

Lean Meets Theoretical Computer Science: Scalable Synthesis of Theorem Proving Challenges in Formal-Informal Pairs

Anonymous Authors¹

Abstract

Formal theorem proving (FTP) stands at the forefront of LLM reasoning, yet existing datasets are largely limited in scope due to their dependence on manual curation. We identify theoretical computer science (TCS) as a novel testbed with promising potential, as many TCS problems can be algorithmically modeled and automatically verified, thereby enabling scalable formalization with rigorous formal-informal alignment guarantees. We showcase this approach via two TCS problem modules: the Busy Beaver (BB) challenge and the Mixed BooleanArithmetic (MBA) challenge. We evaluate leading reasoning models and finetuned theorem provers on our dataset, where the best-performing DeepSeekProver-v2-671B attains only 57.5% on our BB challenge and no more than 12% on our MBA challenge. These results reveal substantial reasoning gaps beyond conventional static benchmarks. We call on the Lean community to further study the potential of TCS problems in the formal reasoning domain where our approach enables fully automated synthesis of arbitrarily many problems in their strictly aligned formal-informal pairs that are universally **easy to verify** yet systematically **hard to prove**.

1. Introduction

While Large Reasoning Models (LRMs) (Besta et al., 2025) are evolving at an unprecedented speed, this promising progress is increasingly shadowed by a critical asymmetry: the advancing momentum of frontier models is quickly outpacing our ability to accurately gauge the boundary of their reasoning capability, as evidenced by the rapid saturation of legacy benchmarks such as MATH (Hendrycks et al.,

2021) and GPQA (Rein et al., 2023).

Beyond conventional static benchmarks that evaluate on a final answer only, a natural next stage is to examine detailed machine proof to see if models are genuinely capable of making meaningful logical derivations. Formal Theorem Proving (FTP) (Chen et al., 2025) is emerging as a defining frontier for automating such proof validation by harnessing the power of interactive theorem provers such as Lean (de Moura & Ullrich, 2021) and Isabelle (Nipkow et al., 2002).

As a result, there’s a critical need for large-scale, high-quality datasets that bridge theorem proving problems between their formal-informal forms. Previous work such as MiniF2F (Zheng et al., 2022), ProofNet (Azerbaiyev et al., 2023), PutnamBench (Tsoukalas et al., 2024), and CombiBench (Liu et al., 2025) leveraged labor-intensive expert curation to formalize problems from math competitions and college courses, yet such public data sources may suffer from inter-annotator inconsistency and data contamination of the training corpora (Magar & Schwartz, 2022; Dong et al., 2024), where even slight data leakage can inflate model performance by up to 30% points (Dekoninck et al., 2024).

In light of these concerns, Theoretical Computer Science (TCS) emerges as a promising new testbed where many problems can be **algorithmically modeled and verified**, which enables scalable fully automated problem formalization using a modularized TCS problem such as the Halting Problem of Turing Machines. Furthermore, we can adjust the difficulty of synthesized problems by tuning complexity parameters such as the state count of a Turing Machine or the variable count in an algebraic expression. This parameterized setting also enables us to synthesize fresh problems from a infinitely vast problem space, which acts as a robust mitigation against data contamination. It’s particularly worth noting that our synthesis does not rely on any LLM-generated content in any steps, which would be unavoidably impacted by language model hallucination and potential regurgitation of training data.

In light of these concerns, our main contribution in this paper is threefold:

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

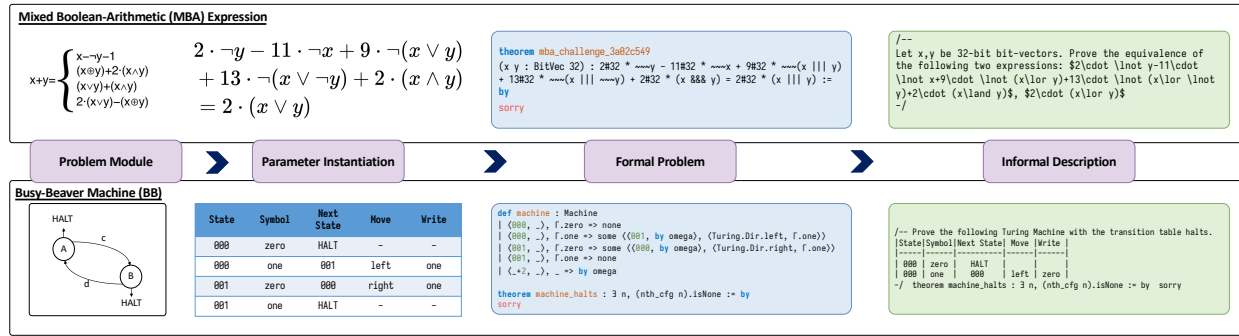


Figure 1. Complete Pipeline of our Synthesis Framework: We start with a problem module that can be instantiated by a rule-based system to generate theorem proving problem description in Lean code shown in the blue block and the same description in Markdown shown in the green blocks.

- We propose a **contamination-resistant synthesis framework with adjustable difficulty** that generates arbitrarily many theorem proving challenges in their rigorously aligned formal-informal pairs by leveraging TCS-inspired problem modules as shown in Figure 1.
- We evaluate a variety of frontier LRMs and LLM-based theorem provers. The best-performing DeepSeek-Prover-v2-671B attained best performance of 57.5% on our BB challenge and only 12% on our MBA challenge respectively, yet other theorem provers such as Kimina-Prover-Preview (Wang et al., 2025) and Goedel-Prover-SFT (Lin et al., 2025) struggle to resolve any of our synthesized problems, which may hint at how base model scaling could enhance capability of LLM-based theorem provers.
- We design a **novel step-level task decomposition setting** to gauge if models can effectively use a customized **Out-Of-Distribution** library of Lean lemma to resolve unseen problems, where best-performing model, OpenAI-o4-mini, showed 98.88% success rate. This reveals the critical bottleneck of the proof task lies in whole-proof level understanding and strategic planning as opposed to step-level problem-solving.

2. Related Work

Formal Theorem Proving (FTP). LLM has demonstrated promising potential in advancing the field of mathematics (Polu & Sutskever, 2020) with the help of interactive theorem provers as endorsed by many world-leading mathematicians such as Peter Scholze (Scholze, 2022) and Terence Tao (Tao, 2025). By integrating FTP in various contexts such as LLM-based Agents (Li et al., 2024a) and Retrieval-Augmented Generation (RAG) (Zayyad & Adi, 2024), many downstream applications can be empowered such as code copilot (Murphy et al., 2024) and operations research (Yang et al., 2024).

Autoformalization. Autoformalization refers to the process of translating natural language into formal language, which is inherently challenging for human due to the complex syntax of FTP languages. While the synthetic data approach (Wu et al., 2021; Huang et al., 2024; Ying et al., 2024) attempted to leverage LLMs for this process (Li et al., 2024b; Wu et al., 2022; Lu et al., 2024), current progress is largely limited by expert annotation to proofread and correct LLM-generated formal contents, which is labor-intensive and exposed to inter-annotator inconsistencies. Therefore, we aim to automate the entire process of (in)formalization via a template-based approach to ensure rigorous alignment, rendering manual curation *an optional quality check* as opposed to *an expensive necessity*.

Benchmark Contamination. In light of rising concerns about whether our benchmarks reflect genuine reasoning capability of LLMs (Magar & Schwartz, 2022; Dong et al., 2024), various attempts were made to detect and mitigate data contamination (Dekoninck et al., 2024). For instance, LiveBench (White et al., 2024) aims to mitigate contamination by refreshing benchmark questions every 6 months, and MixEval (Ni et al., 2024) attempted to remove preference bias by adopting a mixture of benchmarks. Nevertheless, data contamination remains a persistent issue so long as the questions are sourced from the public Internet at large, including synthetic data reliant on LLM generation.

3. Methodology

3.1. The Busy Beaver (BB) Challenge

We showcase our approach based on the halting problem of the famous Busy Beaver Machine. Busy Beaver machines are simple yet powerful computation models as any Turing-complete problems can be reduced to the Busy Beaver problem. Many research questions in modern mathematics, such

as the Collatz Conjecture (Michel, 1993) and the famous Goldbach Conjecture (Aaronson, 2020) have been proven to be reducible to BB(6) and BB(27) respectively.

The Problem Module. A busy beaver machine is an N -state, 2-symbol Turing machine, where the Busy Beaver function $\mathbf{BB}(N)$ (Harland, 2022) is defined as the maximum number of steps an N -state halting machine can execute before halt. In this challenge we instantiate with $N \leq 5$ as BB(5) is currently the limit of frontier mathematics research (Aaronson, 2020), we note that this parameter N can be any natural number in principle. By setting this parameter N we enable adjustable difficulty of synthesized problems and the number of machines available grows exponentially with the N as detailed in Table 1, thereby ensuring the infinite scalability of our synthesis.

n	Number of Machines
1	64
2	20,736
3	16,777,216
4	4,294,967,296

Table 1. Number of Busy Beaver (BB) Machines for Increasing Values of N .

Challenge definition. We ask models to prove whether a given BB machine described by a state transition table halts or not.

Generate easy-to-verify ground-truth. In order to know exactly whether a given Busy Beaver machine halts, we execute the machine for $BB(N) + 1$ steps, where $BB(N)$ is the Busy Beaver number for N -state machines, defined as the maximum steps for any N -state Turing Machine to run without halting. The numerical value of $\mathbf{BB}(N)$ ($N \leq 5$) are proven results in mathematics (Xu, 2024). BB Machines that reaches halt state within $BB(N) + 1$ steps are halting, otherwise non-halting by the definition of $BB(N)$.

Synthesize hard-to-prove challenge. We instantiate a Busy Beaver machine with a state transition table, which is then used to generate both Lean and Markdown statements using an expert-defined template. This synthetic framework ensures scalable generation of arbitrarily many proof questions in rigorously aligned formal-informal pairs as illustrated in Figure 1. Specifically, each BB machine is uniquely defined by a State Transition Table containing tuples in the form State \times Symbol \times State \times Move \times Symbol, where State $\in \{0, 1, \dots, N - 1, \text{HALT}\}$, Symbol $\in \{\text{zero}, \text{one}\}$, Move $\in \{\text{left}, \text{right}\}$, indicating a BB machine at a certain state with head pointing to certain symbol needs to change to another state, write a symbol to the

original head position and move to left or right. To generate a N state BB machine, we randomly sample the remaining elements for all $2N$ transition tuples.

With the state transition table, we apply an expert-defined template to generate the proof challenge in Lean. The template contains all necessary definitions and two blanks to fill in for the machine definition and the halting/non-halting theorem, a `sorry` at the end indicates where models should fill in their proof. We also generate a Markdown description based on a rigorously aligned template as detailed in Appendix D. We provide a complete example of our BB Challenge with problem descriptions in both Lean-Markdown and model responses in Appendix C.

3.2. The Mixed Boolean-Arithmetic (MBA) Challenge

The Problem Module Mixed Boolean-Arithmetic (MBA) expressions combine arithmetic operators ($+$, $-$, \times) with bitwise operations (\wedge (AND), \vee (OR), \oplus (XOR), \neg (NOT)) over integer variables, which form universal representations of polynomial functions over bitvectors (Reichenwallner & Meerwald-Stadler, 2022) used in cryptographic obfuscation (Liu et al., 2021).

MBA equations can be categorized into two types according to their construction:

Linear MBA equations. Equations of the form

$$\underbrace{\sum_k a_k \times e_k(x, y)}_{\text{Linear MBA expression A}} = \underbrace{\sum_k b_k \times e'_k(x, y)}_{\text{Linear MBA expression B}},$$

where each coefficient $a_k, b_k \in \{\pm 1, \pm 2, \dots, \pm 11\}$. and each atomic boolean expression $e_k(x, y), e'_k(x, y)$ have 16 options as below .

$$\left\{ \begin{array}{cccc} x & x \vee y & x \wedge y & x \oplus y \\ \neg x & x \vee \neg y & x \wedge \neg y & \neg(x \oplus y) \\ y & \neg(x \vee y) & x \wedge \neg x & \neg(x \wedge y) \\ \neg y & \neg(x \vee \neg y) & \neg(x \wedge \neg y) & \neg(x \wedge \neg x) \end{array} \right\}$$

Nonlinear MBA equations. Equations in which each side may include products of MBA sub-expressions ($e_{k_1}(x, y) \times e_{k_2}(x, y)$) and arbitrarily nested bitwise operations, e.g. $10 \times \neg(1 \times (x \wedge y) \vee \neg((-1) \times (x \vee \neg y)))$.

Dataset Generation In the first stage, we generate two-variable MBA equations following the MBA obfuscator protocol (Liu et al., 2021). Specifically, each equation is constructed by first generating a linear MBA expression that is identically zero, and then randomly moving a subset of terms to the right-hand side with their signs flipped. We then translate each equation into a Lean theorem paired with an informal Markdown description as illustrated in

Figure 1. Every variable and constant coefficient are declared as `BitVec 32` so that all arithmetic and bitwise operations are on 32-bit bitvectors modulo 2^{32} . We ensure infinite scalability by incrementing the number of variables in the MBA expression, which guarantees a sufficiently large problem space to mitigate contamination via combinatorial explosion.

Challenge definition. We ask models to prove the equivalence of two MBA expressions. Proofs must be constructed step by step using lemmas for arithmetic, bitwise, and distributive identities to ensure genuine symbolic reasoning. We explicitly discourage tactics such as `bv_decide` (which invokes a SAT solver) and other automatic approaches.

SAT solver. Lean provides a built-in tactic called `bv_decide` for automatically proving or disproving equalities over fixed-width bitvectors (`BitVec n`). Internally, `bv_decide` reduces the goal to a Boolean satisfiability problem (SAT) by applying bit-blasting, a process that translates bitvector operations into equivalent propositional logic formulas over individual bits. The resulting Boolean formula is then handed to an external SAT solver to check satisfiability. By default, `bv_decide` uses a 10-second timeout for the solver.

Generate easy-to-verify ground-truth. For any linear MBA equation

$$E_1(x, y) = E_2(x, y),$$

we normalize both sides into their *Weighted 2-DNF* (W2DNF):

$$\text{W2DNF}(E) = \sum_{i,j \in \{0,1\}} c_{ij}(E) (\ell_i(x) \wedge \ell_j(y)),$$

where $\ell_0(z) = \neg z$, $\ell_1(z) = z$ and each $c_{ij}(E) \in \mathbb{Z}_{32}$. Thus the original equality $E_1 = E_2$ holds if and only if

$$c_{ij}(E_1) = c_{ij}(E_2) \quad \text{for all } i, j \in \{0, 1\}.$$

Normalization steps are directly translated into a sequence of Lean tactics that Lean can verify as the ground-truth solution for the MBA challenge.

Customized lemma library. We curated a custom library containing all necessary lemmas to reduce any 2-variable linear MBA expression into weighted 2-DNF form without relying on any external theorem or lemma, which enabled us to automatically generate and verify all ground-truth proofs for our MBA challenge. The complete library is provided in Appendix F.

Step-level task decomposition. In order to test if models perform differently in handling step-level task comparing

to entire proof, we decompose each ground-truth proof into a sequence of step-level tasks. After each tactic is applied, Lean reports an unsolved goal. Whenever the left-hand side (LHS) or right-hand side (RHS) of the goal changes, we extract a step lemma according to the following rules:

- **For LHS changes:** The lemma’s left-hand side matches the *previous* goal’s LHS, and its right-hand side matches the *current* goal’s LHS.
- **For RHS changes:** The lemma’s left-hand side matches the *previous* goal’s RHS, and its right-hand side matches the *current* goal’s RHS.

For each problem, the final step (and *only* the final step) is always `simp`. All non-final tactics are exclusively one of three types: `simp only`, `rw`, or `nth_rewrite`, and crucially, each must use a lemma from our custom library. We provide a complete example of the step lemma challenge with model responses in Appendix C and results for this setting in 5.

4. Experimental

We report Pass@16 on various frontier models, which counts a case as successful if and only if at least one of the 16 attempts could pass Lean compiler verification, that is, no errors are reported and no ‘sorry’ is used in the proof.

