

# A LANGUAGE-AGENT APPROACH TO FORMAL THEOREM-PROVING

Amitayush Thakur, Yeming Wen & Swarat Chaudhuri

The University of Texas at Austin

{amitayush, ywen}@utexas.edu, swarat@cs.utexas.edu

## ABSTRACT

Language agents, which use a large language model (LLM) capable of in-context learning to interact with an external environment, have recently emerged as a promising approach to control tasks. We present the first language-agent approach to formal theorem-proving. Our method, COPRA, uses a high-capacity, black-box LLM (GPT-4) as part of a policy for a stateful backtracking search. During the search, the policy can select proof tactics and retrieve lemmas and definitions from an external database. Each selected tactic is executed in the underlying proof framework, and the execution feedback is used to build the prompt for the next policy invocation. The search also tracks selected information from its history and uses it to reduce hallucinations and unnecessary LLM queries.

We evaluate COPRA on the `miniF2F` benchmark for Lean and a set of Coq tasks from the CompCert project. On these benchmarks, COPRA is significantly better than one-shot invocations of GPT-4, as well as state-of-the-art models fine-tuned on proof data, at finding correct proofs quickly.

## 1 INTRODUCTION

Automatically proving formal theorems (Newell et al., 1957) is a longstanding challenge in computer science. Autoregressive language models (Polu & Sutskever, 2020; Han et al., 2021; Yang et al., 2023) have recently emerged as an effective approach to this problem. Such models are trained on proofs written in frameworks like Coq (Huet et al., 1997) or Lean (de Moura et al., 2015), which allows proof goals to be iteratively simplified using a set of *tactics*. Theorem-proving then amounts to generating a sequence of tactics that iteratively “discharges” a given proof goal.

A weakness of this method is that it does not model the *interaction* between the model and the underlying proof framework. The application of a tactic is an *action* that changes the state of the proof and the interpretation of future tactics. By ignoring these game-like dynamics, autoregressive models miss out on a valuable source of feedback and end up being more susceptible to hallucinations.

In this paper, we show that the nascent paradigm of *large-language-model (LLM) agents* (Yao et al., 2022; Wang et al., 2023; Shinn et al., 2023) can help address this weakness. Here, one uses an LLM as a *agent* that interacts with an external environment. Information gathered through interaction is used to update the LLM’s prompt, eliciting new agent behavior because of in-context learning.

Our approach, called COPRA<sup>1</sup> (Figure 1), uses an off-the-shelf, high-capacity LLM (GPT-4 (OpenAI, 2023)) as part of a policy in that interacts with a proof environment like Coq or Lean. At each time step, the policy consumes a textual prompt and chooses to use an available tactic, or backtrack, or retrieve relevant lemmas and definitions from an external corpus. When the policy selects a tactic, we “execute” it using the underlying proof assistant. The feedback from the execution is used to construct a new prompt for the policy, and the process repeats.

COPRA goes beyond prior language-agent methods in using domain knowledge and information from the search history to use LLM queries frugally. When tactics fail, the policy records this information and uses it to avoid future failures. The policy also has access to a symbolic procedure

<sup>1</sup>COPRA is an acronym for “In-context **P**rover **A**gent”.

