# USEFOR: USEFULNESS-DRIVEN LEARNING OF FOR-MAL MATHEMATICS

**Anonymous authors**Paper under double-blind review

000

001

002003004

010 011

012

013

014

016

017

018

019

021

024

025

026

027

028

029

031 032 033

034

037

038

040

041

042 043

044

046

047

048

051

052

#### **ABSTRACT**

Creating an AI that can truly "do" mathematics requires more than just solving isolated problems. It must mimic the creative, progressive nature of human mathematicians, who build upon previous work to generate new knowledge. A crucial part of this process is proposing theorems that serve as useful building blocks for proving more advanced theorems. In this paper, we introduce UseFor, a novel framework that formalizes this notion of usefulness and demonstrates how it can be used to train a usefulness-driven AI mathematician. UseFor determines a theorem's usefulness based on two criteria: its reusability in subsequent proofs and its contribution to increasing proof likelihood. We integrate UseFor into the self-play conjecturing-and-proving setting of Minimo ((Poesia et al., 2024)). That is, starting from only axioms, we iteratively train conjecturers to propose useful formal statements and provers that explicitly reuse them when generating formal proofs. We experimentally evaluate this usefulness-driven self-play approach across three mathematical domains: arithmetic, propositional logic, and group theory. Our evaluation considers two metrics: intrinsic usefulness, which measures how often our trained provers reuse theorems, and extrinsic usefulness, judged by a state-ofthe-art large language model and external provers like SMT solvers. Our results demonstrate that our usefulness-trained model effectively generates a large number of intrinsically and extrinsically useful formal theorems. For instance, our approach outperforms the original Minimo by 2.9 times in extrinsic usefulness for arithmetic. Our work highlights the significant potential of integrating usefulness in AI-driven mathematical discovery.

# 1 Introduction

Mathematical reasoning has long stood as a frontier challenge for artificial intelligence. (Newell & Simon, 1956) While large language models (LLMs) have achieved rapid progress in formal theorem proving (Yang et al., 2024), most approaches depend on extensive human-written corpora of proofs and conjectures (Yang et al., 2023; Ying et al., 2025). This dependence limits the domains in which they can operate and prevents them from advancing beyond existing human knowledge. By contrast, human mathematicians build knowledge by conjecturing new statements and proving them, gradually extending their theoretical landscape without external supervision.

A natural question, therefore, arises: can we replicate this process automatically? Specifically, can an artificial agent, starting only from axioms, learn via self-play between conjecturing and proving, bootstrapping its own knowledge and progressively discovering new mathematics? This paradigm (McAllester, 2020) would eliminate the need for human data, enable exploration of domains where no proofs exist, and produce a scalable source of synthetic theorems for training future provers. Recent work such as (Poesia et al., 2024) shows early signs of this vision: by iterating between conjecture generation and Monte Carlo tree search (MCTS)-guided proof search, Minimo gradually learns to prove increasingly more difficult statements from scratch. However, as we argue below, merely increasing difficulty as defined in Minimo is not enough to drive true theory building.

Minimo, like most self-play provers, evaluates conjectures purely by their difficulty—the negative log-probability of successful proofs under the current prover policy. This encourages the generation of challenging statements, but overlooks whether they actually help prove other results. In practice,

we find that conjectures promoted solely by difficulty are rarely reused in later proofs and offer little leverage for solving harder targets.

We are inspired here by Bengio & Malkin (2024)'s perspective that the value of a theorem lies in its connection to other theorems. They noted that "a crucial component of the usefulness of a new proven theorem t (in the context of previous theorems  $\mathcal{T}(S)$ ) is how efficiently  $\mathcal{T}(S) \cup \{t\}$  compresses the set of all provable mathematical statements  $\mathcal{M}$ ". Tao (2007) makes a similar point, observing that the strength of a theorem is best judged by "testing it against a class of questions and problems that the theorem is intended to assist with solving". Both perspectives highlight that valuable conjectures should be judged not in isolation but advance the prover's qualitative capabilities and accelerate cumulative theory building.

This motivates our core idea: to advance automated theory building, we need a tractable metric that approximates this relational value of conjectures. We introduce such a metric, UseFor, which defines a conjecture as *useful* if it both appears in the proof of a downstream target and increases the prover's success likelihood (log-probability) on that target. This dual criterion excludes trivial tautologies with no proving power and narrowly phrased statements with little applicability in other proofs, yielding a practical proxy for Bengio & Malkin (2024)'s compression perspective.

We leverage UseFor to develop a usefulness-aware self-play loop that builds directly on Minimo (Poesia et al., 2024). After each round of conjecturing and proving, we use UseFor to identify proven conjectures with the most relational value for proving other conjectures. These useful theorems are made available to the proving model as lemmas and weighted more heavily in training, guiding future conjectures and proofs toward structures that accelerate cumulative theory building.

We evaluate our usefulness-driven self-play framework on three domains: arithmetic, propositional logic, and group theory. Our results show that our approach generates a significant number of theorem usages and more useful conjectures compared to Minimo based on LLM-as-a-judge. For instance, for arithmetics, our approach outperforms Minimo by 2.9 times in producing useful conjectures. These findings indicate that usefulness is a stronger intrinsic signal than difficulty for guiding conjecture generation, and that usefulness-aware self-play offers a scalable path toward data-free theory exploration.

#### **Main contributions.** Our paper makes the following contributions:

- We formalize the notion of theorem usefulness as a dual criterion of *usage* and *improvement*, and propose a tractable procedure for measuring it within self-play (Section 3.2.1).
- We introduce a usefulness-aware self-play loop that augments Minimo by selecting conjectures according to relational usefulness rather than difficulty (Section 3.2.2).
- We present stabilization techniques (triviality filtering and novelty bias) that keep the loop from collapsing into tautologies or memorized variants (Section 3.2.3).
- We provide empirical results across arithmetic, propositional logic, and group theory, demonstrating that usefulness-driven conjectures are more reusable and lead to higher prover success rates than difficulty-based baselines (Section 4).

# 2 RELATED WORK

Our work is primarily related to prior bodies of work on mathematical conjecturing, tactic discovery, and theory exploration. Our approach is distinguished by the fact that our model is trained in a tabula rasa fashion, without any pre-existing examples, and evaluated on the theory exploration task.