While it’s possible to set other n for Pass@ n , we note that $n = 16$ is the common best practice for existing formal theorem proving benchmarks (Zheng et al., 2022; Tsoukalas et al., 2024; Liu et al., 2025), which strikes a balance of reasonable budget and non-deterministic inference. For unexpected failures such as Timeout, Out-Of-Memory and API error, we re-run each trial without counting the failed case into $n = 16$.

4.1. Evaluation Setup

Evaluated models. We evaluated both frontier general-purpose reasoning models and LLM-based theorem provers. Dedicated formal theorem provers are fine-tuned using FTP data on base models such as DeepSeekMath (Shao et al., 2024). For general-purpose LLMs, we chose a mixture of advanced open-source and proprietary models: DeepSeek-R1-671B (DeepSeek-AI et al., 2025), QwQ-32B (Qwen, 2025), OpenAI-o3 and o4-mini (OpenAI, 2025). Both OpenAI models carry default settings of 200K context window, 100K max output tokens, and a knowledge cut-off date on June 1, 2024. For LLM-based theorem provers, we evaluated various frontier models including Goedel-Prover-SFT (Lin et al., 2025), Leanabelle-Prover (Zhang et al., 2025), Kimina-Prover-Preview-7B (Wang et al., 2025), DeepSeek-Prover-v2-7B, and DeepSeek-Prover-v2-671B (Ren et al., 2025).

Model	BB(1)	BB(2)	BB(3)	BB(4)	Total
<i>General LRMs</i>					
DeepSeek-R1	21/50	12/50	10/50	12/50	55/200 (27.5%)
QwQ-32B-Preview	26/50	21/50	26/50	25/50	98/200 (49.0%)
OpenAI-o3	15/50	17/50	19/50	17/50	68/200 (34.0%)
OpenAI-o4-mini	19/50	10/50	12/50	13/50	54/200 (27.0%)
<i>Theorem Provers</i>					
Goedel-Prover-SFT	0/50	0/50	0/50	0/50	0/200 (0.0%)
Leanabell-Prover-GD-RL	0/50	0/50	0/50	0/50	0/200 (0.0%)
Kimina-Prover-Preview-7B	0/50	0/50	0/50	0/50	0/200 (0.0%)
DeepSeek-Prover-v2-7B	9/50	9/50	10/50	10/50	38/200 (19.0%)
DeepSeek-Prover-v2-671B	34/50	30/50	26/50	25/50	115/200 (57.5%)

Table 2. Pass@16 Results on BB Challenge, where each column represents an incremental setting for state count N as a complexity parameter.

Inference setting. For open-source models, we use greedy sampling to aid reproducibility, which is equivalent to fixing the sampling temperature to $T = 0.0$ as recommended by the model developers (DeepSeek-AI et al., 2025; Qwen, 2025) for math and coding tasks. Furthermore, We enable thinking mode for all models whenever possible and follow default settings on Reasoning Efforts, Context Length, Max Tokens, and System Prompts in model metadata.

We implemented our code under a Linux environment with Python 3.10. All open-source models are deployed using `transformer` library and `PyTorch` on a 4xGH200 server, with open-source models downloaded via `huggingface-cli`. The proprietary models are used directly via the OpenAI API. We prompted models to follow detailed instruction as in Appendix C.

After generation, model outputs are parsed to extract Lean proof, which is validated using a regular expression. The proof is then verified using `Kimina-Lean-Server` (Santos et al., 2025) built with `Lean` and `Mathlib` v4.19.0 via a RESTful API.

4.2. Results

BB challenge. We report Pass@16 in Table 2 and analyze detailed failure modes in Section 4.4. Frontier LRMs OpenAI-o3 and o4-mini have struggled to score past 35% on our BB challenge, marking a critical gap in formal reasoning capability despite rapid saturation on static benchmarks from public sources.

Among LLM-based theorem provers, best-performing DeepSeek-Prover-v2-671B (Ren et al., 2025) scored 57.5%, but we notice rapid performance degradation as models scale down to 7B, where advanced theorem provers such as Goedel-Prover-SFT (Lin et al., 2025), Leanabelle-Prover-

GD-RL (Zhang et al., 2025) and Kimina-Prover-Preview-7B (Wang et al., 2025) struggled to solve even 1 instance of our challenges.

Paradoxically, QwQ-32B (Qwen, 2025) attained best performance among general LRMs with only 32B parameters, which raises an intriguing question as to how scaling affects the performance of general reasoning models vs. dedicated theorem provers in a potentially different manner.

MBA challenge. We report Pass@16 in Table 5, revealing even more pronounced difficulties than BB challenge. Further, we decompose 10 randomly selected MBA challenge into 365 steps according to Section 3.2 and report step-level performance in Table 4.

Frontier LRMs including DeepSeek-R1 (DeepSeek-AI et al., 2025) and QwQ-32B (Qwen, 2025), struggled to solve any of our challenges, while OpenAI-o4-mini (OpenAI, 2025) solved only 6 problems.

Specialized theorem provers exhibited extreme stratification: DeepSeek-Prover-v2-671B (Ren et al., 2025) solved 12% of linear problems with lemmas, but failed all problems without lemmas and all non-linear problems. Smaller provers showed near-universal failure, with only Kimina-Prover-7B (Wang et al., 2025) solving a single instance. This demonstrates that even state-of-the-art formal systems struggle with fundamental algebraic reasoning when deprived of lemma support or faced with non-linear complexity.

4.3. Success Cases

We present a successful case of a proof of non-halting BB(4) machine and more cases in Appendix E. The model DeepSeekProver-v2-671B leveraged proof by induction and attempted to use various automated tactics, such as `aesop`

Model	Linear w/ Lemma	Linear w/o Lemma	Non-Linear
<i>Reasoning LLMs</i>			
DeepSeek-R1	0/100	0/100	0/40
OpenAI-o4-mini	2/100	4/100	0/40
QwQ-32B-Preview	0/100	0/100	0/40
<i>Theorem Provers</i>			
DeepSeek-Prover-v2-7B	0/100	0/100	0/40
Goedel-Prover	0/100	0/100	0/40
Kimina-Prover-7B	1/100	0/100	0/40
Leanabelle-Prover	0/100	0/100	0/40
DeepSeek-Prover-v2-671B	12/100	0/100	0/40

Table 3. Pass@16 Performance on MBA challenge. We evaluated all models using 2 parallel settings where models are tested with vs. without our customized lemma library as detailed in Section 3.2. This parallel setting aims to gauge their capability of in-context learning unseen lemma to tackle novel problems out of their training corpora.

Model	#Correct Steps
OpenAI-o4-mini	352/356
DeepSeek-Prover-v2-671B	258/356

Table 4. Step-level task performance on linear MBA where models are asked to choose an appropriate lemma from our customized lemma library for each step of the proof.

and omega, showing that a considerable manipulation of the automated tactic taxonomy.

```

DeepSeek-Prover-v2-671B:
theorem machine_never_halts : ∀ n, (nth_cfg n).isSome
  => by
  intro n
  have h : ∀ n, (nth_cfg n).isSome := by
    intro n
    induction n with
    | zero => simp [nth_cfg, init]
    | succ n ih =>
      cases h1 : nth_cfg n with
      | none => simp_all [nth_cfg, step, machine, Cfg,
        ↪ Stmt, Γ, Fin.forall_fin_succ, Fin.
        ↪ forall_fin_succ']
      | some cfg =>
        simp_all [nth_cfg, step, machine, Cfg, Stmt, Γ,
        ↪ Fin.forall_fin_succ, Fin.forall_fin_succ
        ↪']
  <|>
  (try cases cfg <|> simp_all [nth_cfg, step,
    ↪ machine, Cfg, Stmt, Γ, Fin.
    ↪ forall_fin_succ, Fin.forall_fin_succ'])
  <|>
  (try aesop)
  <|>
  (try omega)
  <|>
  (try aesop)
exact h n
    
```

In this successful MBA case, the problem is simple enough that bv_decide is able to solve it within its default 10-

second timeout by leveraging a SAT solver.

```

o4-mini:
theorem mba_challenge_88282d89 (x y : BitVec 32) : 1#32
  ↪ * ~~~(x &&& ~~~y) + 1#32 * ~~~(x ^^^ y) - 3#32 *
  ↪ (x ||| ~~~y) + 1#32 * ~~~(x ||| y) + 3#32 * (x
  ↪ &&& ~~~y) + 1#32 * (x &&& y) = 1#32 * ~~~(x |||
  ↪ ~~~y) := by
  bv_decide
    
```

This case is more complex such that bv_decide alone cannot solve. However, the combination of norm_num and bv_decide is sufficient to complete the proof. Notably, the sequence of have statements introducing previously proven lemmas is entirely unused, and the tactics following the first bv_decide are never executed, as the goal is already solved at that point. This suggests that the model remains heavily reliant on automated tactics and lacks a clear understanding of what these tactics actually do.

```

DeepSeek-Prover-v2-671B:
theorem mba_challenge_3165c251 (x y : BitVec 32) :
  -3#32 * (x ||| ~~~y) - 1#32 * x - 2#32 * y + 1#32 *
  ↪ ~~~(x ||| y) + 2#32 * ~~~(x ||| ~~~y) + 2#32 *
  ↪ (x &&& ~~~y) + 6#32 * (x &&& y) = -2#32 * ~~~y
  := by
  have h1 := bv32_not_xor_eq_or x y
  ...
  have h21 := bv32_sum_all y x
  norm_num [Nat.add_assoc, Nat.add_comm, Nat.
    ↪ add_left_comm, Nat.mul_assoc, Nat.mul_comm, Nat.
    ↪ mul_left_comm] at *
  <|> simp_all (config := {decide := true})
  <|> bv_decide
  <|> simp_all (config := {decide := true})
  <|> bv_decide
    
```

4.4. Failure Cases

Our experiments reveal several key failure modes of frontier models. These failure can be categorized into 4 major types with signature patterns as shown in Table 5. More cases in

Appendix E.

Error Message Type	%(#Errors)
Irrelevant Hallucination	67.27% (122404)
Tactic Misuse	23.22% (42247)
Voluntary Give Up	4.88% (8871)
Type Mismatch	4.47% (8136)
Others	0.16% (300)

Table 5. Summary of error type distribution for both BB and MBA Challenge.

Hallucination. Hallucination consistently appears within 7B theorem provers, where they often tried to use non-existent theorems or tactics. DeepSeek-Prover-v2-671B has exhibited much stronger anti-hallucination capability thanks to its large base model.

Here we show 2 cases where Goedel-Prover-SFT attempted to prove by induction on a completely unrelated problem. Similarly, Kimina-Prover-Preview tried to rename the given theorem head `nth_cfg` to `nth_configuration`, which led to an unsolved goal at the end of the proof.

```
Goedel-Prover-SFT:
theorem sequence_formula {a : ℕ → ℕ} (h1 : a 1 = 1)
  (h_rec : ∀ n : ℕ, n ≥ 1 → a (n + 1) = 2 * a n + 1)
  ↪ :
  ∀ n : ℕ, n ≥ 1 → a n = 2^n - 1 := by
```

```
Kimina-Prover-Preview:
theorem machine_never_halts : ∀ n, (nth cfigurations n)
  ↪.isSome := by
  intro n
  induction n with
  | zero =>
  | simp [nth configurations]
  | succ n ih =>
  | simp [nth configurations, ih]
  </> aesop
```

For the MBA challenge, we show another failure case of OpenAI-o4-mini where model hallucinates a non-existent tactic `bv_simp`.

```
o4-mini:
theorem mba_challenge_14b066fa (x y : BitVec 32) :
  ↪ -1#32 * ~~~(x ^^^ y) + 2#32 * ~~~x + 2#32 * ~~~(x
  ↪ &&& y) + 1#32 * (x ||| y) - 3#32 * x + 3#32 * (x
  ↪ ||| ~~~y) - 2#32 * (x ^^^ y) + 1#32 * ~~~(x &&&
  ↪ ~~~x) - 5#32 * y - 8#32 * ~~~(x ||| y) - 2#32 * (
  ↪ x &&& ~~~y) + 3#32 * (x &&& y) = -1#32 * ~~~(x
  ↪ &&& ~~~y) := by
  bv_simp
  ring
```

Tactic misuse. Models showed a dangerous tendency to blindly apply automated proof tactics such as `aesop` (which automatically looks for search tactics recursively in a proof

tree) without understanding the necessary conditions for `aesop` to be applicable, resulting in the failure of exhaustive search.

```
DeepSeek-Prover-v2-671B:
theorem machine_never_halts : ∀ n, (nth configurations
  ↪ n).isSome := by
  intro n
  ...
  </> aesop
```

In this MBA challenge example, this tendency is even more pronounced. The model gathers many lemmas into one `simp_only` call without checking their relevance and then blindly invokes `ring`. The proof is bloated with useless rewrites and heavy automation, yet still fails demonstrating a failure to reason about lemma and tactic applicability.

```
o4-mini with lemmas:
theorem mba_challenge_080db093 (x y : BitVec 32) :
  -5#32 * ~~~y + 1#32 * ~~~(x &&& ~~~x) + 4#32 * (x |||
  ↪ y) - 6#32 * (x ^^^ y) - 6#32 * ~~~(x &&& ~~~y)
  ↪ - 1#32 * ~~~(x &&& y) - 3#32 * y - 11#32 * ~~~
  ↪ x + 24#32 * ~~~(x ||| y) + 22#32 * ~~~(x |||
  ↪ ~~~y) + 13#32 * (x &&& ~~~y) + 6#32 * (x &&& y
  ↪) = -4#32 * ~~~(x ^^^ y) + 6#32 * (x ||| ~~~y)
  ↪ := by
  simp only [bv32_and_not_self, bv32_not_and_self,
  ↪ bv32_or_not_self, bv32_not_or_self] at *
  simp only [bv32_x_distr, bv32_y_distr, ...] at *
  simp only [bv32_sum_all] at *
  ring
```

This sample illustrates a tactic misuse where `bv_decide` fails to solve the goal within the 10-second timeout. The failure highlights the limitations of relying solely on automated tactics for large or intricate goals.

```
o4-mini:
import Std.Tactic.BVDecide
/--
Let x,y be 32-bit bit-vectors. Prove the equivalence of
  ↪ the following two expressions: 7·x - 5·y - 2·(x⊕
  ↪ y) - 6·¬(x ∧ ¬x) - 5·(x∨y) - 2·¬(x∧y) - (x∨¬y) +
  ↪ 4·¬y - 7·¬(x⊕y) + 13·¬(x∨y) + 21·¬(x∨¬y) + 17·(x∧
  ↪ y) and -5·(x∧¬y) + 1·¬x
-/
theorem mba_challenge_6f99807f (x y : BitVec 32) :
  ↪ ((7#32 * x) - (5#32 * y) - (2#32 * (x ^^^ y)) -
  ↪ (6#32 * (~~~(x &&& (~~~x)))) - (5#32 * (x ||| y))
  ↪ - (2#32 * (~~~(x &&& y))) - (1#32 * (x ||| (~~~y
  ↪))) + (4#32 * (~~~y)) - (7#32 * (~~~(x ^^^ y))) +
  ↪ (13#32 * (~~~(x ||| y))) + (21#32 * (~~~(x |||
  ↪ (~~~y)))) + (17#32 * (x &&& y))) = ((-5#32) * (x
  ↪ &&& (~~~y))) + (1#32 * (~~~x)) := by
  bv_decide
```

Type mismatch. Type mismatch is a fundamental syntactic error arising from the type-dependent nature of Lean system, which directly reflects models' failure to understand and correctly apply respective Lean tactics in question. In this case, the 7B theorem-proving models failed to interpret and adapt to Leans inherent type rules, producing expressions whose inferred types did not match the expected ones.