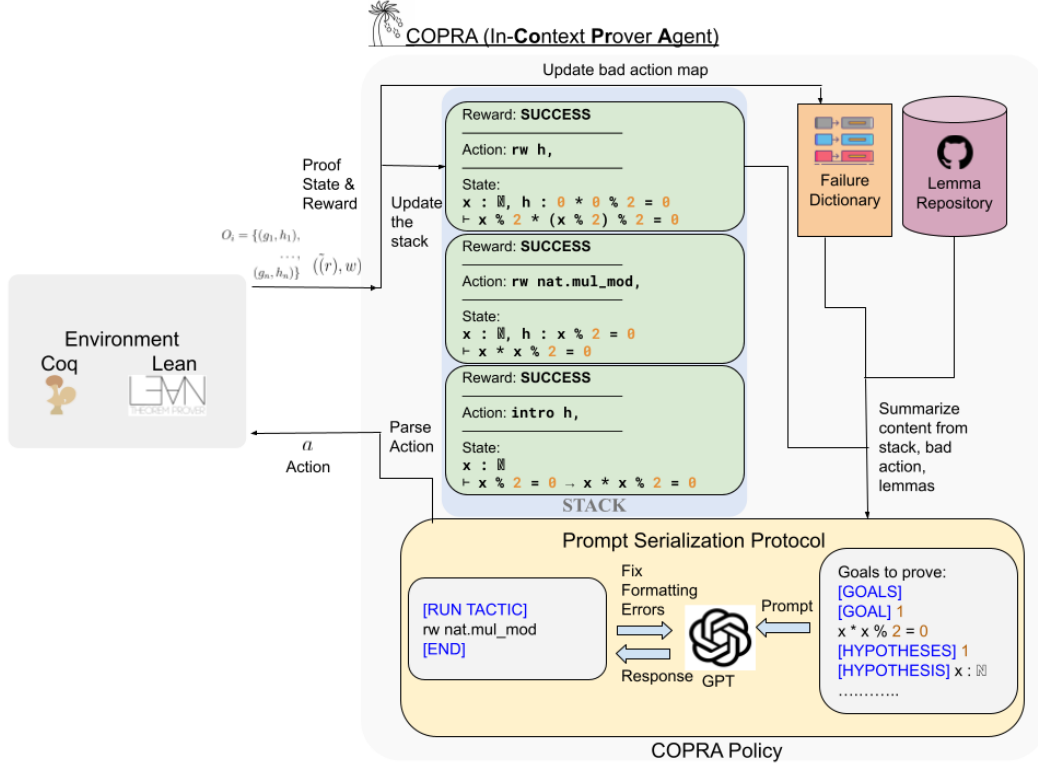


Figure 1: An overview of COPRA. The system implements a policy that interacts with a proof environment (Coq or Lean). Internally, a COPRA policy consists of an LLM (GPT-4), a stack-based backtracking search, a retrieval mechanism, a dictionary tracking past failures, and a prompt serialization protocol that constructs LLM prompts using the stack and environment feedback and parse LLM outputs into actions.

that checks if one goal is “simpler” than another. A tactic is only used when it simplifies the agent’s proof obligations (ruling out, among other things, cyclic tactic sequences).

We have integrated COPRA with both the Coq and the Lean environments. We evaluate the system using the *miniF2F* (Zheng et al., 2021) benchmark for competition-level mathematical reasoning in Lean and a set of Coq proof tasks (Sanchez-Stern et al., 2020) from the CompCert (Leroy, 2009) project on verified compilation. Using a new metric called *prove-at-k-inferences*, we show that COPRA can converge to correct proofs faster than competing approaches, including the state-of-the-art models (Yang et al., 2023; Sanchez-Stern et al., 2020) trained on formal proof data. We also show that when COPRA fails, it fails quicker than the baseline methods.

To summarize our contributions, we offer: (i) The first approach to formal theorem-proving that leverages LLMs while also modeling interactions between the model and the underlying proof framework; (ii) the first language agent, from any domain, to integrate LLM policies with a search that minimizes LLM queries and hallucinations by tracking domain-specific information from the past; and (iii) an implementation of COPRA that interacts with the Coq and Lean proof environments, and an evaluation on two domains — mathematics competition problems and formal verification — that shows COPRA to find proofs faster than competing approaches.

## 2 THEOREM-PROVING AS A CONTROL PROBLEM

### 2.1 BACKGROUND ON THEOREM-PROVING

A *formal proof* starts with a set of unmet *obligations* stated in a formal language and applies a sequence of *proof tactics* to progressively eliminate these obligations. Each obligation  $o$  consists of a *goal*  $g$  and a *hypothesis*  $h$ . The goal  $g$  consists of the propositions that need to be proved in order

to meet  $o$ ; the hypothesis  $h$  captures assumptions that can be made in the proof of  $g$ . The prover’s long-term objective is to reduce the obligations to the empty set.