Mathematical conjecturing. Our work is most closely based on Minimo (Poesia et al., 2024), which proposes a theorem-proving model in the Peano (Poesia & Goodman, 2023) formal language that is trained through iterative conjecturing and proving from scratch. (Polu & Sutskever, 2020) also propose a model that is trained via self-play, while (Dong & Ma, 2025) demonstrate the ability of the iterative conjecturing-proving paradigm to enhance a pretrained theorem prover. However, these works only use conjecturing as a means to improve the proof-search capabilities of the model, and do not attempt to evaluate the conjecturing abilities of the model directly. LeanConjecturer (Onda et al., 2025) proposes a model specifically designed for the conjecturing task, but uses a

pretrained LLM; in doing so, the ability of the LeanConjecturer model to generate novel conjectures cannot be faithfully evaluated due to inevitable contamination from pre-training data. Compared to these works, our approach evaluates conjecturing as a stand-alone task, while our tabula rasa setting allows us to definitively confirm the novelty of conjectures generated by our model.

**Tactic and premise discovery.** There is also a body of work concerning the task of tactic discovery, which aims to construct tactics in an interactive theorem prover setting that simplify proofs or otherwise enhance proving capabilities. TacMiner (Xin et al., 2025) proposes a method to find tactic simplifications in RCoq, given an existing high-quality corpus of proofs. Lego-Prover (Wang et al., 2023) and Seed-Prover (Chen et al., 2025) use already proven lemmas as a way to strengthen a theorem proving model, in the Isabelle and Lean 4 settings, respectively. However, all of these approaches require a dataset of high-quality, human-generated proofs, while our approach generates useful premises from scratch.

Theory exploration using machine learning. Finally, a third body of work is theory exploration using ML methods, the task of formulating interesting conjectures about a given problem domain (Johansson & Smallbone, 2021). We consider this problem to be the one our work addresses most closely. While a number of classical and neural approaches have been proposed for this task, existing neural methods work by training or finetuning a model based on an existing proof corpus (Urban & Jakubův, 2020). Lemmanaid (Alhessi et al., 2025) uses neuro-symbolic methods by finetuning a model with a subset of an existing proof library, and then evaluating it on another subset of conjectures. In search of a purely intrinsic approach in order to discover how a model could discover this usefulness without relying on human data, we distinguish ourselves by not training on external data, an approach similar to what has been done for SMT solvers (Gauthier & Urban, 2025).

# 3 METHODOLOGY

#### 3.1 BASE SELF-PLAY FRAMEWORK (MINIMO)

A central challenge in building autonomous theorem-proving agents is the lack of human-labeled data. Unlike natural language or code, formal mathematics has limited corpora, and many target domains have essentially no prior datasets. There are two main approaches to address this bottleneck. One line of work leverages autoformalization, which translates large volumes of informal mathematics into formal statements, an approach that has already shown promise in practice. A more first-principles alternative is to remove reliance on pre-existing data entirely by using self-play: coupling<sup>1</sup> a conjecturer, which proposes candidate statements, with a prover, which attempts to establish them. Through repeated interaction, both components improve jointly in a closed loop, enabling progress even in domains with no human supervision or existing corpus.

Minimo (Poesia et al., 2024), implemented in the Peano environment (Poesia & Goodman, 2023), instantiates this idea. Starting from axioms alone, it alternates between conjecture generation and proof search. Over time, the prover strengthens by training on successful proof traces, while the conjecturer adapts toward statements near the boundary of provability. This process yields an automatically generated curriculum of increasing difficulty, with no reliance on human annotations. We summarize its core components next, as they provide the foundation on which our method builds.

#### 3.1.1 Conjecturing

The conjecturer  $C_{\theta}$ , with  $\theta$  denoting the model parameters, generates statements in the Peano language, a dependently typed tree-based formal system with a finite action space (Poesia & Goodman, 2023). Each formula is represented as a well-formed term tree. To prevent invalid formulas, Minimo employs constrained decoding Poesia et al. (2021): at each step, candidate tokens are filtered so that only those extending the current tree into a valid continuation remain. This guarantees both syntactic and semantic validity, and prevents wasting prover efforts on malformed statements.

<sup>&</sup>lt;sup>1</sup>In practice, coupling means using a single underlying model for both conjecturer and prover roles, as proposed by Poesia et al. (2024).

#### 3.1.2 PROOF SEARCH

The prover  $\mathcal{P}_{\theta}$  attempts to establish each conjecture using Monte Carlo tree search (MCTS), which is well suited to the combinatorial branching structure of formal proofs. A proof state s encodes the current context, including the set of assumptions, the active subgoal, and any partially completed steps. At each state, the Peano environment defines the finite set of admissible inference actions. Guided by the prover's policy  $\pi_{\theta}(a \mid s)$  and value estimates, MCTS expands trajectories

$$\tau = (s_0, a_0, s_1, a_1, \dots, s_T)$$

starting from the initial state  $s_0$ . A completed trajectory corresponds to a valid proof of the conjecture. Its score is the log-likelihood under the prover's policy,

$$\ell(c) = \log p_{\theta}(\tau \mid c) = \sum_{t=0}^{T-1} \log \pi_{\theta}(a_t \mid s_t).$$

Less negative values of  $\ell(c)$  indicate that the proof was expected under the current policy, while more negative values correspond to surprising but ultimately valid proofs. To exploit partial progress, Minimo applies *hindsight relabeling*: even when a conjecture cannot be proved in full, explored search trees are decomposed into valid subtraces corresponding to intermediate lemmas, which are then incorporated as additional training data (Poesia et al., 2024). This enlarges the training set and recycles computation that would otherwise be wasted on failed proofs.

# 3.1.3 Conjecturing-Proving Self-Play Loop

The conjecturer  $C_{\theta_i}$  and prover  $\mathcal{P}_{\theta_i}$  interact in an iterative loop. At iteration i, the conjecturer samples a batch of N candidate statements

$$Q_i = \{c_1, \ldots, c_N\} \sim C_{\theta_i}(\cdot \mid \mathcal{T}_i),$$

where  $\mathcal{T}_i$  is the current theory consisting of axioms and previously promoted lemmas. For each  $c \in \mathcal{Q}_i$ , the prover attempts to establish it via MCTS:

$$(\operatorname{proof}(c), \ell(c), \operatorname{trace}(c)) \leftarrow \operatorname{MCTS\_PROVE}(c; \mathcal{T}_i, \mathcal{P}_{\theta_i}),$$

where  $\operatorname{proof}(c)$  is a complete proof trajectory  $\tau$  if one is found (or  $\varnothing$  otherwise),  $\operatorname{trace}(c)$  is the explored search tree, and  $\ell(c)$  is the log-likelihood of the trajectory under the prover's policy.