Fortunately, this type of error occurs rarely at 4.47% in the big picture, which speaks to the promising progress that frontier models have a considerably correct understanding of the Lean tactic taxonomy.

```
Kimina-Prover-Preview:
theorem machine_never_halts : ∀ n, (nth_cfg n).isSome
  ↪ := by
  intro n
  induction n with
  ...
  | succ n ih =>
    simp [nth_cfg, step, machine]
  exact ih
```

Voluntary give-up. In this case, the model voluntarily chose to give up by leaving `sorry` under the proof goal without even attempting to solve the problem with any tactics at all. This type of error makes up 4.88% of total errors, which could be seen as a failure of instruction-following often depicted by reasoning models comparing to general-purpose LLMs.

```
Kimina-Prover-Preview:
theorem machine_never_halts : ∀ n, (nth_cfg n).isSome
  ↪ := by
  sorry
```

5. Discussion

The Automation Trap in Theorem Proving Our empirical analysis identifies a systematic overreliance on automated tactics as primary proof mechanisms across frontier models. Models predominantly deploy tactics such as `aesop` for recursive proof search and `bv_decide` for SAT-based solving as black-box oracles, demonstrating minimal comprehension of their formal operational boundaries. This dependency manifests in two critical failure modalities: (1) proofs that circumvent genuine deductive understanding when tactics succeed coincidentally, and (2) complete proof collapse when tactics fail without fallback mechanisms. Notably, models frequently generate elaborate lemma declarations while defaulting to brute-force tactic sequencing a pattern indicating a fundamental disconnect between syntactic manipulation and strategic reasoning that fundamentally undermines verification integrity.

The Reasoning Gap between Step-Level vs. Holistic Proof Task Our step-level task decomposition setting reveals a significant discrepancy between step-level performance and holistic proof synthesis with a simple overhead proof goal. Models demonstrate near-perfect capability in choosing unseen out-of-distribution lemma for atomic inference tasks yet exhibit catastrophic failure rates in composing these operations into complete proofs. This divergence indicates that the primary bottleneck resides not in local operations but in global proof planning and strategic integration.

The persistent inability to reconcile stepwise correctness with end-to-end whole proof generation suggests reasoning capabilities of frontier models remain constrained in long contexts by insufficient strategic planning for orchestrating atomic transformations into coherent proof strategies when challenged with lengthy proofs containing multiple steps.

Contamination-Resistant Evaluation Framework Our methodology establishes a novel evaluation paradigm through the integration of Lean with Theoretical Computer Science. Our framework achieves infinite scalability via algorithmic problem generation from parameterized TCS modules, with granular difficulty modulation through computational parameters (e.g., Turing Machine state complexity, MBA expression depth). Automated verifiability provides ground-truth validation without human intervention, while dynamically generated problem spaces ensure intrinsic resistance to dataset contamination. By enabling rigorous formal-informal alignment absent expert curation, this TCS-inspired synthesis creates a sustainable evaluation ecosystem where benchmark freshness and complexity scales along with the progress of frontier theoretical computer science research.

6. Conclusion

We propose Theoretical Computer Science as a promising testbed for formal theorem proving, which by design is capable of utilizing modularized, fully automatic problem generation to systematically overcome contamination and scalability limitations in legacy benchmarks. Our TCS-inspired modular synthesis paradigm enables the creation of infinitely scalable, formally rigorous problem spaces with adjustable complexity through computational parameters such as Turing Machine state configurations and MBA expression depth. Our experimental results demonstrate a profound reasoning gap: while models achieve near-perfect scores of 98.88% on atomic step-level tasks involving out-of-distribution lemma selection, they collapse to a mere 12% success rate when synthesizing complete lengthy proofs. This drastic performance degradation highlights the devastating role of hallucination during long-context reasoning sequences, which may act as a roadblock for models to form any systematic strategies for tasks requiring lengthy whole proofs. We also reveal the dangerous tendency of frontier models' over-reliance on automated tactics like `aesop` and `bv_decide` without understanding their applicability with limitations. To tackle these challenges, we call on the Lean community to further explore the massive potential of theoretical computer science and its interplay with the realm of formal theorem proving, thereby leveraging the joint effort to advance both domains synergistically.

References

- Aaronson, S. The busy beaver frontier. *SIGACT News*, 51(3):32–54, 2020. doi: 10.1145/3427361.3427369. URL <https://doi.org/10.1145/3427361.3427369>.
- Anthropic. Claude 3.7 sonnet system card. <https://assets.anthropic.com/m/785e231869ea8b3b/original/claude-3-7-sonnet-system-card.pdf>, 2025. Accessed: 2025-05-10.
- Azerbaiyev, Z., Piotrowski, B., Schoelkopf, H., Ayers, E. W., Radev, D., and Avigad, J. Proofnet: Autoformalizing and formally proving undergraduate-level mathematics. *CoRR*, abs/2302.12433, 2023. doi: 10.48550/ARXIV.2302.12433. URL <https://doi.org/10.48550/arXiv.2302.12433>.
- Besta, M., Barth, J., Schreiber, E., Kubicek, A., Catarino, A., Gerstenberger, R., Nyczyk, P., Iff, P., Li, Y., Houliston, S., et al. Reasoning language models: A blueprint. *arXiv preprint arXiv:2501.11223*, 2025.
- Chen, S., Yu, W., Dou, G., and Zhang, Q. A review on mechanical proving and formalization of mathematical theorems. *IEEE Access*, 13:50672–50686, 2025. doi: 10.1109/ACCESS.2025.3552634. URL <https://doi.org/10.1109/ACCESS.2025.3552634>.
- de Moura, L. and Ullrich, S. The lean 4 theorem prover and programming language. In Platzer, A. and Sutcliffe, G. (eds.), *Automated Deduction - CADE 28 - 28th International Conference on Automated Deduction, Virtual Event, July 12-15, 2021, Proceedings*, volume 12699 of *Lecture Notes in Computer Science*, pp. 625–635. Springer, 2021. doi: 10.1007/978-3-030-79876-5_37. URL https://doi.org/10.1007/978-3-030-79876-5_37.
- DeepSeek-AI, Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., Zhang, X., Yu, X., Wu, Y., Wu, Z. F., Gou, Z., Shao, Z., Li, Z., Gao, Z., Liu, A., Xue, B., Wang, B., Wu, B., Feng, B., Lu, C., Zhao, C., Deng, C., Zhang, C., Ruan, C., Dai, D., Chen, D., Ji, D., Li, E., Lin, F., Dai, F., Luo, F., Hao, G., Chen, G., Li, G., Zhang, H., Bao, H., Xu, H., Wang, H., Ding, H., Xin, H., Gao, H., Qu, H., Li, H., Guo, J., Li, J., Wang, J., Chen, J., Yuan, J., Qiu, J., Li, J., Cai, J. L., Ni, J., Liang, J., Chen, J., Dong, K., Hu, K., Gao, K., Guan, K., Huang, K., Yu, K., Wang, L., Zhang, L., Zhao, L., Wang, L., Zhang, L., Xu, L., Xia, L., Zhang, M., Zhang, M., Tang, M., Li, M., Wang, M., Li, M., Tian, N., Huang, P., Zhang, P., Wang, Q., Chen, Q., Du, Q., Ge, R., Zhang, R., Pan, R., Wang, R., Chen, R. J., Jin, R. L., Chen, R., Lu, S., Zhou, S., Chen, S., Ye, S., Wang, S., Yu, S., Zhou, S., Pan, S., and Li, S. S. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *CoRR*, abs/2501.12948, 2025. doi: 10.48550/ARXIV.2501.12948. URL <https://doi.org/10.48550/arXiv.2501.12948>.
- Dekoninck, J., Müller, M. N., and Vechev, M. T. Constat: Performance-based contamination detection in large language models. In Globersons, A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J. M., and Zhang, C. (eds.), *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. URL http://papers.nips.cc/paper_files/paper/2024/hash/a7f89793b9e6f8c6568dbbb66ff727b9b-Abstract-Conference.html.
- Dong, Y., Jiang, X., Liu, H., Jin, Z., Gu, B., Yang, M., and Li, G. Generalization or memorization: Data contamination and trustworthy evaluation for large language models. In Ku, L., Martins, A., and Srikumar, V. (eds.), *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pp. 12039–12050. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.FINDINGS-ACL.716. URL <https://doi.org/10.18653/v1/2024.findings-acl.716>.
- Google DeepMind. Gemini 2.5: Our most intelligent ai model. Google AI Blog, Mar 2025; <https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/>, 2025.
- Harland, J. Generating candidate busy beaver machines (or how to build the zany zoo). *Theor. Comput. Sci.*, 922:368–394, 2022. doi: 10.1016/J.TCS.2022.04.040. URL <https://doi.org/10.1016/j.tcs.2022.04.040>.
- Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., and Steinhardt, J. Measuring mathematical problem solving with the MATH dataset. In Vanschoren, J. and Yeung, S. (eds.), *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*, 2021. URL <https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/be83ab3ecd0db773eb2dc1b0a17836a1-Abstract-round2.html>.
- Huang, Y., Lin, X., Liu, Z., Cao, Q., Xin, H., Wang, H., Li, Z., Song, L., and Liang, X. MUSTARD: mastering uni-

- form synthesis of theorem and proof data. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=8xliOUg9EW>.
- Li, Z., Hua, W., Wang, H., Zhu, H., and Zhang, Y. Formal-llm: Integrating formal language and natural language for controllable llm-based agents, 2024a. URL <https://arxiv.org/abs/2402.00798>.
- Li, Z., Wu, Y., Li, Z., Wei, X., Zhang, X., Yang, F., and Ma, X. Autoformalize mathematical statements by symbolic equivalence and semantic consistency. In Globersons, A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J. M., and Zhang, C. (eds.), *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024b. URL http://papers.nips.cc/paper_files/paper/2024/hash/6034a661584af6c28fd97a6f23e56c0a-Abstract-Conference.html.
- Lin, Y., Tang, S., Lyu, B., Wu, J., Lin, H., Yang, K., Li, J., Xia, M., Chen, D., Arora, S., and Jin, C. Goedel-prover: A frontier model for open-source automated theorem proving, 2025. URL <https://arxiv.org/abs/2502.07640>.
- Liu, B., Feng, W., Zheng, Q., Li, J., and Xu, D. Software obfuscation with non-linear mixed boolean-arithmetic expressions. In Gao, D., Li, Q., Guan, X., and Liao, X. (eds.), *Information and Communications Security - 23rd International Conference, ICICS 2021, Chongqing, China, November 19-21, 2021, Proceedings, Part I*, volume 12918 of *Lecture Notes in Computer Science*, pp. 276–292. Springer, 2021. doi: 10.1007/978-3-030-86890-1_16. URL https://doi.org/10.1007/978-3-030-86890-1_16.
- Liu, J., Lin, X., Bayer, J., Dillies, Y., Jiang, W., Liang, X., Soletskyi, R., Wang, H., Xie, Y., Xiong, B., Yang, Z., Zhang, J., Zhi, L., Li, J., and Liu, Z. Combibench: Benchmarking llm capability for combinatorial mathematics, 2025. URL <https://arxiv.org/abs/2505.03171>.
- Lu, J., Liu, Z., Wan, Y., Huang, Y., Wang, H., Yang, Z., Tang, J., and Guo, Z. Process-driven autoformalization in lean 4. *CoRR*, abs/2406.01940, 2024. doi: 10.48550/ARXIV.2406.01940. URL <https://doi.org/10.48550/arXiv.2406.01940>.
- Magar, I. and Schwartz, R. Data contamination: From memorization to exploitation. In Muresan, S., Nakov, P., and Villavicencio, A. (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pp. 157–165. Association for Computational Linguistics, 2022. doi: 10.18653/V1/2022.ACL-SHORT.18. URL <https://doi.org/10.18653/v1/2022.acl-short.18>.
- Michel, P. Busy beaver competition and collatz-like problems. *Arch. Math. Log.*, 32(5):351–367, 1993. doi: 10.1007/BF01409968. URL <https://doi.org/10.1007/BF01409968>.
- Murphy, W., Holzer, N., Qiao, F., Cui, L., Rothkopf, R., Koenig, N., and Santolucito, M. Combining LLM code generation with formal specifications and reactive program synthesis. *CoRR*, abs/2410.19736, 2024. doi: 10.48550/ARXIV.2410.19736. URL <https://doi.org/10.48550/arXiv.2410.19736>.
- Ni, J., Xue, F., Yue, X., Deng, Y., Shah, M., Jain, K., Neubig, G., and You, Y. Mixeval: Deriving wisdom of the crowd from LLM benchmark mixtures. In Globersons, A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J. M., and Zhang, C. (eds.), *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. URL http://papers.nips.cc/paper_files/paper/2024/hash/b1f34d7b4a03a3d80be8e72eb430dd81-Abstract-Conference.html.
- Nipkow, T., Paulson, L. C., and Wenzel, M. *Isabelle/HOL - A Proof Assistant for Higher-Order Logic*, volume 2283 of *Lecture Notes in Computer Science*. Springer, 2002. ISBN 3-540-43376-7. doi: 10.1007/3-540-45949-9. URL <https://doi.org/10.1007/3-540-45949-9>.
- OpenAI. Openai o3 and o4-mini system card. <https://cdn.openai.com/pdf/2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf>, 2025.
- Polu, S. and Sutskever, I. Generative language modeling for automated theorem proving. *CoRR*, abs/2009.03393, 2020. URL <https://arxiv.org/abs/2009.03393>.
- Qwen. Qwq - 32b: The power of scaling rl, March 2025. URL <https://qwenlm.github.io/blog/qwq-32b/>. Accessed: 2025 - 05 - 14.
- Reichenwallner, B. and Meerwald-Stadler, P. Efficient de-obfuscation of linear mixed boolean-arithmetic expressions. In III, G. G. R. and Blazytko, T. (eds.), *Proceedings of the 2022 ACM Workshop on Research on*

- offensive and defensive techniques in the context of Man At The End (MATE) attacks, Los Angeles, CA, USA, 11 November 2022, pp. 19–28. ACM, 2022. doi: 10.1145/3560831.3564256. URL <https://doi.org/10.1145/3560831.3564256>.
- Rein, D., Hou, B. L., Stickland, A. C., Petty, J., Pang, R. Y., Dirani, J., Michael, J., and Bowman, S. R. GPQA: A graduate-level google-proof q&a benchmark. *CoRR*, abs/2311.12022, 2023. doi: 10.48550/ARXIV.2311.12022. URL <https://doi.org/10.48550/arXiv.2311.12022>.
- Ren, Z. Z., Shao, Z., Song, J., Xin, H., Wang, H., Zhao, W., Zhang, L., Fu, Z., Zhu, Q., Yang, D., Wu, Z. F., Gou, Z., Ma, S., Tang, H., Liu, Y., Gao, W., Guo, D., and Ruan, C. Deepseek-prover-v2: Advancing formal mathematical reasoning via reinforcement learning for subgoal decomposition, 2025. URL <https://arxiv.org/abs/2504.21801>.
- Santos, M. D., Wang, H., de Saxcé, H., Wang, R., Baksys, M., Unsal, M., Liu, J., Liu, Z., and Li, J. Kimina lean server: Technical report. *arXiv preprint arXiv:2504.21230*, 2025.
- Scholze, P. Liquid tensor experiment. *Exp. Math.*, 31(2):349–354, 2022. doi: 10.1080/10586458.2021.1926016. URL <https://doi.org/10.1080/10586458.2021.1926016>.
- Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang, H., Zhang, M., Li, Y. K., Wu, Y., and Guo, D. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL <https://arxiv.org/abs/2402.03300>.
- Tao, T. ams.org. <https://www.ams.org/notices/202501/rnoti-p6.pdf>, 2025. [Accessed 10-05-2025].
- 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. In Globersons, A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J. M., and Zhang, C. (eds.), *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. URL http://papers.nips.cc/paper_files/paper/2024/hash/1582eaf9e0cf349e1e5a6ee453100aa1-Abstract-Datasets_and_Benchmarks_Track.html.
- Wang, H., Unsal, M., Lin, X., Baksys, M., Liu, J., Santos, M. D., Sung, F., Vinyes, M., Ying, Z., Zhu, Z., Lu, J., de Saxcé, H., Bailey, B., Song, C., Xiao, C., Zhang, D., Zhang, E., Pu, F., Zhu, H., Liu, J., Bayer, J., Michel, J., Yu, L., Dreyfus-Schmidt, L., Tunstall, L., Pagani, L., Machado, M., Bourigault, P., Wang, R., Polu, S., Barroyer, T., Li, W.-D., Niu, Y., Fleureau, Y., Hu, Y., Yu, Z., Wang, Z., Yang, Z., Liu, Z., and Li, J. Kimina-prover preview: Towards large formal reasoning models with reinforcement learning, 2025. URL <https://arxiv.org/abs/2504.11354>.
- White, C., Dooley, S., Roberts, M., Pal, A., Feuer, B., Jain, S., Schwartz-Ziv, R., Jain, N., Saifullah, K., Naidu, S., Hegde, C., LeCun, Y., Goldstein, T., Neiswanger, W., and Goldblum, M. Livebench: A challenging, contamination-free LLM benchmark. *CoRR*, abs/2406.19314, 2024. doi: 10.48550/ARXIV.2406.19314. URL <https://doi.org/10.48550/arXiv.2406.19314>.
- Wu, Y., Jiang, A., Ba, J., and Grosse, R. B. {INT}: An inequality benchmark for evaluating generalization in theorem proving. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=O6LPudowNQm>.
- Wu, Y., Jiang, A. Q., Li, W., Rabe, M. N., Staats, C., Jamnik, M., and Szegedy, C. Autoformalization with large language models. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/d0c6bc641a56bebee9d985b937307367-Abstract-Conference.html.
- Xu, C. Skelet 17 and the fifth busy beaver number, 2024. URL <https://arxiv.org/abs/2407.02426>.
- Yang, Z., Wang, Y., Huang, Y., Guo, Z., Shi, W., Han, X., Feng, L., Song, L., Liang, X., and Tang, J. Optibench meets resocratic: Measure and improve llms for optimization modeling, 2024. URL <https://arxiv.org/abs/2407.09887>.
- 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. In Globersons, A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J. M., and Zhang, C. (eds.), *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. URL http://papers.nips.cc/paper_files/paper/2024/hash/1582eaf9e0cf349e1e5a6ee453100aa1-Abstract-Datasets_and_Benchmarks_Track.html.

nips.cc/paper_files/paper/2024/hash/
bf236666a2cc5f3ae05d2e08485efc4c-
Abstract-Datasets_and_Benchmarks_
Track.html.