We illustrate this process with the example in Figure 2-(a). This example shows a Lean (de Moura et al., 2015) proof, automatically generated using COPRA, of a basic theorem about modular arithmetic. The proof first applies the `intro` tactic, which changes a goal  $P \rightarrow Q$  to a hypothesis  $P$  and a goal  $Q$ . Next, it applies the `rw` (rewrite) tactic, which gives a way to apply substitutions to goals and hypotheses, several times. It ends with the application of the `refl` (reflexivity) tactic, which eliminates goals that say that a value is equal to itself.

Existing LLM-based approaches to automatic theorem-proving view such proofs as purely syntactic artifacts. However, the rigorous semantics of proofs can be difficult to learn using such an approach, leading to the generation of incorrect proofs. Figure 2-(c) shows a GPT-4-generated incorrect proof of our theorem.

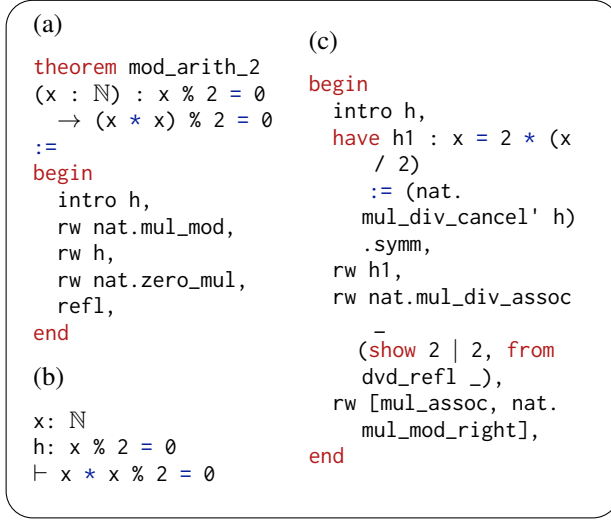


Figure 2: (a) A Lean theorem and a correct proof found by COPRA. (b) Proof state after the first tactic. (c) An incorrect proof generated by GPT-4.

## 2.2 A MARKOV DECISION PROCESS FORMULATION

By contrast, COPRA is based on a view of automatic theorem-proving as a *control problem*. Like prior work on reinforcement learning (RL) for proof synthesis (Wu et al., 2021), we view a theorem-prover as a *policy* that interacts with a stateful proof environment (e.g., Lean) and model the interaction between the policy and the environment as a deterministic Markov Decision Process (MDP). We depart from prior RL-based work for theorem-proving by imposing a partial order on MDP states, allowing rewards to have a textual component, and allowing history-dependent policies.

Now we describe the different components of our *proof MDP*.

**States.** As before, let an *obligation* be a pair  $(g, h)$ , where  $g$  is a goal and  $h$  a hypothesis. A *state* of the MDP is either a special symbol called *error* or a set  $O = \{o_1, \dots, o_k\}$  of obligations  $o_i$ .

The MDP has a unique *initial state*  $o_{in}$  with a single obligation  $(g_{in}, h_{in})$ , where the goal  $g_{in}$  and the hypothesis  $h_{in}$  are extracted from the user-provided theorem that we are trying to prove. Its unique *final state* QED is the empty obligation set.

Following Sanchez-Stern et al. (2020), we define a partial order  $\sqsubseteq$  over states that defines when a state is “at least as hard” than another and use it to avoid actions that do not lead to progress in the proof. Formally, for states  $O_1$  and  $O_2$  with  $O_1 \neq error$  and  $O_2 \neq error$ ,  $O_1 \sqsubseteq O_2$  iff

$$\forall o_i = (g_i, h_i) \in O_1. \exists o_k = (g_k, h_k) \in O_2. g_k = g_i \wedge (h_k \rightarrow h_i).$$

Intuitively,  $O_1 \sqsubseteq O_2$  if for every obligation in  $O_1$ , there is a stronger obligation in  $O_2$ . We assume we have an efficient symbolic procedure that can check this relationship for any pair of states. The procedure is *sound*, meaning that if it reports  $O_1 \sqsubseteq O_2$ , the relationship actually holds. However, it is *incomplete*, i.e., it may not detect all relationships of the form  $O_1 \sqsubseteq O_2$ .