Conjectures are then stratified by empirical difficulty. Let  $S_i = \{c \in Q_i : \operatorname{proof}(c) \neq \emptyset\}$  be the set of successful conjectures, and let  $q_{20}$  and  $q_{50}$  denote the 20th and 50th percentiles of  $\{\ell(c) : c \in S_i\}$ . Labels are assigned to each conjecture c as

$$\mathrm{label}(c) = \begin{cases} \mathrm{``fail''}, & \mathrm{if\ proof}(c) = \varnothing, \\ \mathrm{``hard''}, & \mathrm{if\ } \ell(c) < q_{20}, \\ \mathrm{``easy''}, & q_{20} \leq \ell(c) < q_{50}, \\ \mathrm{``trivial''}, & \ell(c) \geq q_{50}. \end{cases}$$

The dataset for iteration i is then

$$\mathcal{E}_i = \{(\operatorname{trace}(c), \operatorname{label}(c)) : c \in \mathcal{Q}_i\},\$$

which aggregates conjectures, proofs when available, and hindsight-relabeled subproofs extracted from failed searches. Both  $C_{\theta}$  and  $\mathcal{P}_{\theta}$  are updated on  $\mathcal{E}_{i}$ , creating a feedback loop: the conjecturer shifts toward generating statements just beyond the prover's current reach, while the prover expands its competence from the resulting proofs. This difficulty-driven loop is the foundation upon which we build in Section 3.2, where difficulty is replaced with a more relational signal of usefulness.

#### 3.2 USEFULNESS-AWARE SELF-PLAY LOOP

The self-play framework of Minimo provides a compelling basis for data-free theory exploration: starting from axioms, conjecturing and proving improve together in a bootstrapping loop. However, its training signal is limited to conjectural difficulty, measured as the negative log-probability of a

proof under the current prover. While effective for generating a curriculum of harder statements, this signal is ultimately syntactic. It rewards conjectures that are improbable under the model's local policy, but does not account for whether they connect meaningfully to other theorems explored so far. As a result, the system often promotes conjectures that are labeled as "hard" but not necessarily useful statements: isolated identities that stretch the prover temporarily but are rarely reused and add little structure to the theory. The log-probability score treats difficulty as an end in itself, overlooking the relational role that lemmas play in enabling further proofs and sustaining theory growth.

To address this limitation, we introduce a usefulness-based self-play loop. Instead of ranking conjectures solely by syntactic hardness, we ask whether incorporating a new lemma changes the prover's future behavior, specifically, whether it makes other statements easier to prove. Conjectures that are both provable and demonstrably beneficial in downstream proofs are promoted into the growing library, and their traces are used to train both the conjecturer and the prover. This shifts the learning objective from accumulating difficult but isolated statements to building a network of reusable ones, better aligned with the cumulative nature of mathematical discovery.

#### 3.2.1 Definition of the Usefulness Metric

The perspectives of Bengio & Malkin (2024) and Tao (2007) converge on the idea that the value of a theorem is relational: it derives its significance not from truth alone, but from its effect on subsequent reasoning. Yet they articulate this in complementary registers. Bengio & Malkin (2024) frames usefulness in information-theoretic terms, proposing that a theorem acts as a *compression primitive*—its addition to a base theory reduces the description length of other proofs. Tao (2007) instead emphasizes the pragmatic dimension: the *strength* of a theorem is revealed only by confronting new problems and observing the range of arguments it simplifies.

While these views are philosophically aligned, neither directly yields a metric implementable within a self-play loop. Compression, though elegant, requires comparing description lengths over the unbounded space  $\mathcal{M}$  of all provable statements, which is an intractable quantity in practice. Tao (2007)'s criterion, by contrast, presupposes a human mathematician's judgment in selecting "a class of questions and problems" against which to test strength. What is missing is a procedure that preserves the spirit of both notions while remaining computable for a prover–conjecturer system.

Our contribution is to bridge this gap by constructing an operational proxy for usefulness that can be applied iteratively inside the self-play loop. At a high level, the metric estimates a conjecture's capacity to expand the prover's effective reach: conjectures are useful insofar as their availability systematically reduces the effort of proving a benchmark set of targets.

Formally, let  $\mathcal B$  be a benchmark set consisting of theorems that are difficult, but not impossible, for the prover to prove. For each  $b \in \mathcal B$ , let  $p_{\theta}(\tau_b \mid b)$  denote the prover's probability of producing a proof trajectory  $\tau_b$  under theory  $\mathcal T$ , and let  $p'_{\theta}(\tau_b \mid b)$  denote the same quantity when a candidate lemma  $\ell$  is available. We say that  $\ell$  is useful if there exists  $b \in \mathcal B$  such that

(i) 
$$\ell$$
 is invoked in the proof trace of  $b$ , and (ii)  $\log p'_{\theta}(\tau_b \mid b) - \log p_{\theta}(\tau_b \mid b) > 0$ .

Both conditions are essential: usage without improvement admits trivial tautologies such as  $\forall x. \ x = x$ , which the prover may frequently attempt but which yield no real progress. Only when the two conditions coincide do we identify lemmas that are genuinely structural.

To evaluate this criterion efficiently, we do not re-prove  $\mathcal{B}$  for every lemma in isolation. Instead, given a set of newly proved conjectures  $\mathcal{C}$ , we subsample a subset of size  $\lceil \sqrt{|\mathcal{C}|} \rceil$  and temporarily add them to the context. Each  $b \in \mathcal{B}$  is then re-proved once under this extended theory. If a candidate  $\ell$  appears in the proof of b and the resulting log-likelihood improves relative to baseline, the gain is attributed to  $\ell$ . The aggregate score

$$U(\ell) = \sum_{b \in \mathcal{B}} \mathbf{1} \{ \ell \in \operatorname{proof}(b) \} \cdot \max\{0, \log p_{\theta}'(\tau_b \mid b) - \log p_{\theta}(\tau_b \mid b) \}$$

is then used to rank candidates. Only the top  $\rho$  fraction are promoted to the library, together with their associated proofs and hindsight traces. This provides a tractable mechanism for selecting conjectures that repeatedly demonstrate both reuse and measurable downstream gains.

