldots, g_n as subterms.

This leads to issues when the proof bifurcates into two branches. For example, consider applying the tactic have k := sorry in the REPL's tactic mode.¹⁰ This corresponds to the assignment of the term ((fun $k \mapsto ?_{-}$) sorry) to the metavariable corresponding to the main goal, where ?_ is a new unassigned metavariable. Lean REPL then enables solving the ?_ and sorry parts independently. However, assignments made in one of these two branches are not propagated to the other one, meaning that, from the perspective of each branch, the root proof term will always contain at least one hole introduced by an unassigned metavariable or the sorry keyword. This leads to the kernel always failing during a type check for the duration of this proof. The end effect is that the current verification strategy leads to false negatives when the proof being verified contains at least one syntactic branching.

In contrast, the proposed strategy implemented in LeanTree type-checks only the new assignments $g_1 : G_1, \ldots, g_n : G_n$ introduced by a tactic application. The type correctness of the whole proof term t : T then follows transitively. Most importantly, this new verification strategy mitigates the false negatives described in the previous paragraph. Additionally,

⁸https://lean-lang.org/doc/reference/latest/Tactic-Proofs/Tactic-Reference/#apply__

⁹Published on GitHub on Apr 10, 2025.

¹⁰This is the standard way to use the have tactic in the REPL – see for example https://github.com/leanprover-community/repl/blob/master/test/have_by_sorry.in

the new strategy also decreases the time complexity of verifying a proof from quadratic to linear since each sub-assignment is only verified once.

Note that all proofs found by our prover are independently verified directly using Lean, bypassing Lean REPL together with the new verification algorithm described in this section.

C. Experiments setup

In all experiments, we prompt the Llema-7B model using the following prompt.

```
Complete the given Lean 4 code.
Directly output the completed Lean 4 code without any additional text, comments, or reasoning.
```

When performing linear rollouts, we manually select only the first tactic from the model output.

During generation, we disable sampling of the sorry and admit keywords, together with syntax for comment sections.

D. Dataset schema

The dataset is distributed in two files corresponding to the two data sources:

- lean-trees_mathlib.jsonl
- lean-trees_deepseek-prover-v1.jsonl

Both of these files share the following unified schema.

```
<sample> ::= {
1
     "path": <string>,
2
     "imports": [<string>],
3
     "theorems": [<error> | {
4
5
       "span": <span>,
       "name": <string?>,
6
       "context": [<string>],
7
8
       "by_blocks": [{
9
          "tree": <error> | {
            "root": <proof_node>
10
11
          }
       }]
12
13
     }]
14
   }
15
   <proof_node> ::= {
16
     "id": <string>,
17
     "proof_size": <int>,
18
     "proof_depth": <int>,
19
     "tactic": {
20
       "tactic_string": <string>,
21
       "span": <span>,
22
       "children": [<string>],
23
       "tactic_depends_on": [<string>]
24
     }
25
     "state": {
26
27
       "goals": [{
         "tag": <string?>,
28
          "type": <string>,
29
          "hypotheses": [{
30
            "type": <string>,
31
            "user_name": <string>,
32
            "value": <string?>
33
34
         }]
       }]
35
     }
36
  }
37
38
   <span> ::= {
39
40
     "start": <int>,
     "finish": <int>
41
   }
42
43
   <error> ::= {
44
45
     "error": <string>
46
   }
```