**Actions and Transitions.** The actions in our MDP are the proof environment’s *tactics*.

The transition function  $T(O, a)$  determines the result of applying an action  $a$  to a state  $O$ . When  $a$  is a tactic, we assume the underlying proof environment to return a state  $O'$  that results from applying  $a$  to  $O$ . If  $a$  is a “bad” tactic, then  $O'$  equals *error*; otherwise,  $O'$  is a new set of obligations. We assume that our agent can evaluate  $T(O, a)$  for any state  $O$  and action  $a$ . While this assumption is unacceptable in many MDP problems, it is reasonable in the theorem-proving setting.

**Rewards.** As usual, we assume a *reward function*  $R(O, a)$  that evaluates an action  $a$  at a state  $O$ . Historically, such functions are scalar-valued; however, because we use LLMs as policies, we allow rewards to also include rich textual feedback from the proof environment. Concretely, we consider rewards of the form  $R(O, a) = (\tilde{r}, w)$ , where:

(1)  $\tilde{r}$  is a very high positive value if  $T(O, a) = \text{QED}$ , a negative value if  $T(O, a) = \text{error}$ , and 0 otherwise, and (2)  $w$  is the feedback from the proof environment when  $a$  is executed from  $O$ .

**Histories and Policies.** A *history* of length  $N$  is a sequence

$$h = \langle (O_0, a_0, O'_0, r_0), (O_1, a_1, O'_1, r_1), \dots, (O_{N-1}, a_{N-1}, O'_N, r_N) \rangle$$

such that  $O_0 = O_{in}$  and for all  $i$ ,  $r_i = R(O_i, a_i)$  and  $O'_i = T(O_i, a_i)$ . Intuitively, a history records the interactions between the prover agent and the proof environment up to a point of time. We denote by  $h_i$  the  $i$ -th prefix of  $h$ . For example,  $h_0 = \langle \rangle$ ,  $h_1 = \langle (O_0, a_0, O'_0, r_0) \rangle$ , and so on.

A *policy* is a probabilistic function  $\pi$  that maps histories to distributions over pairs  $(O, a)$ , where  $O$  is a state and  $a$  is an action. Intuitively, at each point, the policy determines the next query to make to the proof environment.

A policy can have an internal state as well as access to external knowledge (specifically, a lemma database). A *trajectory* of a policy  $\pi$  is a history  $h$  as above such that for each  $i$ ,

$$\Pr[\pi(h_i) = (O_i, a_i)] > 0.$$

Letting each  $r_i = (\tilde{r}_i, w_i)$ , the *scalar reward* from a trajectory is simply the average  $\frac{1}{N} \sum_i \tilde{r}_i$ . We define the *aggregate (scalar) reward* of  $\pi$  as the expected scalar reward from trajectories sampled from  $\pi$ .

**Language Agents.** Given our setup, one can naturally pose the problem of reinforcement-learning a policy with optimal aggregate reward. In this paper, we do not take on this problem. Instead, we consider a fixed policy — a wrapper around a pretrained LLM (GPT-4) that can learn in-context — and show that this policy can achieve a high reward. It is this policy that defines our *language agent*.

### 3 THE COPRA AGENT

A COPRA policy has access to an LLM (in practice, GPT-4) and performs a depth-first search. During the search, it records information about failed actions. It also uses the  $\sqsubseteq$  relation over states to checks that it is making progress on the proof.

Figure 3 shows pseudocode for such a policy. The policy maintains a stack of MDP states and a “failure dictionary”  $Bad$  that maps a state to a set of actions that are known to be “unproductive” at the state. At each search step, the algorithm pushes the current state on the stack and retrieves external lemmas and definitions relevant to the state. After this, it repeatedly serializes the stack and  $Bad(O)$  into a prompt and feeds it to the LLM. The LLM’s output is parsed into an action, and the agent executes it in the environment.