**Illustrative scenario.** Consider arithmetic with multiplication defined inductively. Early in training, the prover may not yet know lemmas such as  $x \times (y+1) = xy+x$ . Without this fact, even simple

targets like  $(x+1) \times (y+1) = xy + x + y + 1$  require long derivations by repeatedly unfolding the definition of multiplication. Once  $x \times (y+1) = xy + x$  is conjectured and proved, however, it can be applied directly, and many small multiplication-addition identities shorten dramatically. Our metric marks this lemma as useful precisely because it is both *used* in subsequent proofs and its presence *improves* the prover's success probability. Later discoveries, such as distributivity, compound this effect across broader families, but it is these intermediate stepping-stone lemmas that first enable steady cumulative progress.

#### 3.2.2 Training Loop with the Usefulness Metric

We now describe how the usefulness metric is integrated into the conjecturing–proving loop. The outer structure mirrors Minimo (Poesia et al., 2024): in each iteration the agent generates conjectures, the prover attempts proofs via MCTS, and traces are collected. The crucial difference lies in how conjectures are filtered, promoted, and fed back into training. Whereas Minimo labels conjectures by proof log-probability percentiles and emphasizes those deemed "hard," our framework evaluates conjectures by their relational usefulness.

At iteration i, the conjecturer first proposes a batch  $C_i$ , which is passed through a triviality filter to remove vacuous identities. Each conjecture is then attempted under the current theory  $\mathcal{T}_i$  using MCTS, producing proofs, log-likelihoods, and hindsight examples. Following Minimo, conjectures are provisionally bucketed into "hard," "easy," and "trivial" categories by percentile of log-likelihood. Nonfailing conjectures are collected as candidate lemmas.

The key departure comes in how the "hard" subset is treated. Rather than promoting them directly, we apply the *usefulness test* using the current pool of non-failing lemmas as context. A random subsample  $L_i \subseteq \mathcal{H}_i$  of size  $\lceil \sqrt{|\mathcal{H}_i|} \rceil$  is drawn, and each benchmark  $b \in \mathcal{B}$  is re-proved under both the baseline theory  $\mathcal{T}_i$  and the augmented theory  $\mathcal{T}_i \cup L_i$ . If a lemma  $\lambda \in L_i$  is invoked in the augmented proof of b and improves its log-likelihood relative to the baseline, the gain is added to its cumulative usefulness score  $U_i(\lambda)$ . Candidates are then ranked by  $U_i(\lambda)$ , and only the top  $\rho$  fraction are promoted into the persistent theory  $\mathcal{T}_{i+1}$ . Because  $L_i$  is resampled at every iteration, different subsets of candidates are tested over time, so all conjectures eventually receive usefulness credit.

Finally, the training dataset  $\mathcal{E}_i$  is assembled. It contains the conjectures, their proofs, and percentile labels, and the hindsight traces from the base loop, as well as the useful lemmas and re-proving examples produced during usefulness testing. The agent is updated on  $\mathcal{E}_i$ , and the promoted lemmas are added to  $\mathcal{T}_{i+1}$  for future iterations.

In summary, Minimo's curriculum is driven by proof difficulty under the current prover, whereas our loop is driven by demonstrable downstream impact. Only conjectures that are both *used* and *improve* benchmark proofs are promoted, producing a library that is not just deeper but more interconnected, with lemmas reappearing across proofs and compounding overall success.

#### 3.2.3 OTHER IMPROVEMENT TECHNIQUES

Although the usefulness loop provides a stronger supervisory signal than difficulty alone, we observed recurrent failure modes in practice. In particular, the conjecturer may become trapped in local minima, repeatedly generating trivial identities or minor variants of existing statements. This behavior resembles exploiting shortcuts rather than genuine advancements. To enhance robustness and better align training with the intended objectives, we incorporate two additional techniques.

**Triviality filtering.** Before proof attempts, we remove conjectures that match heuristic patterns such as tautologies (x = x) or constant-only identities (0 + 1 = 1). Such statements can be discharged immediately without reasoning, yet they satisfy the usage criterion and would otherwise dominate the usefulness signal. Filtering them prevents the prover cycles from being wasted and prevents the model from collapsing toward vacuous but spuriously rewarding lemmas.

**Novelty bias.** To encourage structural diversity, we penalize the conjecturer for generating statements that share long prefixes with previously generated conjectures. This discourages local memorization and pushes exploration toward unexplored syntactic regions, thereby increasing the likelihood of uncovering genuinely new lemmas that expand the theorem library.

These techniques do not alter the usefulness metric itself, but regularize the conjecturer's proposal distribution. By filtering trivialities and discouraging near-duplicates, the system avoids spurious short-term rewards and maintains pressure toward conjectures that are both novel and reusable. Empirically, they improve the stability of the usefulness-aware loop, allowing the training distribution to shift steadily toward conjectures that promote cumulative theory building.

#### 4 EXPERIMENTAL EVALUATION

We now evaluate UseFor on three mathematical domains: (i) arithmetic, (ii) propositional logic, and (iii) group theory. The axioms of these domain are presented in Appendix A. Our goal is to assess whether UseFor demonstrates the essential qualities of a desirable reasoning system: (a) the ability to accumulate knowledge across iterations, (b) the ability to generate conjectures of intrinsic value, and (c) whether our usefulness-driven training is necessary. In more detail, we would like to address the following research questions:

- **RQ1**: Can the prover reuse theorems proven in previous iterations to prove current conjectures? Reuse is essential for cumulative theory building: without it, a system risks repeatedly rediscovering tautologies or isolated results, rather than developing an interconnected body of theory.
- **RQ2**: Do likelihoods of theorem-reusing proofs increase across multiple iterations? This would signify that during the training process, the prover is gradually gaining more capabilities and confidence in theorem reuse.
- **RQ3**: Are the conjectures useful beyond self-play? Extrinsic usefulness tests whether the system discovers theorems a mathematician would value, rather than artifacts of the training loop. In contrast, the difficulty metric evaluated in Poesia et al. (2024) is not suited for this purpose.
- **RQ4**: *Is the usefulness metric essential for both conjecturing and proving (ablation study)?* Without it, do either component struggle to discover interesting theorems? This matters because the entire training loop relies on this metric as its guiding signal.