Zayyad, M. and Adi, Y. Formal language knowledge corpus
for retrieval augmented generation, 2024. URL <https://arxiv.org/abs/2412.16689>.

Zhang, J., Wang, Q., Ji, X., Liu, Y., Yue, Y., Zhang, F.,
Zhang, D., Zhou, G., and Gai, K. Leanabell-prover:
Posttraining scaling in formal reasoning, 2025. URL
<https://arxiv.org/abs/2504.06122>.

Zheng, K., Han, J. M., and Polu, S. minif2f: a cross-
system benchmark for formal olympiad-level mathemat-
ics. In *The Tenth International Conference on Learn-
ing Representations, ICLR 2022, Virtual Event, April
25-29, 2022*. OpenReview.net, 2022. URL <https://openreview.net/forum?id=9ZPegFuFTFv>.

A. Limitations

While this work has been conducted according to best practices of previous work on formal theorem proving, there are several potential limitations:

Firstly, our evaluation for proprietary OpenAI models are conducted with their default configuration, combined with our zero-temperature setting for open-sourced models to enhance reproducibility, all evaluated models may not have shared completely identical settings due to the undisclosed parameter settings of proprietary models.

Meanwhile, in order to strike a balance between budgetary restraints and performance representativeness, we only selected best-performing models on mainstream formal theorem proving benchmarks such as PutnamBench (Tsoukalas et al., 2024) and miniF2F (Zheng et al., 2022) for evaluation. There may be merit in exploring how other frontier proprietary models such as Gemini-2.5 (Google DeepMind, 2025) and Claude-3.7-Sonnet (Anthropic, 2025) would perform on our dataset. Therefore, we plan to open-source our codebase and dataset to the community for researchers with ample resources to explore more models.

Lastly, while we leverage Lean compiler to automate proof validation, the proof is graded on a pass/fail basis as a whole. However, we note it’s possible to implement a stepwise grading metric by dynamically masking each step in a whole proof and ask the model to fill in each step respectively, thereby gauging models’ capability to implement key intermediate steps.

B. Automated Proof Validation

In our validation pipeline, we adopt the open-source Lean 4 server implementation from Kimina Lean Server (Santos et al., 2025). The Kimina Lean Server provides a Python interface allowing real-time feedback for generated proofs. With a Python function call `verify`, the server receives a list of Lean proofs and returns validation results for each proof. The server itself handles multiple requests efficiently by spreading verification across multiple Lean REPL processes.

C. Complete Example with Our Prompting Strategy

```

'''lean4
{lean_code}
'''
You can make your own auxiliary corollaries and theorems to support the proof, instead of only completing the part
→with the sorry. Please output the entire program and not just the last part. Please output only the program
→and add no other comment, such that your answer is a compilable lean code. Make sure to reason enough to make
→ your code correct.
'''

```

System Prompt We used default system prompts as in the evaluated models’ metadata.

A complete BB example. Below we demonstrate a complete record of evaluation, including the problem statement in Lean 4 and Markdown, our evaluation prompt and the model response:

```

import Mathlib.Computability.TuringMachine

inductive Γ
| zero
| one
deriving DecidableEq

instance : Inhabited Γ := ⟨ Γ.zero ⟩

structure Stmt where
  move : Turing.Dir
  write : Γ

def Machine := Fin 2 → Γ → Option (Fin 2 Stmt)

structure Cfg where
  q : Fin 2
  Tape : Turing.Tape Γ