One outcome of the action could be that the agent arrives at QED. Alternatively, the new state could be an error or represent obligations that are at least as hard

```

COPRA( $O$ )
1  PUSH( $st, O$ )
2   $\rho \leftarrow \text{RETRIEVE}(O)$ 
3  for  $j \leftarrow 1$  to  $k$ 
4      do  $p \leftarrow \text{PROMPTIFY}(st, Bad(O), \rho, r)$ 
5           $a \sim \text{PARSEACTION}(\text{LLM}(p))$ 
6           $O' \leftarrow T(O, a), r \leftarrow R(O, a)$ 
7          if  $O' = \text{QED}$ 
8              then terminate successfully
9          else if  $O' = \text{error}$  or
               $\exists O'' \in st. O'' \sqsubseteq O'$ 
10             then add  $a$  to  $Bad(O)$ 
11             else COPRA( $O'$ )
12  POP( $st$ )

```

Figure 3: The search procedure in COPRA.  $T$  is the environment’s transition function and  $R$  is the reward function.  $st$  is a stack, initialized to be empty.  $Bad(O)$  is a set of actions, initialized to  $\emptyset$ , that are known to be bad at  $O$ . LLM is an LLM, PROMPTIFY generates a prompt, PARSEACTION parses the output of the LLM into an action (repeatedly querying the LLM in case there are formatting errors in its output), and RETRIEVE gathers relevant lemmas and definitions from an external source. The procedure is initially called with argument  $O_{in}$ .

Agent Prompt (State)	Goals to prove: [GOALS] [GOAL] 1 $x * x \% 2 = 0$ [HYPOTHESES] 1 [HYPOTHESIS] $x : \mathbb{N}$ [HYPOTHESIS] $h : x \% 2 = 0$			Goals to prove: [GOALS] [GOAL] 1 $x \% 2 * (x \% 2) \% 2 = 0$ [HYPOTHESES] 1 [HYPOTHESIS] $x : \mathbb{N}$ [HYPOTHESIS] $h : x \% 2 = 0$
Agent Prompt (Stack)	[LAST STEP] intro h,	[STEPS] [STEP] intro h,  [LAST STEP] rw h,	[STEPS] [STEP] intro h,  [INCORRECT STEPS] [STEP] rw h,  [LAST STEP] apply nat.mul_mod_right,	[STEPS] [STEP] intro h,  [LAST STEP] rw nat.mul_mod,
Agent Prompt (Reward)	[SUCCESS] [END]	[ERROR MESSAGE] Got error in 'rw h': error: rewrite tactic failed, did not find instance of the pattern in the target expression $x \% 2$ [END]	[ERROR MESSAGE] Got error in 'apply nat.mul_mod_right': error: invalid apply tactic, failed to unify $x * x \% 2 = 0$ with $?m\_1 * ?m\_2 \% ?m\_1 = 0$ [END]	[SUCCESS] [END]
↑ Requests # ↓ Response #	Seq # 1	Seq # 2	Seq # 3	Seq # 4
LLM Response	[RUN TACTIC] rw h, [END]	[RUN TACTIC] apply nat.mul_mod_right, [END]	[RUN TACTIC] rw nat.mul_mod, [END]	[RUN TACTIC] rw h, [END]

Figure 4: The prompt serialization protocol. We highlight the different parts of the prompts to show how we use the state stack and the textual reward from the environment.

as what is currently on the stack (for example, this could be because of a cycle in a tactic). In this case, the agent rejects the new state. Otherwise, it recursively continues the proof from the new state. After issuing a few queries to the LLM, the agent backtracks.

**Prompt Serialization Protocol.** The routines PROMPTIFY and PARSEACTION together constitute the *prompt serialization protocol* and are critical to the success of the policy. Now we elaborate on these procedures.

PROMPTIFY carefully places the different pieces of information relevant to the proof in the prompt. It also includes logic for trimming this information to fit the most relevant parts in the LLM’s context window. Every prompt has two parts: the “system prompt” and the “agent prompt”.

The agent prompts are synthetically generated using a context-free grammar and contain information about the state stack (including the current proof state), the textual reward for the previous action, and the set of actions we know to avoid at the current proof state.

The system prompt describes the rules of engagement for the LLM. It contains a grammar (distinct from the one for agent prompts) that we expect the LLMs to follow when it proposes a course of action. The grammar carefully incorporates cases when the response is incomplete because of the LLM’s token limits. We parse partial responses to extract the next action using the PARSEACTION routine. PARSEACTION also identifies formatting errors (if any) in the LLM’s responses, possibly communicating with the LLM multiple times until these errors are resolved.