**Evaluation Metrics.** In light of the above interesting research questions, we employ two complementary metrics designed to capture structural usefulness:

- Intrinsic usefulness: measured as the number of times a previously proven theorem is reused during usefulness testing. A high score indicates that the system is both conjecturing and successfully reusing theorems in its own proving process.
- Extrinsic usefulness: measured via an LLM-as-judge (GPT-4.1), which rates conjectures for mathematical value after a deduplication step that removes near-duplicates (details can be found in Appendix B). We also require that the conjectures can be proved by an external automated prover based on the Z3 STM solver (De Moura & Bjørner, 2008), which is effective on our self-playgenerated conjectures. This metric evaluates whether conjectures would be judged useful by a human mathematician, beyond the system's internal dynamics.

Prior work, such as MINIMO, evaluated progress using proof success rates. Yet this metric is easily inflated by trivial statements (e.g., tautologies) that are provable but add nothing to cumulative reasoning, which has been consistently observed in experiments. Our intrinsic and extrinsic usefulness metrics avoid this issue by directly testing whether conjectures are structurally valuable. Since MINIMO did not consider using previously proven theorems, its intrinsic metric is identically zero; we therefore include it only as a baseline in RQ2 and RQ3.

**Experiment Setup.** All models were trained for 10 iterations, with 200 conjectures generated per iteration. Each experiment was repeated three times, and we report mean values with standard deviations across runs. Proof search was conducted using Monte Carlo Tree Search (MCTS) with a budget of 1000 expansions per conjecture. We repeat all experiments for three times and report averaged results to account for stochastic variations.

(RQ1) The model reuses previously proven conjectures. Reuse is a key indicator of cumulative reasoning: a system that fails to apply previously proven theorems risks stagnating in isolated rediscoveries, rather than developing an interconnected theory. In our experiments, UseFor shows a

steady increase in lemma usage during usefulness testing (Figure 1). Although the first few iterations provide little signal, usage accelerates in later iterations, demonstrating that the model progressively conjectures more useful theorems and becomes increasingly capable of applying them. This trend is consistent across all domains, and we expect it to persist with additional iterations.

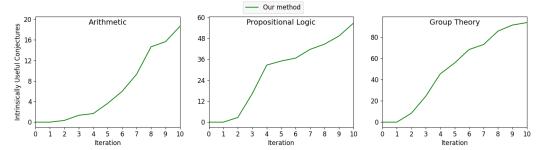


Figure 1: Intrinsic Evaluation: Total theorem use count with increasing iterations.

(RQ2) The model grows increasingly more confident in theorem reuse. As training progresses, our model grows increasingly confident in its use of previously conjectured lemmas, as evaluated by the average log-probability of proofs where at least one previously conjecture was used (Figure 4). This aligns with the significant increase in intrinsically useful conjectures across multiple iterations, as shown in Figure 1, and shows that the UseFor training objective is effective in encouraging the model to use previously proven lemmas. As we notice an upwards trend as iterations continue, this also demonstrates that our lemmas become more difficult to prove as time goes on, as earlier provers assign low probabilities to them.

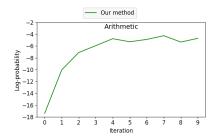


Figure 2: Average log-probabilities of proofs where a previously conjecture was used, across prover iterations.

(RQ3) The conjectures are extrinsically useful. Extrinsic usefulness tests whether conjectures would be judged valuable outside the self-play loop (Figure 3). In early iterations, UseFor quickly identifies many "easy" theorems accessible through shallow search. Crucially, usefulness continues to increase in later iterations, indicating that the system discovers progressively deeper and less trivial results. The growth of this metric is super-linear, suggesting that UseFor not only maintains but amplifies its ability to generate theorems a mathematician would regard as useful. In Appendix B.4, we provide examples of extrinsically useful theorems conjectured by our model.

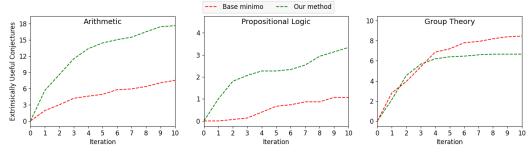


Figure 3: Extrinsic Evaluation: Number of deduplicated useful theorems per iteration, as determined by GPT-4.1 as a judge and proved by an SMT solver.

(**RQ4**) **Usefulness training is necessary.** This experiment evaluates how our usefulness training signal affects performance (Figure 4). We focus here on the domain of arithmetic, though the same pattern holds in the other domains. We compare UseFor with an improved version of Minimo, named "No usefulness training", which has access to prior proven conjectures but is not trained

to reuse them. As shown in Figure 4, if training is omitted, the system performs markedly worse: extrinsically, fewer theorems are judged to be useful by LLM-as-a-judge and SMT solver. This demonstrates the importance of training for updating the conjecturer with usefulness feedback steers it toward generating conjectures that are genuinely valuable for future proofs.

# 5 CONCLUSION

We studied automated conjecturing from minimal axioms as a prerequisite to scalable, self-improving theorem proving. While prior self-play systems such as (Poesia et al., 2024) use difficulty (low proof log-probability) as the sole training signal, we argued that difficulty alone is insufficient for theory building. We introduced a usefulness-aware self-play framework that evaluates conjectures by their downstream impact: whether they are actually reused in subsequent proofs and whether their inclusion increases the success likelihood of proving other targets. This dual criterion operationalizes the intuition that valuable theorems function as *compression primitives* for mathematics, turning isolated wins into reusable structure. Integrated into the self-play loop, the metric selects, promotes, and trains on lemmas that reshape future proof search through usefulness.

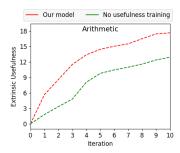


Figure 4: Ablation showing the necessity of usefulness training.

Across arithmetic, propositional logic, and group theory, UseFor steadily increases both the intrinsic reuse of conjectured theorems and the extrinsic usefulness of its discoveries. Performance improves over successive iterations, with the system progressing from "easy" lemmas to conjectures requiring deeper proofs. Ablation studies show that training both the prover and conjecturer with usefulness feedback is necessary: removing either sharply reduces both intrinsic and extrinsic metrics. Together, these findings confirm that usefulness-aware self-play can build coherent and cumulative theories directly from axioms.