```

```

715 def init (l : List Γ) : Cfg := ⟨⟨0, by omega⟩, Turing.Tape.mk 1⟩
716
717 def step (M : Machine) : Cfg → Option Cfg :=
718   fun ⟨q, T⟩      (M q T.head).map fun ⟨q', a⟩      ⟨q', (T.write a.write).move a.move⟩
719
720 def machine : Machine
721 | ⟨000, _⟩, Γ.zero => some ⟨⟨000, by omega⟩, ⟨Turing.Dir.left, Γ.zero⟩⟩
722 | ⟨000, _⟩, Γ.one  => none
723 | ⟨001, _⟩, Γ.zero  => none
724 | ⟨001, _⟩, Γ.one   => some ⟨⟨000, by omega⟩, ⟨Turing.Dir.left, Γ.one⟩⟩
725 | ⟨_+2, _⟩, _      => by omega
726
727 def nth_cfg : (n : Nat) -> Option Cfg
728 | 0 => init []
729 | Nat.succ n => match (nth_cfg n) with
730 | none => none
731 | some cfg => step machine cfg
732
733 /--
734 Prove the following Turing Machine with the transition table never halts.
735 | State | Symbol | Next State | Move | Write |
736 |-----|-----|-----|-----|-----|
737 | 000 | zero | 000 | left | zero |
738 | 000 | one  | HALT | - | - |
739 | 001 | zero | HALT | - | - |
740 | 001 | one  | 000 | left | one  |
741 -/
742 theorem machine_never_halts : ∀ n, (nth_cfg n).isSome := by
743   sorry
744   ``
745
746 You can make your own auxiliary corollaries and theorems to support the proof, instead of only completing the part
747 ↳with the sorry. Please output the entire program and not just the last part. Please output only the program
748 ↳and add no other comment, such that your answer is a compilable lean code. Make sure to reason enough to make
749 ↳your code correct.
750
751 ### Response:
752
753 ```lean4
754 import Mathlib
755 import Aesop
756
757 set_option maxHeartbeats 0
758
759 open BigOperators Real Nat Topology Rat
760
761 /-- Given a sequence {a_n} with a_1 = 1 and a_{n+1} = 2a_n + 1,
762   prove that a_n = 2^n - 1 for all n ≥ 1 -/
763 theorem sequence_formula {a : ℕ → ℕ} (h1 : a 1 = 1)
764 (h_rec : ∀ n : ℕ, n ≥ 1 → a (n + 1) = 2 * a n + 1) :
765   ∀ n : ℕ, n ≥ 1 → a n = 2^n - 1 := by
766   /-
767   We need to prove that for a sequence \(\{a_n\}\) defined by \((a_1 = 1)\) and \((a_{n+1} = 2a_n + 1)\), the
768   ↳general term \((a_n)\) is given by \((a_n = 2^n - 1)\) for all \((n \geq 1)\).
769   1. **Base Case**: For \((n = 1)\), we have \((a_1 = 1)\). This matches \((2^1 - 1 = 1)\), so the base case holds.
770   2. **Inductive Step**: Assume \((a_n = 2^n - 1)\) holds for some \((n \geq 1)\). We need to show it holds for \((n + 1)\).
771   ↳ + 1 \).
772   - By the recurrence relation, \((a_{n+1} = 2a_n + 1)\).
773   - Using the inductive hypothesis \((a_n = 2^n - 1)\), we get:
774   ↳
775   ↳
776   ↳
777   ↳
778   ↳
779   ↳
780   ↳
781   ↳
782   ↳
783   ↳
784   ↳
785   ↳
786   ↳
787   ↳
788   ↳
789   ↳
790   ↳
791   ↳
792   ↳
793   ↳
794   ↳
795   ↳
796   ↳
797   ↳
798   ↳
799   ↳
800   ↳
801   ↳
802   ↳
803   ↳
804   ↳
805   ↳
806   ↳
807   ↳
808   ↳
809   ↳
810   ↳
811   ↳
812   ↳
813   ↳
814   ↳
815   ↳
816   ↳
817   ↳
818   ↳
819   ↳
820   ↳
821   ↳
822   ↳
823   ↳
824   ↳
825   ↳
826   ↳
827   ↳
828   ↳
829   ↳
830   ↳
831   ↳
832   ↳
833   ↳
834   ↳
835   ↳
836   ↳
837   ↳
838   ↳
839   ↳
840   ↳
841   ↳
842   ↳
843   ↳
844   ↳
845   ↳
846   ↳
847   ↳
848   ↳
849   ↳
850   ↳
851   ↳
852   ↳
853   ↳
854   ↳
855   ↳
856   ↳
857   ↳
858   ↳
859   ↳
860   ↳
861   ↳
862   ↳
863   ↳
864   ↳
865   ↳
866   ↳
867   ↳
868   ↳
869   ↳
870   ↳
871   ↳
872   ↳
873   ↳
874   ↳
875   ↳
876   ↳
877   ↳
878   ↳
879   ↳
880   ↳
881   ↳
882   ↳
883   ↳
884   ↳
885   ↳
886   ↳
887   ↳
888   ↳
889   ↳
890   ↳
891   ↳
892   ↳
893   ↳
894   ↳
895   ↳
896   ↳
897   ↳
898   ↳
899   ↳
900   ↳
901   ↳
902   ↳
903   ↳
904   ↳
905   ↳
906   ↳
907   ↳
908   ↳
909   ↳
910   ↳
911   ↳
912   ↳
913   ↳
914   ↳
915   ↳
916   ↳
917   ↳
918   ↳
919   ↳
920   ↳
921   ↳
922   ↳
923   ↳
924   ↳
925   ↳
926   ↳
927   ↳
928   ↳
929   ↳
930   ↳
931   ↳
932   ↳
933   ↳
934   ↳
935   ↳
936   ↳
937   ↳
938   ↳
939   ↳
940   ↳
941   ↳
942   ↳
943   ↳
944   ↳
945   ↳
946   ↳
947   ↳
948   ↳
949   ↳
950   ↳
951   ↳
952   ↳
953   ↳
954   ↳
955   ↳
956   ↳
957   ↳
958   ↳
959   ↳
960   ↳
961   ↳
962   ↳
963   ↳
964   ↳
965   ↳
966   ↳
967   ↳
968   ↳
969   ↳
970   ↳
971   ↳
972   ↳
973   ↳
974   ↳
975   ↳
976   ↳
977   ↳
978   ↳
979   ↳
980   ↳
981   ↳
982   ↳
983   ↳
984   ↳
985   ↳
986   ↳
987   ↳
988   ↳
989   ↳
990   ↳
991   ↳
992   ↳
993   ↳
994   ↳
995   ↳
996   ↳
997   ↳
998   ↳
999   ↳
1000  ↳
1001  ↳
1002  ↳
1003  ↳
1004  ↳
1005  ↳
1006  ↳
1007  ↳
1008  ↳
1009  ↳
1010  ↳
1011  ↳
1012  ↳
1013  ↳
1014  ↳
1015  ↳
1016  ↳
1017  ↳
1018  ↳
1019  ↳
1020  ↳
1021  ↳
1022  ↳
1023  ↳
1024  ↳
1025  ↳
1026  ↳
1027  ↳
1028  ↳
1029  ↳
1030  ↳
1031  ↳
1032  ↳
1033  ↳
1034  ↳
1035  ↳
1036  ↳
1037  ↳
1038  ↳
1039  ↳
1040  ↳
1041  ↳
1042  ↳
1043  ↳
1044  ↳
1045  ↳
1046  ↳
1047  ↳
1048  ↳
1049  ↳
1050  ↳
1051  ↳
1052  ↳
1053  ↳
1054  ↳
1055  ↳
1056  ↳
1057  ↳
1058  ↳
1059  ↳
1060  ↳
1061  ↳
1062  ↳
1063  ↳
1064  ↳
1065  ↳
1066  ↳
1067  ↳
1068  ↳
1069  ↳
1070  ↳
1071  ↳
1072  ↳
1073  ↳
1074  ↳
1075  ↳
1076  ↳
1077  ↳
1078  ↳
1079  ↳
1080  ↳
1081  ↳
1082  ↳
1083  ↳
1084  ↳
1085  ↳
1086  ↳
1087  ↳
1088  ↳
1089  ↳
1090  ↳
1091  ↳
1092  ↳
1093  ↳
1094  ↳
1095  ↳
1096  ↳
1097  ↳
1098  ↳
1099  ↳
1100  ↳
1101  ↳
1102  ↳
1103  ↳
1104  ↳
1105  ↳
1106  ↳
1107  ↳
1108  ↳
1109  ↳
1110  ↳
1111  ↳
1112  ↳
1113  ↳
1114  ↳
1115  ↳
1116  ↳
1117  ↳
1118  ↳
1119  ↳
1120  ↳
1121  ↳
1122  ↳
1123  ↳
1124  ↳
1125  ↳
1126  ↳
1127  ↳
1128  ↳
1129  ↳
1130  ↳
1131  ↳
1132  ↳
1133  ↳
1134  ↳
1135  ↳
1136  ↳
1137  ↳
1138  ↳
1139  ↳
1140  ↳
1141  ↳
1142  ↳
1143  ↳
1144  ↳
1145  ↳
1146  ↳
1147  ↳
1148  ↳
1149  ↳
1150  ↳
1151  ↳
1152  ↳
1153  ↳
1154  ↳
1155  ↳
1156  ↳
1157  ↳
1158  ↳
1159  ↳
1160  ↳
1161  ↳
1162  ↳
1163  ↳
1164  ↳
1165  ↳
1166  ↳
1167  ↳
1168  ↳
1169  ↳
1170  ↳
1171  ↳
1172  ↳
1173  ↳
1174  ↳
1175  ↳
1176  ↳
1177  ↳
1178  ↳
1179  ↳
1180  ↳
1181  ↳
1182  ↳
1183  ↳
1184  ↳
1185  ↳
1186  ↳
1187  ↳
1188  ↳
1189  ↳
1190  ↳
1191  ↳
1192  ↳
1193  ↳
1194  ↳
1195  ↳
1196  ↳
1197  ↳
1198  ↳
1199  ↳
1200  ↳
1201  ↳
1202  ↳
1203  ↳
1204  ↳
1205  ↳
1206  ↳
1207  ↳
1208  ↳
1209  ↳
1210  ↳
1211  ↳
1212  ↳
1213  ↳
1214  ↳
1215  ↳
1216  ↳
1217  ↳
1218  ↳
1219  ↳
1220  ↳
1221  ↳
1222  ↳
1223  ↳
1224  ↳
1225  ↳
1226  ↳
1227  ↳
1228  ↳
1229  ↳
1230  ↳
1231  ↳
1232  ↳
1233  ↳
1234  ↳
1235  ↳
1236  ↳
1237  ↳
1238  ↳
1239  ↳
1240  ↳
1241  ↳
1242  ↳
1243  ↳
1244  ↳
1245  ↳
1246  ↳
1247  ↳
1248  ↳
1249  ↳
1250  ↳
1251  ↳
1252  ↳
1253  ↳
1254  ↳
1255  ↳
1256  ↳
1257  ↳
1258  ↳
1259  ↳
1260  ↳
1261  ↳
1262  ↳
1263  ↳
1264  ↳
1265  ↳
1266  ↳
1267  ↳
1268  ↳
1269  ↳
1270  ↳
1271  ↳
1272  ↳
1273  ↳
1274  ↳
1275  ↳
1276  ↳
1277  ↳
1278  ↳
1279  ↳
1280  ↳
1281  ↳
1282  ↳
1283  ↳
1284  ↳
1285  ↳
1286  ↳
1287  ↳
1288  ↳
1289  ↳
1290  ↳
1291  ↳
1292  ↳
1293  ↳
1294  ↳
1295  ↳
1296  ↳
1297  ↳
1298  ↳
1299  ↳
1300  ↳
1301  ↳
1302  ↳
1303  ↳
1304  ↳
1305  ↳
1306  ↳
1307  ↳
1308  ↳
1309  ↳
1310  ↳
1311  ↳
1312  ↳
1313  ↳
1314  ↳
1315  ↳
1316  ↳
1317  ↳
1318  ↳
1319  ↳
1320  ↳
1321  ↳
1322  ↳
1323  ↳
1324  ↳
1325  ↳
1326  ↳
1327  ↳
1328  ↳
1329  ↳
1330  ↳
1331  ↳
1332  ↳
1333  ↳
1334  ↳
1335  ↳
1336  ↳
1337  ↳
1338  ↳
1339  ↳
1340  ↳
1341  ↳
1342  ↳
1343  ↳
1344  ↳
1345  ↳
1346  ↳
1347  ↳
1348  ↳
1349  ↳
1350  ↳
1351  ↳
1352  ↳
1353  ↳
1354  ↳
1355  ↳
1356  ↳
1357  ↳
1358  ↳
1359  ↳
1360  ↳
1361  ↳
1362  ↳
1363  ↳
1364  ↳
1365  ↳
1366  ↳
1367  ↳
1368  ↳
1369  ↳
1370  ↳
1371  ↳
1372  ↳
1373  ↳
1374  ↳
1375  ↳
1376  ↳
1377  ↳
1378  ↳
1379  ↳
1380  ↳
1381  ↳
1382  ↳
1383  ↳
1384  ↳
1385  ↳
1386  ↳
1387  ↳
1388  ↳
1389  ↳
1390  ↳
1391  ↳
1392  ↳
1393  ↳
1394  ↳
1395  ↳
1396  ↳
1397  ↳
1398  ↳
1399  ↳
1400  ↳
1401  ↳
1402  ↳
1403  ↳
1404  ↳
1405  ↳
1406  ↳
1407  ↳
1408  ↳
1409  ↳
1410  ↳
1411  ↳
1412  ↳
1413  ↳
1414  ↳
1415  ↳
1416  ↳
1417  ↳
1418  ↳
1419  ↳
1420  ↳
1421  ↳
1422  ↳
1423  ↳
1424  ↳
1425  ↳
1426  ↳
1427  ↳
1428  ↳
1429  ↳
1430  ↳
1431  ↳
1432  ↳
1433  ↳
1434  ↳
1435  ↳
1436  ↳
1437  ↳
1438  ↳
1439  ↳
1440  ↳
1441  ↳
1442  ↳
1443  ↳
1444  ↳
1445  ↳
1446  ↳
1447  ↳
1448  ↳
1449  ↳
1450  ↳
1451  ↳
1452  ↳
1453  ↳
1454  ↳
1455  ↳
1456  ↳
1457  ↳
1458  ↳
1459  ↳
1460  ↳
1461  ↳
1462  ↳
1463  ↳
1464  ↳
1465  ↳
1466  ↳
1467  ↳
1468  ↳
1469  ↳
1470  ↳
1471  ↳
1472  ↳
1473  ↳
1474  ↳
1475  ↳
1476  ↳
1477  ↳
1478  ↳
1479  ↳
1480  ↳
1481  ↳
1482  ↳
1483  ↳
1484  ↳
1485  ↳
1486  ↳
1487  ↳
1488  ↳
1489  ↳
1490  ↳
1491  ↳
1492  ↳
1493  ↳
1494  ↳
1495  ↳
1496  ↳
1497  ↳
1498  ↳
1499  ↳
1500  ↳
1501  ↳
1502  ↳
1503  ↳
1504  ↳
1505  ↳
1506  ↳
1507  ↳
1508  ↳
1509  ↳
1510  ↳
1511  ↳
1512  ↳
1513  ↳
1514  ↳
1515  ↳
1516  ↳
1517  ↳
1518  ↳
1519  ↳
1520  ↳
1521  ↳
1522  ↳
1523  ↳
1524  ↳
1525  ↳
1526  ↳
1527  ↳
1528  ↳
1529  ↳
1530  ↳
1531  ↳
1532  ↳
1533  ↳
1534  ↳
1535  ↳
1536  ↳
1537  ↳
1538  ↳
1539  ↳
1540  ↳
1541  ↳
1542  ↳
1543  ↳
1544  ↳
1545  ↳
1546  ↳
1547  ↳
1548  ↳
1549  ↳
1550  ↳
1551  ↳
1552  ↳
1553  ↳
1554  ↳
1555  ↳
1556  ↳
1557  ↳
1558  ↳
1559  ↳
1560  ↳
1561  ↳
1562  ↳
1563  ↳
1564  ↳
1565  ↳
1566  ↳
1567  ↳
1568  ↳
1569  ↳
1570  ↳
1571  ↳
1572  ↳
1573  ↳
1574  ↳
1575  ↳
1576  ↳
1577  ↳
1578  ↳
1579  ↳
1580  ↳
1581  ↳
1582  ↳
1583  ↳
1584  ↳
1585  ↳
1586  ↳
1587  ↳
1588  ↳
1589  ↳
1590  ↳
1591  ↳
1592  ↳
1593  ↳
1594  ↳
1595  ↳
1596  ↳
1597  ↳
1598  ↳
1599  ↳
1600  ↳
1601  ↳
1602  ↳
1603  ↳
1604  ↳
1605  ↳
1606  ↳
1607  ↳
1608  ↳
1609  ↳
1610  ↳
1611  ↳
1612  ↳
1613  ↳
1614  ↳
1615  ↳
1616  ↳
1617  ↳
1618  ↳
1619  ↳
1620  ↳
1621  ↳
1622  ↳
1623  ↳
1624  ↳
1625  ↳
1626  ↳
1627  ↳
1628  ↳
1629  ↳
1630  ↳
1631  ↳
1632  ↳
1633  ↳
1634  ↳
1635  ↳
1636  ↳
1637  ↳
1638  ↳
1639  ↳
1640  ↳
1641  ↳
1642  ↳
1643  ↳
1644  ↳
1645  ↳
1646  ↳
1647  ↳
1648  ↳
1649  ↳
1650  ↳
1651  ↳
1652  ↳
1653  ↳
1654  ↳
1655  ↳
1656  ↳
1657  ↳
1658  ↳
1659  ↳
1660  ↳
1661  ↳
1662  ↳
1663  ↳
1664  ↳
1665  ↳
1666  ↳
1667  ↳
1668  ↳
1669  ↳
1670  ↳
1671  ↳
1672  ↳
1673  ↳
1674  ↳
1675  ↳
1676  ↳
1677  ↳
1678  ↳
1679  ↳
1680  ↳
1681  ↳
1682  ↳
1683  ↳
1684  ↳
1685  ↳
1686  ↳
1687  ↳
1688  ↳
1689  ↳
1690  ↳
1691  ↳
1692  ↳
1693  ↳
1694  ↳
1695  ↳
1696  ↳
1697  ↳
1698  ↳
1699  ↳
1700  ↳
1701  ↳
1702  ↳
1703  ↳
1704  ↳
1705  ↳
1706  ↳
1707  ↳
1708  ↳
1709  ↳
1710  ↳
1711  ↳
1712  ↳
1713  ↳
1714  ↳
1715  ↳
1716  ↳
1717  ↳
1718  ↳
1719  ↳
1720  ↳
1721  ↳
1722  ↳
1723  ↳
1724  ↳
1725  ↳
1726  ↳
1727  ↳
1728  ↳
1729  ↳
1730  ↳
1731  ↳
1732  ↳
1733  ↳
1734  ↳
1735  ↳
1736  ↳
1737  ↳
1738  ↳
1739  ↳
1740  ↳
1741  ↳
1742  ↳
1743  ↳
1744  ↳
1745  ↳
1746  ↳
1747  ↳
1748  ↳
1749  ↳
1750  ↳
1751  ↳
1752  ↳
1753  ↳
1754  ↳
1755  ↳
1756  ↳
1757  ↳
1758  ↳
1759  ↳
1760  ↳
1761  ↳
1762  ↳
1763  ↳
1764  ↳
1765  ↳
1766  ↳
1767  ↳
1768  ↳
1769  ↳
1770  ↳
1771  ↳
1772  ↳
1773  ↳
1774  ↳
1775  ↳
1776  ↳
1777  ↳
1778  ↳
1779  ↳
1780  ↳
1781  ↳
1782  ↳
1783  ↳
1784  ↳
1785  ↳
1786  ↳
1787  ↳
1788  ↳
1789  ↳
1790  ↳
1791  ↳
1792  ↳
1793  ↳
1794  ↳
1795  ↳
1796  ↳
1797  ↳
1798  ↳
1799  ↳
1800  ↳
1801  ↳
1802  ↳
1803  ↳
1804  ↳
1805  ↳
1806  ↳
1807  ↳
1808  ↳
1809  ↳
1810  ↳
1811  ↳
1812  ↳
1813  ↳
1814  ↳
1815  ↳
1816  ↳
1817  ↳
1818  ↳
1819  ↳
1820  ↳
1821  ↳
1822  ↳
1823  ↳
1824  ↳
1825  ↳
1826  ↳
1827  ↳
1828  ↳
1829  ↳
1830  ↳
1831  ↳
1832  ↳
1833  ↳
1834  ↳
1835  ↳
1836  ↳
1837  ↳
1838  ↳
1839  ↳
1840  ↳
1841  ↳
1842  ↳
1843  ↳
1844  ↳
1845  ↳
1846  ↳
1847  ↳
1848  ↳
1849  ↳
1850  ↳
1851  ↳
1852  ↳
1853  ↳
1854  ↳
1855  ↳
1856  ↳
1857  ↳
1858  ↳
1859  ↳
1860  ↳
1861  ↳
1862  ↳
1863  ↳
1864  ↳
1865  ↳
1866  ↳
1867  ↳
1868  ↳
1869  ↳
1870  ↳
1871  ↳
1872  ↳
1873  ↳
1874  ↳
1875  ↳
1876  ↳
1877  ↳
1878  ↳
1879  ↳
1880  ↳
1881  ↳
1882  ↳
1883  ↳
1884  ↳
1885  ↳
1886  ↳
1887  ↳
1888  ↳
1889  ↳
1890  ↳
1891  ↳
1892  ↳
1893  ↳
1894  ↳
1895  ↳
1896  ↳
1897  ↳
1898  ↳
1899  ↳
1900  ↳
1901  ↳
1902  ↳
1903  ↳
1904  ↳
1905  ↳
1906  ↳
1907  ↳
1908  ↳
1909  ↳
1910  ↳
1911  ↳
1912  ↳
1913  ↳
1914  ↳
1915  ↳
1916  ↳
1917  ↳
1918  ↳
1919  ↳
1920  ↳
1921  ↳
1922  ↳
1923  ↳
1924  ↳
1925  ↳
1926  ↳
1927  ↳
1928  ↳
1929  ↳
1930  ↳
1931  ↳
1932  ↳
1933  ↳
1934  ↳
1935  ↳
1936  ↳
1937  ↳
1938  ↳
1939  ↳
1940  ↳
1941  ↳
1942  ↳
1943  ↳
1944  ↳
1945  ↳
1946  ↳
1947  ↳
1948  ↳
1949  ↳
1950  ↳
1951  ↳
1952  ↳
1953  ↳
1954  ↳
1955  ↳
1956  ↳
1957  ↳
1958  ↳
1959  ↳
1960  ↳
1961  ↳
1962  ↳
1963  ↳
1964  ↳
1965  ↳
1966  ↳
1967  ↳
1968  ↳
1969  ↳
1970  ↳
1971  ↳
1972  ↳
1973  ↳
1974  ↳
1975  ↳
1976  ↳
1977  ↳
1978  ↳
1979  ↳
1980  ↳
1981  ↳
1982  ↳
1983  ↳
1984  ↳
1985  ↳
1986  ↳
1987  ↳
1988  ↳
1989  ↳
1990  ↳
1991  ↳
1992  ↳
1993  ↳
1994  ↳
1995  ↳
1996  ↳
1997  ↳
1998  ↳
1999  ↳
2000  ↳
2001  ↳
2002  ↳
2003  ↳
2004  ↳
2005  ↳
2006  ↳
2007  ↳
2008  ↳
2009  ↳
2010  ↳
2011  ↳
2012  ↳
2013  ↳
2014  ↳
2015  ↳
2016  ↳
2017  ↳
2018  ↳
2019  ↳
2020  ↳
2021  ↳
2022  ↳
2023  ↳
2024  ↳
2025  ↳
2026  ↳
2027  ↳
2028  ↳
2029  ↳
2030  ↳
2031  ↳
2032  ↳
2033  ↳
2034  ↳
2035  ↳
2036  ↳
2037  ↳
2038  ↳
2039  ↳
2040  ↳
2041  ↳
2042  ↳
2043  ↳
2044  ↳
2045  ↳
2046  ↳
2047  ↳
2048  ↳
2049  ↳
2050  ↳
2051  ↳
2052  ↳
2053  ↳
2054  ↳
2055  ↳
2056  ↳
2057  ↳
2058  ↳
2059  ↳
2060  ↳
2061  ↳
2062  ↳
2063  ↳
2064  ↳
2065  ↳
2066  ↳
2067  ↳
2068  ↳
2069  ↳
2070  ↳
2071  ↳
2072  ↳
2073  ↳
2074  ↳
2075  ↳
2076  ↳
2077  ↳
2078  ↳
2079  ↳
2080  ↳
2081  ↳
2082  ↳
2083  ↳
2084  ↳
2085  ↳
2086  ↳
2087  ↳
2088  ↳
2089  ↳
2090  ↳
2091  ↳
2092  ↳
2093  ↳
2094  ↳
2095  ↳
2096  ↳
2097  ↳
2098  ↳
2099  ↳
2100  ↳
2101  ↳
2102  ↳
2103  ↳
2104  ↳
2105  ↳
2106  ↳
2107  ↳
2108  ↳
2109  ↳
2110  ↳
2111  ↳
2112  ↳
2113  ↳
2114  ↳
2115  ↳
2116  ↳
2117  ↳
2118  ↳
2119  ↳
2120  ↳
2121  ↳
2122  ↳
2123  ↳
2124  ↳
2125  ↳
2126  ↳
2127  ↳
2128  ↳
2129  ↳
2130  ↳
2131  ↳
2132  ↳
2133  ↳
2134  ↳
2135  ↳
2136  ↳
2137  ↳
2138  ↳
2139  ↳
2140  ↳
2141  ↳
2142  ↳
2143  ↳
2144  ↳
2145  ↳
2146  ↳
2147  ↳
2148  ↳
2149  ↳
2150  ↳
2151  ↳
2152  ↳
2153  ↳
2154  ↳
2155  ↳
2156  ↳
2157  ↳
2158  ↳
2159  ↳
2160  ↳
2161  ↳
2162  ↳
2163  ↳
2164  ↳
2165  ↳
2166  ↳
2167  ↳
2168  ↳
2169  ↳
2170  ↳
2171  ↳
2172  ↳
2173  ↳
2174  ↳
2175  ↳
2176  ↳
2177  ↳
2178  ↳
2179  ↳
2180  ↳
2181  ↳
2182  ↳
2183  ↳
2184  ↳
2185  ↳
2186  ↳
2187  ↳
2188  ↳
2189  ↳
2190  ↳
2191  ↳
2192  ↳
2193  ↳
2194  ↳
2195  ↳
2196  ↳
2197  ↳
2198  ↳
2199  ↳
2200  ↳
2201  ↳
2202  ↳
2203  ↳
2204  ↳
2205  ↳
2206  ↳
2207  ↳
2208  ↳
2209  ↳
2210  ↳
2211  ↳
2212  ↳
2213  ↳
2214  ↳
2215  ↳
2216  ↳
2217  ↳
2218  ↳
2219  ↳
2220  ↳
2221  ↳
2222  ↳
2223  ↳
2224  ↳
2225  ↳
2226  ↳
2227  ↳
2228  ↳
2229  ↳
2230  ↳
2231  ↳
2232  ↳
2233  ↳
2234  ↳
2235  ↳
2236  ↳
2237  ↳
2238  ↳
2239  ↳
2240  ↳
2241  ↳
2242  ↳
2243  ↳
2244  ↳
2245  ↳
2246  ↳
2247  ↳
2248  ↳
2249  ↳
2250  ↳
2251  ↳
2252  ↳
2253  ↳
2254  ↳
2255  ↳
2256  ↳
2257  ↳
2258  ↳
2259  ↳
2260  ↳
2261  ↳
2262  ↳
2263  ↳
2264  ↳
2265  ↳
2266  ↳
2267  ↳
2268  ↳
2269  ↳
2270  ↳
2271  ↳
2272  ↳
2273  ↳
2274  ↳
2275  ↳
2276  ↳
2277  ↳
2278  ↳
2279  ↳
2280  ↳
2281  ↳
2282  ↳
2283  ↳
2284  ↳
2285  ↳
2286  ↳
2287  ↳
2288  ↳
2289  ↳
2290  ↳
2291  ↳
2292  ↳
2293  ↳
2294  ↳
2295  ↳
2296  ↳
2297  ↳
2298  ↳
2299  ↳
2300  ↳
2301  ↳
2302  ↳
2303  ↳
2304  ↳
2305  ↳
2306  ↳
2307  ↳
2308  ↳
2309  ↳
2310  ↳
2311  ↳
2312  ↳
2313  ↳
2314  ↳
2315  ↳
2316  ↳
2317  ↳
2318  ↳
2319  ↳
2320  ↳
2321  ↳
2322  ↳
2323  ↳
2324  ↳
2325  ↳
2326  ↳
2327  ↳
2328  ↳
2329  ↳
2330  ↳
2331  ↳
2332  ↳
2333  ↳
2334  ↳
2335  ↳
2336  ↳
2337  ↳
2338  ↳
2339  ↳
2340  ↳
2341  ↳
2342  ↳
2343  ↳
2344  ↳
2345  ↳
2346  ↳
2347  ↳
2348  ↳
2349  ↳
2350  ↳
2351  ↳
2352  ↳
2353  ↳
2354  ↳
2355  ↳
2356  ↳
2357  ↳
2358  ↳
2359  ↳
2360  ↳
2361  ↳
2362  ↳
2363  ↳
2364  ↳
2365  ↳
2366  ↳
2367  ↳
2368  ↳
2369  ↳
2370  ↳
2371  ↳
2372  ↳
2373  ↳
2374  ↳
2375  ↳
2376  ↳
2377  ↳
2378  ↳
2379  ↳
2380  ↳
2381  ↳
2382  ↳
2383  ↳
2384  ↳
2385  ↳
2386  ↳
2387  ↳
2388  ↳
2389  ↳
2390  ↳
2391  ↳
2392  ↳
2393  ↳
2394  ↳
2395  ↳
2396  ↳
2397  ↳
2398  ↳
2399  ↳
2400  ↳
2401  ↳
2402  ↳
2403  ↳
2404  ↳
2405  ↳