**Example.** Figure 4 illustrates the prompt serialization protocol at work during the generation of the proof in Figure 2-(b). Seq #1-#4 represent distinct invocations of the LLM. In each invocation, PROMPTIFY first generates the “agent prompt,” which consists of three parts. The first part (“state”) is simply a serialization of the current proof state. The second (“stack”) incorporates information about previous actions as well as the bad actions for the current proof state. The third (“reward”) encodes the feedback from the environment regarding the success or failure of the last action. The response of the LLM to this prompt is then translated into an action using PARSEACTION. This action is then executed on the theorem prover.

## 4 EVALUATION

Our findings about COPRA are that: (i) the approach can find proofs significantly quicker than the state-of-the-art finetuning-based baselines, both in terms of number of LLM queries and wall-clock time; (ii) in problems where all current methods fail, COPRA fails faster; (iii) the use of GPT-4, as opposed to GPT-3.5, within the agent is essential for success; and (iv) backtracking significantly improves the system’s performance on harder problems. Now we elaborate on our experimental methodology and these results.

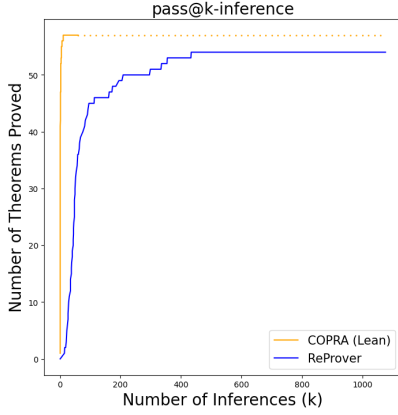


Figure 5: COPRA vs. REPROVER on the miniF2F benchmark

244 Lean formalizations of mathematics competition problems, solved using a range of techniques such as induction, algebraic manipulation, and contradiction; and (ii) a set of Coq problems from the CompCert compiler verification project (Leroy, 2009) that was previously used to evaluate the PROVERBOT9001 system Sanchez-Stern et al. (2020).

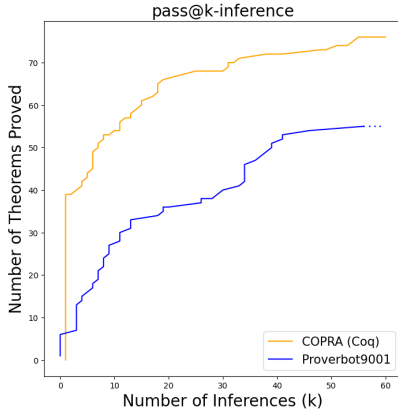


Figure 6: COPRA vs. PROVERBOT9001 on the CompCert benchmark

our evaluation on miniF2F turns off COPRA’s and REPROVER’s retrievers.

In the CompCert domain, we compare with PROVERBOT9001 (Sanchez-Stern et al., 2020), which, while not LLM-based, is the best publicly available model for Coq. Unlike miniF2F, this benchmark comes with a large training set as well as a test set, and we use the training set for retrieving relevant lemmas and definitions. Our retrieval mechanism, in this case, is a simple BM25 search.

**Implementing COPRA.** Our implementation of COPRA has GPT-4 as the underlying LLM and can interact with both the Lean and the Coq proof environments. Because of the substantial cost of GPT-4 queries, we cap the number of LLM queries that COPRA can make by 60. To further reduce costs, COPRA first tries to prove its theorems via a single LLM query (one-shot prompting). It only invokes its agent behavior when the one-shot prompting fails to find a proof.

The “system prompt” in the one-shot approach is slightly different than that for COPRA, containing instructions to generate a proof in one go rather than step by step. For both COPRA and the one-shot baselines, the prompt contains a single proof example that clarifies how proofs need to be formatted. This proof example remains the same for all test cases.

**Benchmarks.** We evaluate our approach on two domains: (i) miniF2F (Zheng et al., 2021), a collection of

**Baselines.** We compare with one-shot invocations of GPT-3.5 and GPT-4 in both the miniF2F and the CompCert domains. We also consider an ablation of COPRA that uses GPT-3.5 as its LLM and another that does not use backtracking. As for fine-tuned baselines, a challenge is that all existing open-source theorem-proving systems only target a single proof environment. As a result, we had to choose different baselines for the Lean (miniF2F) and Coq (CompCert) domains.