Limitations and future work. Our study is confined to relatively small models, limited axioms, and fixed search budgets; scaling to richer foundations (e.g., Lean, Isabelle) and larger models remains an open challenge. However, applying approaches like UseFor on bigger models, such as large language models pretrained large corpora, brings the novel risk of data contamination. Therefore, future work should consider how to mitigate such data leakage concerns in the setting where LLMs are used for conjecturing new mathematical theorems.

Our usefulness test requires additional proof attempts; amortizing this cost (e.g., via certain advanced sampling techniques) is an important engineering direction. While our usage and improvement criterion balances applicability and effect size, it is still a proxy; developing tighter information-theoretic approximations to compression seems promising.

# ETHICS STATEMENT

This work adheres to the ICLR Code of Ethics. In this study, no human subjects was involved. We ensured compliance with all relevant dataset licenses. The experimental outcomes did not have security privacy and security concerns.

#### REPRODUCIBILITY STATEMENT

We are committed to ensuring the reproducibility of our work. We provide detailed experimental setup in Section 4 and Appendices A and B. We will publicly release all our code, models, and evaluation logs upon paper acceptance.

# REFERENCES

Yousef Alhessi, Sólrún Halla Einarsdóttir, George Granberry, Emily First, Moa Johansson, Sorin Lerner, and Nicholas Smallbone. Lemmanaid: Neuro-symbolic lemma conjecturing, 2025. URL

```
486 https://arxiv.org/abs/2504.04942.
```

- Yoshua Bengio and Nikolay Malkin. Machine learning and information theory concepts towards an ai mathematician. *Bulletin of the American Mathematical Society*, 61(3):457–469, 2024.
- Luoxin Chen, Jinming Gu, Liankai Huang, Wenhao Huang, Zhicheng Jiang, Allan Jie, Xiaoran Jin, Xing Jin, Chenggang Li, Kaijing Ma, Cheng Ren, Jiawei Shen, Wenlei Shi, Tong Sun, He Sun, Jiahui Wang, Siran Wang, Zhihong Wang, Chenrui Wei, Shufa Wei, Yonghui Wu, Yuchen Wu, Yihang Xia, Huajian Xin, Fan Yang, Huaiyuan Ying, Hongyi Yuan, Zheng Yuan, Tianyang Zhan, Chi Zhang, Yue Zhang, Ge Zhang, Tianyun Zhao, Jianqiu Zhao, Yichi Zhou, and Thomas Hanwen Zhu. Seed-prover: Deep and broad reasoning for automated theorem proving, 2025. URL https://arxiv.org/abs/2507.23726.
- Leonardo De Moura and Nikolaj Bjørner. Z3: An efficient smt solver. In *International conference* on Tools and Algorithms for the Construction and Analysis of Systems, 2008.
- Kefan Dong and Tengyu Ma. Stp: Self-play llm theorem provers with iterative conjecturing and proving, 2025. URL https://arxiv.org/abs/2502.00212.
- Thibault Gauthier and Josef Urban. Learning conjecturing from scratch, 2025. URL https://arxiv.org/abs/2503.01389.
- Moa Johansson and Nicholas Smallbone. Conjectures, tests and proofs: An overview of theory exploration. *Electronic Proceedings in Theoretical Computer Science*, 341:1–16, September 2021. ISSN 2075-2180. doi: 10.4204/eptcs.341.1. URL http://dx.doi.org/10.4204/EPTCS.341.1.
- David McAllester. Mathzero, the classification problem, and set-theoretic type theory, 2020. URL https://arxiv.org/abs/2005.05512.
- Allen Newell and Herbert A. Simon. The logic theory machine—a complex information processing system. *IRE Transactions on Information Theory*, 2(3):61–79, 1956.
- Naoto Onda, Kazumi Kasaura, Yuta Oriike, Masaya Taniguchi, Akiyoshi Sannai, and Sho Sonoda. Leanconjecturer: Automatic generation of mathematical conjectures for theorem proving, 2025. URL https://arxiv.org/abs/2506.22005.
- Gabriel Poesia and Noah D. Goodman. Peano: learning formal mathematical reasoning. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 381(2251), June 2023. ISSN 1471-2962. doi: 10.1098/rsta.2022.0044. URL http://dx.doi.org/10.1098/rsta.2022.0044.
- Gabriel Poesia, Alex Polozov, Vu Le, Ashish Tiwari, Gustavo Soares, Christopher Meek, and Sumit Gulwani. Synchromesh: Reliable code generation from pre-trained language models. In *International Conference on Learning Representations (ICLR)*, 2021.
- Gabriel Poesia, David Broman, Nick Haber, and Noah D. Goodman. Learning formal mathematics from intrinsic motivation, 2024. URL https://arxiv.org/abs/2407.00695.
- Stanislas Polu and Ilya Sutskever. Generative language modeling for automated theorem proving, 2020. URL https://arxiv.org/abs/2009.03393.
- Terence Tao. On the strength of theorems. https://terrytao.wordpress.com/advice-on-writing-papers/on-the-strength-of-theorems/, 2007. Accessed: 2025-09-15.
- Josef Urban and Jan Jakubův. First neural conjecturing datasets and experiments, 2020. URL https://arxiv.org/abs/2005.14664.
- Haiming Wang, Huajian Xin, Chuanyang Zheng, Lin Li, Zhengying Liu, Qingxing Cao, Yinya Huang, Jing Xiong, Han Shi, Enze Xie, Jian Yin, Zhenguo Li, Heng Liao, and Xiaodan Liang. Lego-prover: Neural theorem proving with growing libraries, 2023. URL https://arxiv.org/abs/2310.00656.

Yutong Xin, Jimmy Xin, Gabriel Poesia, Noah Goodman, Qiaochu Chen, and Isil Dillig. Automated discovery of tactic libraries for interactive theorem proving, 2025. URL https://arxiv.org/abs/2503.24036.

Kaiyu Yang, Aidan M. Swope, Alex Gu, Rahul Chalamala, Peiyang Song, Shixing Yu, Saad Godil, Ryan Prenger, and Anima Anandkumar. Leandojo: Theorem proving with retrieval-augmented language models, 2023. URL https://arxiv.org/abs/2306.15626.

Kaiyu Yang, Gabriel Poesia, Jingxuan He, Wenda Li, Kristin Lauter, Swarat Chaudhuri, and Dawn Song. Formal mathematical reasoning: A new frontier in ai, 2024. URL https://arxiv.org/abs/2412.16075.