```

```

770 -- When n = 0, the induction hypothesis does not apply, so we directly use the initial condition
771 simp_all
772 | succ n =>
773 -- For n = succ n, use the recurrence relation and the inductive hypothesis
774 simp_all [Nat.succ_eq_add_one, Nat.pow_succ, Nat.mul_succ]
775 -- Simplify the expression to match the form  $2^{(n + 1)} - 1$ 
776 ring
777 <|> omega

```

A complete step lemma challenge example. Below we demonstrate a complete record of evaluation, including the problem statement in Lean 4 and Markdown, our evaluation prompt and the model response:

```

781 Theorem library:
782 ```lean4
783 import Mathlib.Tactic.Lemma
784 import Mathlib.Tactic.NthRewrite
785
786 @[simp]
787 theorem bv32_and_not_self (x : BitVec 32) :
788   x &&& ~~~x = 0 := by
789     simp
790
791 @[simp]
792 theorem bv32_not_not (x : BitVec 32) :
793   ~~~~~x = x := by
794     simp
795
796 @[simp]
797 theorem bv32_or_not_self (x : BitVec 32) :
798   x ||| ~~~x = BitVec.allOnes 32 := by
799     simp
800
801 @[simp]
802 theorem bv32_not_or_self (x : BitVec 32) :
803   ~~~x ||| x = BitVec.allOnes 32 := by
804     simp
805
806 @[simp]
807 theorem bv32_neg_mul (x y : BitVec 32) :
808   -x * y = -(x * y) := by
809     simp
810
811 theorem bv32_not_and (x y : BitVec 32) :
812   ~~~(x &&& y) = ~~~x ||| ~~~y := by
813     rw [BitVec.not_and]
814
815 theorem bv32_not_or (x y : BitVec 32) :
816   ~~~(x ||| y) = ~~~x &&& ~~~y := by
817     rw [BitVec.not_or]
818
819 theorem bv32_not_xor_eq_or (x y : (BitVec 32)) :
820   ~~~(x ^^^ y) = (~~~x &&& ~~~y) ||| (x &&& y) := by
821     ext i
822     simp
823     cases h1 : x[i] <|> cases h2 : y[i]
824     simp
825     simp
826     simp
827     simp
828
829 theorem bv32_xor_eq_or (x y : (BitVec 32)) :
830   (x ^^^ y) = (~~~x &&& y) ||| (x &&& ~~~y) := by
831     ext i
832     simp
833     cases h1 : x[i] <|> cases h2 : y[i]
834     simp
835     simp
836     simp
837     simp

```

```

825
826 theorem bv32_x_distr (x y: BitVec 32) :
827   x = (x &&& y) ||| (x &&& ~~~y) := by
828     ext i
829     simp
830     simp [ Bool.and_or_distrib_left]
831
832 theorem bv32_y_distr (x y: BitVec 32) :
833   y = (x &&& y ||| ~~~x &&& y) := by
834     ext i
835     simp
836     simp [ Bool.and_or_distrib_right]
837
838 theorem bv32_add_assoc (x y z : BitVec 32) :
839   x + y + z = x + (y + z) := by
840     rw [BitVec.add_assoc]
841
842 theorem bv32_add_comm (x y : BitVec 32) :
843   x + y = y + x := by
844     rw [BitVec.add_comm]
845
846 theorem bv32_add_neg_eq_sub {x y : BitVec 32} :
847   x + -y = x - y := by
848     rw [BitVec.add_neg_eq_sub]
849
850 theorem bv32_mul_comm (x y : BitVec 32) :
851   x * y = y * x := by
852     rw [BitVec.mul_comm]
853
854 theorem bv32_var_mul_comm (x y z: BitVec 32) :
855   (x &&& y) * z = z * (x &&& y) := by
856     rw [BitVec.mul_comm]
857
858 theorem bv32_mul_add (x y z : BitVec 32) :
859   x * (y + z) = x * y + x * z := by
860     rw [BitVec.mul_add]
861
862 theorem bv32_neg_eq_mul (x : BitVec 32) :
863   -x = x * (BitVec.allOnes 32) := by
864     rw [ BitVec.neg_one_eq_allOnes]
865     rw [BitVec.mul_neg]
866     rw [BitVec.mul_one]
867
868 theorem bv32_add_mul_one (x y : BitVec 32) :
869   x + x * y = x * (1#32 + y) := by
870     rw [BitVec.mul_add]
871     rw [BitVec.mul_one]
872
873 /--
874 1: x &&& y
875 2: ~~~x &&& y
876 3: x &&& ~~~y
877 4: ~~~x &&& ~~~y
878 -/
879 theorem bv32_or_eq_addl2 (x y : BitVec 32) :
880   (x &&& y) ||| (~~~x &&& y) = (x &&& y) + (~~~x &&& y) := by
881     apply Eq.symm
882     apply BitVec.add_eq_or_of_and_eq_zero
883     simp [ BitVec.and_assoc]
884     simp [BitVec.and_comm _ (~~~x)]
885     simp [ BitVec.and_assoc]
886
887 theorem bv32_or_eq_addl3 (x y : BitVec 32) :
888   (x &&& y) ||| (x &&& ~~~y) = (x &&& y) + (x &&& ~~~y) := by
889     apply Eq.symm
890     apply BitVec.add_eq_or_of_and_eq_zero
891     simp [ BitVec.and_assoc]
892     simp [BitVec.and_comm _ x]
893     simp [BitVec.and_assoc]
894
895 theorem bv32_or_eq_addl4 (x y : BitVec 32) :

```



```

880 (x &&& y) ||| (~~~x &&& ~~~y) = (x &&& y) + (~~~x &&& ~~~y) := by
881 apply Eq.symm
882 apply BitVec.add_eq_or_of_and_eq_zero
883 simp [BitVec.and_assoc]
884 simp [BitVec.and_comm _ (~~~x)]
885 simp [BitVec.and_assoc]
886
887 theorem bv32_or_eq_add21 (x y : BitVec 32) :
888 (~~~x &&& y) ||| (x &&& y) = (~~~x &&& y) + (x &&& y) := by
889 apply Eq.symm
890 apply BitVec.add_eq_or_of_and_eq_zero
891 simp [BitVec.and_assoc]
892 simp [BitVec.and_comm _ x]
893 simp [BitVec.and_assoc]
894
895 theorem bv32_or_eq_add23 (x y : BitVec 32) :
896 (~~~x &&& y) ||| (x &&& ~~~y) = (~~~x &&& y) + (x &&& ~~~y) := by
897 apply Eq.symm
898 apply BitVec.add_eq_or_of_and_eq_zero
899 simp [BitVec.and_assoc]
900 simp [BitVec.and_comm _ x]
901 simp [BitVec.and_assoc]
902
903 theorem bv32_or_eq_add31 (x y : BitVec 32) :
904 (x &&& ~~~y) ||| (x &&& y) = (x &&& ~~~y) + (x &&& y) := by
905 apply Eq.symm
906 apply BitVec.add_eq_or_of_and_eq_zero
907 simp [BitVec.and_assoc]
908 simp [BitVec.and_comm _ x]
909 simp [BitVec.and_assoc]
910
911 theorem bv32_or_eq_add32 (x y : BitVec 32) :
912 (x &&& ~~~y) ||| (~~~x &&& y) = (x &&& ~~~y) + (~~~x &&& y) := by
913 apply Eq.symm
914 apply BitVec.add_eq_or_of_and_eq_zero
915 simp [BitVec.and_assoc]
916 simp [BitVec.and_comm _ (~~~x)]
917 simp [BitVec.and_assoc]
918
919 theorem bv32_or_eq_add41 (x y : BitVec 32) :
920 (~~~x &&& ~~~y) ||| (x &&& y) = (~~~x &&& ~~~y) + (x &&& y) := by
921 apply Eq.symm
922 apply BitVec.add_eq_or_of_and_eq_zero
923 simp [BitVec.and_assoc]
924 simp [BitVec.and_comm _ x]
925 simp [BitVec.and_assoc]
926
927 theorem bv32_or_eq_add_three (x y : BitVec 32) :
928 (x ||| y) = (x &&& ~~~y) + (x &&& y) + (~~~x &&& y) := by
929 nth_rw 1 [bv32_y_distr x y]
930 nth_rw 1 [bv32_x_distr x y]
931 simp [BitVec.or_assoc]
932 simp [BitVec.or_comm _ (x &&& y)]
933 simp [BitVec.or_assoc]
934 rw [BitVec.or_comm (x &&& y)]
935 apply Eq.symm
936 rw [BitVec.add_eq_or_of_and_eq_zero]
937 rw [BitVec.add_eq_or_of_and_eq_zero]
938 simp [BitVec.and_assoc]
939 simp [BitVec.and_comm _ x]
940 simp [BitVec.and_assoc]
941 rw [BitVec.add_comm]
942 rw [bv32_or_eq_add13]
943 rw [bv32_x_distr x y]
944 simp [BitVec.and_assoc]
945
946 theorem bv32_sum_all (x y : BitVec 32) :
947 (~~~x &&& ~~~y) + (~~~x &&& y) + (x &&& y) + (x &&& ~~~y) = BitVec.allOnes 32 := by
948 simp [BitVec.add_comm _ (~~~x &&& y)]
949 simp [BitVec.add_comm _ (x &&& _)]
950 simp [BitVec.add_assoc]