Huaiyuan Ying, Zijian Wu, Yihan Geng, Zheng Yuan, Dahua Lin, and Kai Chen. Lean workbook: A large-scale lean problem set formalized from natural language math problems, 2025. URL https://arxiv.org/abs/2406.03847.

# A AXIOMS

We now provide all the axioms for the three domains considered in our experiments in Section 4. They are taken from the Minimo paper (Poesia et al., 2024) and formalized in the Peano languages (Poesia & Goodman, 2023).

#### Arithmetic

```
562
       = : [nat -> nat -> prop].
563
       nat : type.
564
       z : nat.
      s : [nat -> nat].
566
       o : nat.
567
568
       + : [nat -> nat -> nat].
569
       * : [nat -> nat -> nat].
570
       o_s : (= o (s z)).
571
572
       +_z : [('n : nat) -> (= (+ 'n z) 'n)].
573
       +_s: [('n : nat) -> ('m : nat) -> (= (+ 'n (s 'm)) (s (+ 'n 'm)))].
574
575
       *_z : [('n : nat) -> (= (* 'n z) z)].
       *_s : [('n : nat) -> ('m : nat) -> (= (* 'n (s 'm)) (+ 'n (* 'n 'm)))].
576
577
       nat_ind : [('p : [nat -> prop]) -> ('p z) -> [('n : nat) ->
578
                ('p 'n) \rightarrow ('p (s 'n))] \rightarrow [('n : nat) \rightarrow ('p 'n)]].
579
580
       #backward nat_ind.
       \#forward +_z ((+ 'n z) : nat).
581
       #forward +_s ((+ 'n (s 'm)) : nat).
582
       \#forward *_z ((* 'n z) : nat).
583
       \#forward *_s ((* 'n (s 'm)) : nat).
584
585
```

#### **Propositional logic**

```
587    prop : type.
588
589    false : prop.
590    /* Connectives */
591    not : [prop -> prop].
592    and : [prop -> prop -> prop].
593    or : [prop -> prop -> prop].
        iff : [prop -> prop -> prop].
```

```
594
595
       /* Introduction rule for conjunction */
596
       #backward and_i.
597
       and_i : [('P : prop) -> ('Q : prop) -> 'P -> 'Q -> (and 'P 'Q)].
       /* Elimination rules for conjunction */
598
       #forward and_el ('_ : (and 'P 'Q)).
599
       and_el : [('P : prop) \rightarrow ('Q : prop) \rightarrow (and 'P 'Q) \rightarrow 'P].
600
       \#forward and_er ('_ : (and 'P 'Q)).
601
       and_er : [('P : prop) \rightarrow ('Q : prop) \rightarrow (and 'P 'Q) \rightarrow 'Q].
602
       /* Introduction rules for disjunction */
603
       #backward or_il.
       or_il : [('P : prop) -> ('Q : prop) -> 'P -> (or 'P 'Q)].
605
       #backward or_ir.
606
       or_ir : [('P : prop) -> ('Q : prop) -> 'Q -> (or 'P 'Q)].
607
       /* Elimination rule for disjunction */
       #backward or_e infer infer infer subgoal subgoal.
608
       or_e : [('P : prop) -> ('Q : prop) -> ('R : prop) -> (or 'P 'Q) -> ['P -> 'R] -> ['Q -> 'R] -> 'R].
609
610
611
       /\star Introduction rule for negation \star/
612
       #backward not_i.
       not_i : [('P : prop) -> ['P -> false] -> (not 'P)].
613
       /* Elimination rule for negation */
614
       not_e : [('P : prop) -> (not 'P) -> 'P -> false].
615
       #backward exfalso.
616
       exfalso : [false -> ('P : prop) -> 'P].
617
       /* Introduction rules for equivalence */
618
       #backward iff_i.
619
       iff_i : [('P : prop) -> ('Q : prop) -> ['P -> 'Q] -> ['Q -> 'P] -> (iff '
620
           P'Q)].
621
       /* Elimination rules for equivalence */
622
       \#forward iff_el ('_ : (iff 'P 'Q)).
       iff_el : [('P : prop) \rightarrow ('Q : prop) \rightarrow (iff 'P 'Q) \rightarrow ['P \rightarrow 'Q]].
623
       \#forward iff_er ('_ : (iff 'P 'Q)).
624
       iff_er : [('P : prop) -> ('Q : prop) -> (iff 'P 'Q) -> ['Q -> 'P]].
625
626
       /* Excluded middle */
627
       #forward em.
       em : [('P : prop) -> (or 'P (not 'P))].
628
629
       Group theory
630
631
       = : [('t : type) -> 't -> 't -> prop].
632
       G : type.
633
634
       op : [G -> G -> G].
635
       id : G.
636
637
       /* Associativity */
       #forward op_assoc ((op (op 'a 'b) 'c) : G).
638
       op_assoc : [('a : G) -> ('b : G) -> ('c : G) ->
639
                 (= (op (op 'a 'b) 'c) (op 'a (op 'b 'c)))].
640
641
       /* Commutativity */
       #forward op_comm ((op 'a 'b) : G).
       op\_comm : [('a : G) -> ('b : G) -> (= (op 'a 'b) (op 'b 'a))].
643
644
       /* Identity */
645
       #forward id_l.
646
       id_1 : [('a : G) \rightarrow (= (op id 'a) 'a)].
647
       /* Inverse */
```

```
inv : [G -> G].
#forward inv_l.
inv_l : [('a : G) -> (= (op (inv 'a) 'a) id)].
```

# **B** EXTRINSIC EVALUATION

 In order to perform extrinsic evaluation, we run 5 iterations of our extrinisc evaluation pipeline, and take the average of the 5 results in order to mitigate variance from different runs of LLM evals. Our extrinsic evaluation pipeline consists of two steps: usefulness checking (Appendix B.1), deduplication (Appendix B.2), and SMT solving. In usefulness checking, we prompt the model concurrently on all conjectures generated by the model and keep the ones marked as useful by the LLM. As we are concurrently requesting for usefulness, we are likely to get a large amount of duplicate conjectures. We therefore make a second pass, calling the model on the useful conjectures to deduplicate them, keeping only sufficiently different theorems so as to get more reasonable results. Finally, we leverage the Z3 SMT solver (De Moura & Bjørner, 2008) to automatically prove the remaining conjectures and count only the proven ones. We found Z3 to be highly effective in proving these conjectures, as they are derived from axioms.

In the specific case of group theory, we noticed the variance in LLM evaluations was significantly higher than other domains, and the LLM had a very high rate of returning false problems. We solved this by running the SMT solver first, and giving a custom deduplication prompt (Appendix B.3) with examples for group theory.

#### **B.1** Usefulness Checking Prompt

```
You are tasked to judge whether a given lean theorem could be considered useful for an automatic theorem prover to have among its known theorems.

This theorem prover has only access to the following axioms and known theorems:

'''

{known_theorems}

As well as access to the 'rfl' and 'rewrite' commands

Here is the theorem you are to evaluate

'''lean4

{generated_conjecture}

'''

Think through the problem step by step. Translate the problem into natural language, then think of what the possible uses of the theorem could be, whether it's obviously true and whether it means something

On the last line, say either USEFUL or NOT USEFUL and nothing else.
```

#### **B.2** DEDUPLICATION PROMPT

```
I have a set of lean theorems, some of which are very similar to each
   other. I want to use them as tactics for proof generation.
Please remove the duplicates, so that I can have a list of only unique
   theorems.
For example, the following four theorems would be duplicates of each
   other:
'''lean4
theorem problem1 : (v0 : Nat) -> v0 * 1 = v0
theorem problem2 : (v0 : Nat) -> (v1 : Nat) -> v1 * 1 = v1
theorem problem3 : (v0 : Nat) -> (v1 : Nat) -> (v2 : v0 = v1) -> v1 * 1 =
   v1
theorem problem4 : (v0 : Nat) -> v0 * (Nat.succ 0) = v0
```

```
702
      The inclusion of an extra variable in problem 2 doesn't change the fact
703
          that the result is exactly the same, and the different names for the
704
          variable doesn't affect the result.
705
      Problem 3 introduces an irrelevant hypothesis, which doesn't get used in
          the theorem, and the conclusion is still the same.
706
      The last one is a trivial result of the others, as 1 is defined as Nat.
707
          succ 0 in this case.
708
      Here is my list of theorems for you to remove duplicates for.
709
710
      I also have attached an explanation for why each could be useful for a
          theorem prover.
711
      { }
712
      Think it through step by step, and then return the list of unique
713
          theorems from this list in a list format inside of a ```lean4``` code
714
           block. Make sure your answer is inside the very last lean codeblock.
           Please make sure to repeat the theorems exactly as I wrote them.
715
716
717
      B.3 GROUP THEORY SPECIFIC PROMPTS
718
719
      I have a set of lean theorems, some of which are very similar to each
720
          other. I want to use them as lemmas for proof generation.
721
      Please remove the duplicates, so that I can have a list of only unique
          theorems.
722
      For example, the following four theorems would be duplicates of each
723
          other:
724
      ```lean4
725
      theorem problem1 : ((v0 : Group) -> (v1 : (v0 = (v0 * (1^{-1}))))) ->
726
          ((1^{-1}) = 1)
      theorem problem2 : ((v0 : Group) -> (v1 : Group) -> ((1^{-1}) = 1))
727
      theorem problem3 : ((v0 : Group) -> ((1^{-1}) = 1))
728
      theorem problem4 : ((v0 : Group) -> (1 = (1^{-1}))
729
730
      Problem 1 introduces an irrelevant hypothesis as compared to problem 3,
731
          as it makes no mention of v0 in its final claim. Therefore, these two
           problems are duplicates of each other.
732
      Problem 2 is a similar case to problem 1: It introduces an extra variable
733
          , but does nothing with it. This is irrelevant, and makes for the
734
          same problem.
735
      Problem 4 is the same as problem 3, but is flipped. As we are running
736
          this using rw, we can simply call this problem in the inverse
          direction, so these two lemmas are the same.
737
738
      In this case, our final result would likely be:
739
      '''lean4
740
      theorem problem3 : ((v0 : Group) -> ((1^{-1}) = 1))
741
742
      Here is my list of theorems for you to remove duplicates for.
743
```

# B.4 Examples of Extrinsically Useful Conjectures

theorem prover.

744

745

746

747

748

749

750 751 752

753

754

755

Table 1 highlights representative conjectures that our evaluation judged to be extrinsically useful across three domains. These serve as concrete examples of the kinds of results UseFor is capable of producing. As an illustration, UseFor produces a 5-step proof of the first propositional-logic

theorems from this list in a list format inside of a ```lean4``` code

block. Make sure your answer is inside the very last lean codeblock.

Please make sure to repeat the theorems exactly as I wrote them.

I also have attached an explanation for why each could be useful for a

Think it through step by step, and then return the list of unique

conjecture in Table 1, using only the base axioms. However, given the tactic iff\_elim, which reduces an equivalence to two implications, together with the axioms

$$False \implies P,$$
 (1)

$$P \wedge Q \implies P,$$
 (2)

UseFor found the following 5-step proof:

1.Split the problem into cases: by iff\_elim  $- \text{Case 1: } False \implies P \land False$  2.introduce False into hypothesis context  $3.False \implies P \land False$  by (1)  $- \text{Case 2: } P \land False \implies False$  4.introduce  $P \land False$  into hypothesis context  $5.P \land False \implies False$  by (2)

This example demonstrates how UseFor produces lemmas that apply broadly and compress multiple reasoning steps into a single inference step in practice. This ability provides a crucial advantage in Monte Carlo Tree Search, where the search space expands exponentially with depth.

We remark that these proven conjectures are also observed to be very important to the prover in future iterations. For instance,  $P \implies \neg \neg P$  and  $1^{-1} = 1$  often serve as powerful shortcuts, condensing multi-step reasoning into a single step and thereby streamlining longer proofs.

Arithmetic	Propositional Logic	Group Theory
$\forall x \in \mathbb{N}, x(x^2 + 1) = x + x^3$	$False \iff (P \land False)$	$1^{-1} = 1$
$2x = 0 \implies x = 0$	$P \implies \neg \neg P$	
$\forall x \in \mathbb{N}, x*1 = x$	$P \iff P$	$\forall x \in G, x \cdot x = x \implies x = 1$

Table 1: Representative conjectures judged extrinsically useful across three considered domains.

### C THE USE OF LARGE LANGUAGE MODELS

Large Language Models (LLMs) were used to support writing, revision, and other text-focused tasks, such as improving clarity, refining grammar and style, and assisting with the organization of written content.