```

```

935 rw [BitVec.add_eq_or_of_and_eq_zero]
936 rw [bv32_or_eq_add_three x y]
937 nth_rw 1 [bv32_x_distr x y]
938 simp [BitVec.or_comm _ y]
939 nth_rw 1 [bv32_y_distr x y]
940 simp [BitVec.or_assoc]
941 simp [BitVec.or_comm _ (x &&& y)]
942 simp [BitVec.or_assoc]
943 simp [BitVec.or_comm _ (~x &&& _)]
944 simp [BitVec.or_assoc]
945 simp [BitVec.or_comm _ (~x &&& y)]
946 rw [bv32_x_distr (~x) y]
947 rw [BitVec.or_assoc]
948 rw [bv32_x_distr x y]
949 simp
950 rw [bv32_or_eq_add_three x y]
951 simp [BitVec.not_or]
952
953 theorem bv32_self_eq_neg_mul (x : BitVec 32) :
954   x = -x * (BitVec.allOnes 32) := by
955   rw [BitVec.neg_mul]
956   rw [BitVec.mul_comm]
957   rw [BitVec.neg_mul]
958   simp [BitVec.neg_one_eq_allOnes]
959
960 theorem bv32_not_self_and_not (x y : BitVec 32) :
961   ~(x &&& ~x) = (~x &&& ~y) + (~x &&& y) + (x &&& y) + (x &&& ~y) := by
962   rw [bv32_not_and]
963   rw [BitVec.not_not]
964   rw [bv32_not_or_self]
965   rw [bv32_sum_all]
966
967 /--
968 Let x,y be 32-bit bit-vectors. Prove the equivalence of the following two expressions:
969 
$$\neg(x \wedge \neg y) \leftrightarrow \neg(x \vee \neg y) - 7 \cdot \neg(x \oplus y) + 11 \cdot \neg(x \wedge y) - 5 \cdot \neg(x \vee y) - 13 \cdot \neg(x \wedge \neg y) + 6 \cdot \neg(x \wedge y)$$

970 -/
971
972 theorem mba_challenge_02a2f35e (x y : BitVec 32) :
973   2#32 * ~(x &&& ~y) - 1#32 * (x ||| ~y) - 7#32 *
974     (x ^ y) + 11#32 * ~(x &&& y) - 5#32 * ~(x ||| y) - 13#32 *
975     (x ||| ~y) + 6#32 * (x &&& y) = 10#32 * (x &&&
976     ~y) := by
977   simp only [bv32_add_neg_eq_sub] /- step 1 -/
978   simp only [bv32_neg_mul] /- step 2 -/
979   simp only [bv32_not_and] /- step 3 -/
980   simp only [bv32_not_or] /- step 4 -/
981   sorry
982
983 lemma mba_challenge_02a2f35e_lhs_step_1 (x y : BitVec 32) :
984   2#32 * ~(x &&& ~y) - 1#32 * (x ||| ~y) - 7#32 *
985     (x ^ y) + 11#32 * ~(x &&& y) - 5#32 * ~(x ||| y) - 13#32 *
986     (x ||| ~y) + 6#32 * (x &&& y) = 2#32 *
987     (x ||| y) + -(13#32 * ~(x ||| ~y)) + 6#32 * (x &&& y) := by
988   simp only [bv32_add_neg_eq_sub]
989
990 lemma mba_challenge_02a2f35e_lhs_step_2 (x y : BitVec 32) :
991   2#32 * ~(x &&& ~y) + -(1#32 * (x ||| ~y)) + -(7#32 *
992     (x ^ y)) + 11#32 * ~(x &&& y) + -(5#32 * ~(x ||| y)) +
993     -(13#32 * ~(x ||| ~y)) + 6#32 * (x &&& y) = 2#32 *
994     (x ||| y) + -13#32 * ~(x ||| ~y) + 6#32 * (x &&& y) := by
995   simp only [bv32_neg_mul]
996
997 lemma mba_challenge_02a2f35e_lhs_step_3 (x y : BitVec 32) :
998   2#32 * ~(x &&& ~y) + -1#32 * (x ||| ~y) + -7#32 *
999     (x ^ y) + 11#32 * ~(x &&& y) + -5#32 * ~(x ||| y) +
1000     -13#32 * ~(x ||| ~y) + 6#32 * (x &&& y) =
1001     2#32 * (~x ||| ~y) + -1#32 * (x ||| ~y) + -7#32 *
1002     (x ^ y) + 11#32 * (~x ||| ~y) + -5#32 *
1003     (x ||| y) + -13#32 * ~(x ||| ~y) + 6#32 * (x &&& y) := by
1004   simp only [bv32_not_and]
1005
1006 lemma mba_challenge_02a2f35e_lhs_step_4 (x y : BitVec 32) :
1007   2#32 * (~x ||| ~y) + -1#32 * (x ||| ~y) + -7#32 *
1008     (x ^ y) + 11#32 * (~x ||| ~y) + -5#32 * ~(x ||| y) +
1009     -13#32 * ~(x ||| ~y) + 6#32 * (x &&& y) = 2#32 *
1010     (~x ||| ~y) + -1#32 * (x ||| ~y) + -7#32 * (x ^ y) +
1011     11#32 * (~x ||| ~y) + -5#32 * (~x &&& ~y) + -13#32 *
1012     (~x &&& ~y) + 6#32 * (x &&& y) := by
1013   simp only [bv32_not_or]
1014
1015 lemma mba_challenge_02a2f35e_lhs_step_5 (x y : BitVec 32) :
1016   2#32 * (~x ||| ~y) + -1#32 * (x ||| ~y) + -7#32 *
1017     (x ^ y) + 11#32 * (~x ||| ~y) + -5#32 * (~x &&& ~y) +
1018     -13#32 * (~x &&& ~y) + 6#32 * (x &&& y) = (x

```

```

990      ↪ &&& y) = 2#32 * (~~~x ||| ~~~~~y) + -1#32 * (x ||| ~~~y) + -7#32 * (~~~x &&& ~~~y ||| x &&& y) + 11#32 *
991      ↪ (~~~x ||| ~~~y) + -5#32 * (~~~x &&& ~~~y) + -13#32 * (~~~x &&& ~~~~~y) + 6#32 * (x &&& y) := by
992      simp only [<Theorem here>]
993
994 You are proving a single step of theorem 'mba_challenge_02a2f35e'. Based on the previous steps, select exactly one
995      ↪ appropriate theorem from the Theorem library above and insert it in place of <Theorem here> to complete the
996      ↪ proof of lemma 'mba_challenge_02a2f35e_lhs_step_5'. Do not modify any code from previous theorems or lemmas.
997      ↪ You are also not allowed to add any other tactics to the proof body.
998
999 ### Response:
1000 You need the lemma that rewrites ~~~(x ^^^ y) into (~~~x &&& ~~~y) ||| (x &&& y). That is
1001 bv32_not_xor_eq_or
1002 So your step 5 becomes:
1003 lemma mba_challenge_02a2f35e_lhs_step_5 (x y : BitVec 32) :
1004     2#32 * (~~~x ||| ~~~~~y) + -1#32 * (x ||| ~~~y) + -7#32 * (~~~x ^^^ y) + 11#32 * (~~~x ||| ~~~y) + -5#32 * (~~~x
1005     ↪ &&& ~~~y) + -13#32 * (~~~x &&& ~~~~~y) + 6#32 * (x &&& y) = 2#32 * (~~~x ||| ~~~~~y) + -1#32 * (x ||| ~~~
1006     ↪ y) + -7#32 * (~~~x &&& ~~~y ||| x &&& y) + 11#32 * (~~~x ||| ~~~y) + -5#32 * (~~~x &&& ~~~y) + -13#32 *
1007     ↪ (~~~x &&& ~~~~~y) + 6#32 * (x &&& y) := by
1008     simp only [bv32_not_xor_eq_or]
1009 This completes the requested rewrite.

```

D. Template-Based Formal-Informal Translation

We use the following Python program to produce problem descriptions in rigorously aligned Lean-Markdown pair.

```

1010 import srsly
1011 from tqdm import tqdm
1012
1013 output_dir = "busy_beaver/prompts_to_run"
1014
1015 tables = list(srsly.read_jsonl("busy_beaver/bb_with_table.jsonl"))
1016
1017 template = """Complete the following Lean 4 code:
1018
1019 ```lean4
1020 {lean_code}
1021 ```
1022
1023 You can make your own auxiliary corollaries and theorems to support the proof, instead of only completing the part
1024      ↪ with the sorry. Please output the entire program and not just the last part. Please output only the program
1025      ↪ and add no other comment, such that your answer is a compilable lean code. Make sure to reason enough to make
1026      ↪ your code correct.
1027
1028 """
1029
1030 informal_tmplate = """
1031 /--
1032 {informal}
1033 -/
1034 """
1035
1036 for problem in tables:
1037     lean_path = problem["file_name"]
1038     file_name = lean_path.split("/")[-1]
1039     out_path = f"{output_dir}/{file_name.removesuffix('.lean')}.txt"
1040     is_halting = True
1041     if "nonhalting" in file_name:
1042         is_halting = False
1043     informal_head = "Prove the following Turing Machine with the transition table " + (
1044         "halts." if is_halting else "never halts."
1045     )
1046     # print(out_path)
1047     lean_code = None
1048     with open(lean_path, "r") as f:
1049         lean_code = f.read().removeprefix("\n")
1050     assert lean_code != None
1051     informal = "/-- \n" + informal_head + "\n" + problem["table"] + "-/\n"
1052     with open(out_path, "w") as f:
1053         f.write(
1054             template.format(
1055                 lean_code=lean_code.replace(

```

```

1045         "theorem machine_", informal + "theorem machine_"
1046     ).replace("\n\ndef nth_cfg", "\ndef nth_cfg")
1047     .replace("\n\n/--", "\n/--")
1048     .replace("sorry\n", "sorry")
1049 )
1050 pass

```

E. More Success and Failure Cases

We demonstrate additional full record of evaluation, including problem statements in their Lean-Markdown pair and model responses for both successful and failed evaluation cases. The comprehensive record of model responses for cases can be found in the results folder of our dataset.

E.1. Success Case

E.1.1. DEEPSEEK-PROVER-V2-671B/BB-4STATE-CASE107-NONHALTING

This is a trivial non-halting case, meaning there is no "HALT" entry in the transition table. It can therefore be resolved via a single structural induction. Although the proof completes, warning messages appear, such as "aesop: failed to prove the goal after exhaustive search", "(try aesop) tactic does nothing", and "this tactic is never executed". These suggest the model does not fully reason about the proof and instead relies on automated tactics, hoping some may resolve the goal.

```

1066 import Mathlib.Computability.TuringMachine
1067
1068 inductive Γ
1069 | zero
1069 | one
1070 deriving DecidableEq
1071
1072 instance : Inhabited Γ := ⟨ Γ.zero ⟩
1073
1074 structure Stmt where
1075   move : Turing.Dir
1076   write : Γ
1077
1078 def Machine := Fin 4 → Γ → Option (Fin 4 → Stmt)
1079
1080 structure Cfg where
1081   q : Fin 4
1082   Tape : Turing.Tape Γ
1083
1084 def init (l : List Γ) : Cfg := ⟨0, by omega⟩, Turing.Tape.mk l
1085
1086 def step (M : Machine) : Cfg → Option Cfg :=
1087   fun ⟨q, T⟩ => (M q T.head).map fun ⟨q', a⟩ => ⟨q', (T.write a.write).move a.move⟩
1088
1089 def machine : Machine
1090 | ⟨000, _⟩, Γ.zero => some ⟨⟨001, by omega⟩, ⟨Turing.Dir.left, Γ.one⟩⟩
1091 | ⟨000, _⟩, Γ.one => some ⟨⟨001, by omega⟩, ⟨Turing.Dir.right, Γ.one⟩⟩
1092 | ⟨001, _⟩, Γ.zero => some ⟨⟨000, by omega⟩, ⟨Turing.Dir.right, Γ.one⟩⟩
1093 | ⟨001, _⟩, Γ.one => some ⟨⟨003, by omega⟩, ⟨Turing.Dir.left, Γ.one⟩⟩
1094 | ⟨002, _⟩, Γ.zero => some ⟨⟨001, by omega⟩, ⟨Turing.Dir.right, Γ.zero⟩⟩
1095 | ⟨002, _⟩, Γ.one => some ⟨⟨000, by omega⟩, ⟨Turing.Dir.right, Γ.zero⟩⟩
1096 | ⟨003, _⟩, Γ.zero => some ⟨⟨000, by omega⟩, ⟨Turing.Dir.left, Γ.zero⟩⟩
1097 | ⟨003, _⟩, Γ.one => some ⟨⟨001, by omega⟩, ⟨Turing.Dir.left, Γ.zero⟩⟩
1098 | ⟨_+4, _⟩, _ => by omega
1099
1100 def nth_cfg : (n : Nat) -> Option Cfg
1101 | 0 => init []
1102 | Nat.succ n => match (nth_cfg n) with
1103 | none => none
1104 | some cfg => step machine cfg
1105
1106 /--
1107 Prove the following Turing Machine with the transition table never halts.
1108 | State | Symbol | Next State | Move | Write |
1109

```



```

1100 |----|----|----|----|----|
1101 | 000 | zero | 001 | left | one |
1102 | 000 | one  | 001 | right| one |
1103 | 001 | zero | 000 | right| one |
1104 | 001 | one  | 003 | left | one |
1105 | 002 | zero | 001 | right| zero |
1106 | 002 | one  | 000 | right| zero |
1107 | 003 | zero | 000 | left | zero |
1108 | 003 | one  | 001 | left | zero |
1109 -/
1110 theorem machine_never_halts : ∀ n, (nth_cfg n).isSome := by
1111   intro n
1112   have h : ∀ n, (nth_cfg n).isSome := by
1113     intro n
1114     induction n with
1115     | zero => simp [nth_cfg, init]
1116     | succ n ih =>
1117       cases h₁ : nth_cfg n with
1118       | none => simp_all [nth_cfg, step, machine, Cfg, Stmt, Γ, Fin.forall_fin_succ, Fin.forall_fin_succ']
1119       | some cfg =>
1120         simp_all [nth_cfg, step, machine, Cfg, Stmt, Γ, Fin.forall_fin_succ, Fin.forall_fin_succ']
1121         <|>
1122         (try cases cfg <|> simp_all [nth_cfg, step, machine, Cfg, Stmt, Γ, Fin.forall_fin_succ, Fin.forall_fin_succ
1123           ↪'])
1124         <|>
1125         (try aesop)
1126         <|>
1127         (try omega)
1128         <|>
1129         (try aesop)
1130   exact h n

```

E.2. Failure Case

E.2.1. DEEPSEEK-PROVER-V2-671B/BB-2STATE-CASE102-NONHALTING

This is a case where the Turing machines state remains fixed at 000. A correct proof requires showing that the machine consistently executes the first rule of the transition table. The model attempts a single structural induction; however, in the inductive case *succ n ih*, it leaves two unsolved goals *False*. These cannot be resolved via the induction hypothesis, suggesting that the model fails to generalize inductive reasoning in this setting.

```

1133 import Mathlib.Computability.TuringMachine
1134
1135 inductive Γ
1136 | zero
1137 | one
1138   deriving DecidableEq
1139
1140 instance : Inhabited Γ := ⟨ Γ.zero ⟩
1141
1142 structure Stmt where
1143   move : Turing.Dir
1144   write : Γ
1145
1146 def Machine := Fin 2 → Γ → Option (Fin 2 → Stmt)
1147
1148 structure Cfg where
1149   q : Fin 2
1150   Tape : Turing.Tape Γ
1151
1152 def init (l : List Γ) : Cfg := ⟨ 0, by omega, Turing.Tape.mk 1 l ⟩
1153
1154 def step (M : Machine) : Cfg → Option Cfg :=
1155   fun ⟨ q, T ⟩ => (M q T.head).map fun ⟨ q', a ⟩ => ⟨ q', (T.write a.write).move a.move ⟩
1156
1157 def machine : Machine
1158 | ⟨ 000, _ ⟩, Γ.zero => some ⟨ ⟨ 000, by omega ⟩, ⟨ Turing.Dir.left, Γ.zero ⟩ ⟩
1159 | ⟨ 000, _ ⟩, Γ.one  => none
1160 | ⟨ 001, _ ⟩, Γ.zero => none

```

```

1155 | ⟨001, _⟩, Γ.one => some ⟨⟨000, by omega⟩, ⟨Turing.Dir.left, Γ.one⟩⟩
1156 | ⟨_+2, _⟩, _ => by omega
1157
1158 def nth_cfg : (n : Nat) -> Option Cfg
1159 | 0 => init []
1160 | Nat.succ n => match (nth_cfg n) with
1161 | none => none
1162 | some cfg => step machine cfg
1163
1164 /--
1165 Prove the following Turing Machine with the transition table never halts.
1166 | State | Symbol | Next State | Move | Write |
1167 |-----|-----|-----|-----|-----|
1168 | 000 | zero | 000 | left | zero |
1169 | 000 | one | HALT | - | - |
1170 | 001 | zero | HALT | - | - |
1171 | 001 | one | 000 | left | one |
1172 -/
1173
1174 theorem machine_never_halts : ∀ n, (nth_cfg n).isSome := by
1175   intro n
1176   have h : ∀ n, (nth_cfg n).isSome := by
1177     intro n
1178     induction n with
1179     | zero =>
1180       simp [nth_cfg, init]
1181     | succ n ih =>
1182       simp_all [nth_cfg, step, machine, Cfg.q, Cfg.Tape, Turing.Tape.move, Option.isSome]
1183   <|> aesop
1184   exact h n

```

E.2.2. GOEDEL-PROVER/BB-2STATE-CASE102-NONHALTING

The model produced proof of a completely unrelated theorem about sequences. This seems to be a hallucination, where the model matched a common pattern like *theorem...by sorry* and filled in a generic proof without understanding the actual task. This shows a key limitation of small provers: Without symbolic grounding or verification, they can only match surface-level patterns rather than truly reasoning.

```

1184 import Mathlib
1185 import Aesop
1186
1187 set_option maxHeartbeats 0
1188
1189 open BigOperators Real Nat Topology Rat
1190
1191 /-- Given a sequence { a_n } with a_1 = 1 and a_{n+1} = 2 a_n + 1,
1192 prove that a_n = 2^n - 1 for all n ≥ 1 -/
1193
1194 theorem sequence_formula {a : ℕ → ℕ} (h1 : a 1 = 1)
1195 (h_rec : ∀ n : ℕ, n ≥ 1 → a (n + 1) = 2 * a n + 1) :
1196   ∀ n : ℕ, n ≥ 1 → a n = 2^n - 1 := by
1197   /-
1198   We need to prove that for a sequence \{ a_n \} defined by \{ a_1 = 1 \} and \{ a_{n+1} = 2a_n + 1 \}, the
1199   ↪ general term \{ a_n \} is given by \{ a_n = 2^n - 1 \} for all \{ n ≥ 1 \}.
1200   1. **Base Case**: For \{ n = 1 \}, we have \{ a_1 = 1 \}. This matches \{ 2^1 - 1 = 1 \}, so the base case holds.
1201   2. **Inductive Step**: Assume \{ a_n = 2^n - 1 \} holds for some \{ n ≥ 1 \}. We need to show it holds for \{ n + 1 \}.
1202   ↪ + 1 \}.
1203   - By the recurrence relation, \{ a_{n+1} = 2a_n + 1 \}.
1204   - Using the inductive hypothesis \{ a_n = 2^n - 1 \}, we get:
1205     \[
1206     a_{n+1} = 2(2^n - 1) + 1 = 2^{n+1} - 2 + 1 = 2^{n+1} - 1
1207     \]
1208   - Thus, \{ a_{n+1} = 2^{n+1} - 1 \}, completing the inductive step.
1209   By induction, \{ a_n = 2^n - 1 \} holds for all \{ n ≥ 1 \}.
1210   -/
1211   intro n hn
1212   induction' hn with n hn
1213   -- Base case: n = 1
1214   -- We know a 1 = 1, which matches 2^1 - 1 = 1
1215   simp_all
1216   -- Inductive step: Assume a n = 2^n - 1, prove a (n + 1) = 2^(n + 1) - 1
1217   cases n with

```

```

1210 | zero =>
1211 -- When n = 0, the induction hypothesis does not apply, so we directly use the initial condition
1212 simp_all
1213 | succ n =>
1214 -- For n = succ n, use the recurrence relation and the inductive hypothesis
1215 simp_all [Nat.succ_eq_add_one, Nat.pow_succ, Nat.mul_succ]
1216 -- Simplify the expression to match the form 2^(n + 1) - 1
1217 ring
1218 </> omega

```

E.2.3. KIMINI-PROVER-7B/BB-4STATE-CASE122-NONHALTING

The model just outputs *sorry*, which means it didnt complete the proof. This might be because it couldnt figure out how to reason about the problem, or because it has seen *sorry* used as a placeholder in Lean and assumed that was enough to count as a valid answer.

```

1223 import Mathlib.Computability.TuringMachine
1224
1225 inductive Γ
1226 | zero
1227 | one
1228 deriving DecidableEq
1229
1230 instance : Inhabited Γ := ⟨ Γ.zero ⟩
1231
1232 structure Statement where
1233   move : Turing Dir
1234   write : Γ
1235
1236 def Machine := Fin 4 → Γ → Option (Fin 4 → Statement)
1237
1238 structure Cfg where
1239   q : Fin 4
1240   Tape : Turing Tape Γ
1241
1242 def init (l : List Γ) : Cfg := ⟨⟨0, by omega⟩, ⟨Turing Dir.right, Γ.zero⟩⟩
1243
1244 def step (M : Machine) : Cfg → Option Cfg :=
1245   fun ⟨q, T⟩ => (M q T.head).map fun ⟨q', a⟩ => ⟨q', (T.map a).move a⟩
1246
1247 def machine : Machine
1248 | ⟨000, _⟩, Γ.zero => some ⟨⟨001, by omega⟩, ⟨Turing Dir.right, Γ.zero⟩⟩
1249 | ⟨000, _⟩, Γ.one => some ⟨⟨003, by omega⟩, ⟨Turing Dir.left, Γ.zero⟩⟩
1250 | ⟨001, _⟩, Γ.zero => some ⟨⟨003, by omega⟩, ⟨Turing Dir.left, Γ.one⟩⟩
1251 | ⟨001, _⟩, Γ.one => some ⟨⟨002, by omega⟩, ⟨Turing Dir.right, Γ.zero⟩⟩
1252 | ⟨002, _⟩, Γ.zero => some ⟨⟨003, by omega⟩, ⟨Turing Dir.left, Γ.zero⟩⟩
1253 | ⟨002, _⟩, Γ.one => some ⟨⟨001, by omega⟩, ⟨Turing Dir.left, Γ.one⟩⟩
1254 | ⟨003, _⟩, Γ.zero => some ⟨⟨001, by omega⟩, ⟨Turing Dir.left, Γ.one⟩⟩
1255 | ⟨003, _⟩, Γ.one => some ⟨⟨000, by omega⟩, ⟨Turing Dir.right, Γ.zero⟩⟩
1256 | ⟨_+4, _⟩, _ => by omega
1257
1258 def nth configurations (n : Nat) : Option Cfg :=
1259 | 0 => init []
1260 | Nat.succ n => match (nth configurations n) with
1261 | none => none
1262 | some cfg => step machine cfg
1263
1264 /--
1265 Prove the following Turing Machine with the transition table never halts.
1266 | State | Symbol | Next State | Move | Write |
1267 |-----|-----|-----|-----|-----|
1268 | 000 | zero | 001 | right | zero |
1269 | 000 | one | 003 | left | zero |
1270 | 001 | zero | 003 | left | one |
1271 | 001 | one | 002 | right | zero |
1272 | 002 | zero | 003 | left | zero |
1273 | 002 | one | 001 | left | one |
1274 | 003 | zero | 001 | left | one |
1275 | 003 | one | 000 | right | zero |
1276 -/

```

```

1265 theorem machine_never_halts :  $\forall$  n, (nth configurations n).isSome := by
1266   sorry

```

F. Customized Lemma Library.

```

1270 import Mathlib.Tactic.NthRewrite
1271
1272 @[simp]
1273 theorem bv32_and_not_self (x : BitVec 32) :
1274   x &&& ~x = 0 := by
1275     simp
1276
1277 @[simp]
1278 theorem bv32_not_not (x : BitVec 32) :
1279   ~~~x = x := by
1280     simp
1281
1282 @[simp]
1283 theorem bv32_or_not_self (x : BitVec 32) :
1284   x ||| ~x = BitVec.allOnes 32 := by
1285     simp
1286
1287 @[simp]
1288 theorem bv32_not_or_self (x : BitVec 32) :
1289   ~x ||| x = BitVec.allOnes 32 := by
1290     simp
1291
1292 theorem bv32_neg_mul (x y : BitVec 32) :
1293   -x * y = -(x * y) := by
1294     simp
1295
1296 theorem bv32_not_and (x y : BitVec 32) :
1297   ~x &&& y = ~x ||| ~y := by
1298     rw [BitVec.not_and]
1299
1300 theorem bv32_not_or (x y : BitVec 32) :
1301   ~x ||| y = ~x &&& ~y := by
1302     rw [BitVec.not_or]
1303
1304 theorem bv32_not_xor_eq_or (x y : (BitVec 32)) :
1305   ~x ^ y = (~x &&& ~y) ||| (x &&& y) := by
1306     ext i
1307     simp
1308     cases h1 : x[i] <> cases h2 : y[i]
1309     simp
1310     simp
1311     simp
1312     simp
1313
1314 theorem bv32_xor_eq_or (x y : (BitVec 32)) :
1315   x ^ y = (~x &&& y) ||| (x &&& ~y) := by
1316     ext i
1317     simp
1318     cases h1 : x[i] <> cases h2 : y[i]
1319     simp
1320     simp
1321     simp
1322     simp
1323
1324 theorem bv32_x_distr (x y : BitVec 32) :
1325   x &&& (y ||| ~y) = (x &&& y) ||| (x &&& ~y) := by
1326     ext i
1327     simp
1328     simp [Bool.and_or_distrib_left]
1329
1330 theorem bv32_y_distr (x y : BitVec 32) :
1331   (x ||| ~x) &&& y = (x &&& y) ||| (~x &&& y) := by
1332     ext i
1333     simp

```

```

1320   simp [    Bool.and_or_distrib_right]
1321
1322   theorem bv32_add_assoc (x y z : BitVec 32) :
1323     x + y + z = x + (y + z) := by
1324     rw [BitVec.add_assoc]
1325
1326   theorem bv32_add_comm (x y : BitVec 32) :
1327     x + y = y + x := by
1328     rw [BitVec.add_comm]
1329
1330   theorem bv32_add_neg_eq_sub {x y : BitVec 32} :
1331     x + -y = x - y := by
1332     rw [BitVec.add_neg_eq_sub]
1333
1334   theorem bv32_mul_comm (x y : BitVec 32) :
1335     x * y = y * x := by
1336     rw [BitVec.mul_comm]
1337
1338   theorem bv32_var_mul_comm (x y z : BitVec 32) :
1339     (x &&& y) * z = z * (x &&& y) := by
1340     rw [BitVec.mul_comm]
1341
1342   theorem bv32_mul_add (x y z : BitVec 32) :
1343     x * (y + z) = x * y + x * z := by
1344     rw [BitVec.mul_add]
1345
1346   theorem bv32_neg_eq_mul (x : BitVec 32) :
1347     -x = x * (BitVec.allOnes 32) := by
1348     rw [    BitVec.neg_one_eq_allOnes]
1349     rw [BitVec.mul_neg]
1350     rw [BitVec.mul_one]
1351
1352   theorem bv32_add_mul_one (x y : BitVec 32) :
1353     x + x * y = x * (1#32 + y) := by
1354     rw [BitVec.mul_add]
1355     rw [BitVec.mul_one]
1356
1357   /--
1358   1: x &&& y
1359   2: ~~~x &&& y
1360   3: x &&& ~~~y
1361   4: ~~~x &&& ~~~y
1362   -/
1363   theorem bv32_or_eq_add12 (x y : BitVec 32) :
1364     (x &&& y) ||| (~~~x &&& y) = (x &&& y) + (~~~x &&& y) := by
1365     apply Eq.symm
1366     apply BitVec.add_eq_or_of_and_eq_zero
1367     simp [    BitVec.and_assoc]
1368     simp [BitVec.and_comm _ (~~~x)]
1369     simp [    BitVec.and_assoc]
1370
1371   theorem bv32_or_eq_add13 (x y : BitVec 32) :
1372     (x &&& y) ||| (x &&& ~~~y) = (x &&& y) + (x &&& ~~~y) := by
1373     apply Eq.symm
1374     apply BitVec.add_eq_or_of_and_eq_zero
1375     simp [    BitVec.and_assoc]
1376     simp [BitVec.and_comm _ x]
1377     simp [BitVec.and_assoc]
1378
1379   theorem bv32_or_eq_add14 (x y : BitVec 32) :
1380     (x &&& y) ||| (~~~x &&& ~~~y) = (x &&& y) + (~~~x &&& ~~~y) := by
1381     apply Eq.symm
1382     apply BitVec.add_eq_or_of_and_eq_zero
1383     simp [    BitVec.and_assoc]
1384     simp [BitVec.and_comm _ (~~~x)]
1385     simp [    BitVec.and_assoc]
1386
1387   theorem bv32_or_eq_add21 (x y : BitVec 32) :
1388     (~~~x &&& y) ||| (x &&& y) = (~~~x &&& y) + (x &&& y) := by
1389     apply Eq.symm
1390     apply BitVec.add_eq_or_of_and_eq_zero

```

```

1375 simp [ BitVec.and_assoc]
1376 simp [BitVec.and_comm _ x]
1377 simp [ BitVec.and_assoc]
1378 theorem bv32_or_eq_add23 (x y : BitVec 32) :
1379   (~~x &&& y) ||| (x &&& ~y) = (~~x &&& y) + (x &&& ~y) := by
1380   apply Eq.symm
1381   apply BitVec.add_eq_or_of_and_eq_zero
1382   simp [ BitVec.and_assoc]
1383   simp [BitVec.and_comm _ x]
1384   simp [ BitVec.and_assoc]
1385 theorem bv32_or_eq_add31 (x y : BitVec 32) :
1386   (x &&& ~y) ||| (x &&& y) = (x &&& ~y) + (x &&& y) := by
1387   apply Eq.symm
1388   apply BitVec.add_eq_or_of_and_eq_zero
1389   simp [ BitVec.and_assoc]
1390   simp [BitVec.and_comm _ x]
1391   simp [ BitVec.and_assoc]
1392 theorem bv32_or_eq_add32 (x y : BitVec 32) :
1393   (x &&& ~y) ||| (~~x &&& y) = (x &&& ~y) + (~~x &&& y) := by
1394   apply Eq.symm
1395   apply BitVec.add_eq_or_of_and_eq_zero
1396   simp [ BitVec.and_assoc]
1397   simp [BitVec.and_comm _ (~~x)]
1398   simp [ BitVec.and_assoc]
1399 theorem bv32_or_eq_add41 (x y : BitVec 32) :
1400   (~~x &&& ~y) ||| (x &&& y) = (~~x &&& ~y) + (x &&& y) := by
1401   apply Eq.symm
1402   apply BitVec.add_eq_or_of_and_eq_zero
1403   simp [ BitVec.and_assoc]
1404   simp [BitVec.and_comm _ x]
1405   simp [ BitVec.and_assoc]
1406 theorem bv32_or_eq_add_three (x y : BitVec 32) :
1407   (x ||| y) = (x &&& ~y) + (x &&& y) + (~~x &&& y) := by
1408   nth_rw 1 [bv32_y_distr x y]
1409   nth_rw 1 [bv32_x_distr x y]
1410   simp [ BitVec.or_assoc]
1411   simp [BitVec.or_comm _ (x &&& y)]
1412   simp [ BitVec.or_assoc]
1413   rw [BitVec.or_comm (x &&& y)]
1414   apply Eq.symm
1415   rw [BitVec.add_eq_or_of_and_eq_zero]
1416   rw [BitVec.add_eq_or_of_and_eq_zero]
1417   simp [ BitVec.and_assoc]
1418   simp [BitVec.and_comm _ x]
1419   simp [BitVec.and_assoc]
1420   rw [BitVec.add_comm]
1421   rw [ bv32_or_eq_add13]
1422   rw [ bv32_x_distr x y]
1423   simp [ BitVec.and_assoc]
1424 theorem bv32_sum_all (x y : BitVec 32) :
1425   (~~x &&& ~y) + (~~x &&& y) + (x &&& y) + (x &&& ~y) = BitVec.allOnes 32 := by
1426   simp [BitVec.add_comm _ (~~x &&& y)]
1427   simp [BitVec.add_comm _ (x &&& _)]
1428   simp [ BitVec.add_assoc]
1429   rw [BitVec.add_eq_or_of_and_eq_zero]
1430   rw [ bv32_or_eq_add_three x y]
1431   nth_rw 1 [bv32_x_distr x y]
1432   simp [BitVec.or_comm _ y]
1433   nth_rw 1 [bv32_y_distr x y]
1434   simp [ BitVec.or_assoc]
1435   simp [BitVec.or_comm _ (x &&& y)]
1436   simp [ BitVec.or_assoc]
1437   simp [BitVec.or_comm _ (~~x &&& _)]
1438   simp [ BitVec.or_assoc]
1439   simp [BitVec.or_comm _ (~~x &&& y)]

```

```

1430   rw [    bv32_x_distr (~~~x) y]
1431   rw [BitVec.or_assoc]
1432   rw [    bv32_x_distr x y]
1433   simp
1434   rw [    bv32_or_eq_add_three x y]
1435   simp [    BitVec.not_or]
1436
1437   theorem bv32_self_eq_neg_mul (x: BitVec 32):
1438     x = -x * (BitVec.allOnes 32) := by
1439       rw [BitVec.neg_mul]
1440       rw [BitVec.mul_comm]
1441       rw [    BitVec.neg_mul]
1442       simp [    BitVec.neg_one_eq_allOnes]
1443
1444   theorem bv32_not_self_and_not (x y : BitVec 32) :
1445     ~~~(x &&& ~~~x) = (~~~x &&& ~~~y) + (~~~x &&& y) + (x &&& y) + (x &&& ~~~y) := by
1446       rw [bv32_not_and]
1447       rw [BitVec.not_not]
1448       rw [bv32_not_or_self]
1449       rw [bv32_sum_all]

```