Our fine-tuned baseline for the miniF2F domain is REPROVER, a state-of-the-art open-source prover that is part of the Leandjo project (Yang et al., 2023). A challenge with this baseline is that like COPRA, it uses a retrieval mechanism. However, building a comparable retriever for COPRA would require an indexed training corpus on problems relevant to miniF2F. However, miniF2F is only an evaluation set and does not come with a training corpus. As a result, for an apples-to-apples comparison,

Approach	# Theorems proved /# Theorems	% proved	Avg. Inferences in Total	Avg. Inferences on Failure	Avg. Inferences on Pass	Max. Inferences Allowed
miniF2F Test Dataset						
GPT 3.5 One-Shot	7/244	2.8%	1	1	1	1
GPT 4 One-Shot	26/244	10.6%	1	1	1	1
COPRA (GPT-3.5)	29/244	11.89%	12.83	14.23	2.45	60
ReProver	54/244	22.13%	350.7	427.24	81.6	1076
<b>COPRA (GPT-4)</b>	<b>57/244</b>	<b>23.36%</b>	<b>20.94</b>	<b>26.79</b>	<b>1.75</b>	60
CompCert Test Dataset						
GPT 3.5 One-Shot	10/118	8.47%	1	1	1	1
GPT 4 One-Shot	36/118	30.51%	1	1	1	1
Proverbot	<b>98/118</b>	<b>83.05%</b>	184.7	256.8	170.0	2344
<b>COPRA</b>	76/118	64.41%	<b>12.9</b>	<b>10.9</b>	<b>16.57</b>	60

Table 1: Aggregate statistics for COPRA and the baselines on miniF2F and CompCert

Approach	Avg. Time In Seconds					
	Per Proof			Per Inference		
	On Pass	On Fail	All	On Pass	On Fail	All
ReProver (on CPU)	279.19	618.97	543.78	3.42	1.45	1.55
ReProver (on GPU)	267.94	601.35	520.74	2.06	0.44	0.48
COPRA (GPT-3.5)	39.13	<b>134.26</b>	<b>122.21</b>	15.97	9.43	9.53
<b>COPRA (GPT-4)</b>	<b>30.21</b>	191.73	140.86	17.26	7.16	6.73

Table 2: Average time taken by our approach (COPRA) and ReProver on miniF2F dataset.

For cost reasons, our evaluation for CompCert uses 118 out of the 501 theorems used in the original evaluation of PROVERBOT9001 Sanchez-Stern et al. (2020). For fairness, we include all the 98 theorems proved by PROVERBOT9001 in our subset. The remaining theorems are randomly sampled.

**Metric: pass@ $k$ -inferences.** The standard metric for evaluating theorem-provers is *pass@ $k$*  (Lample et al., 2022; Yang et al., 2023). In this metric, a prover is given a budget of  $k$  *proof attempts*; the method is considered successful if one of these attempts leads to success. However, a key objective of our research is to discover proofs *quickly*, with fewer LLM queries and lower wall-clock time. The *pass@ $k$*  metric does not evaluate this characteristic as it does not quantify the number of LLM queries or amount of time needed by a proof attempt.

Approach	# Theorems proved /# Theorems	% proved
miniF2F Test Dataset		
COPRA (GPT-4) w/o backtracking	56/244	22.95%
<b>COPRA (GPT-4)</b>	<b>57/244</b>	<b>23.36%</b>
CompCert Test Dataset		
COPRA (GPT-4) w/o backtracking	52/118	44.06%
<b>COPRA (GPT-4)</b>	<b>76/118</b>	<b>64.41%</b>

Table 3: The effectiveness of backtracking

Figure 6 shows a comparison between COPRA and PROVERBOT9001.

To address this concern, we introduce a new metric, *pass@ $k$ -inferences*, and evaluate COPRA and its competitors using this metric. Here, we measure the number of correct proofs that a prover can generate with a budget of  $k$  or fewer LLM inference queries. One challenge here is that we want this metric to be correlated number of correct proofs that the prover produces within a wall-clock time budget; however, the cost of an inference query is propor-



tional to the number of responses generated per query. To maintain the correlation between the number of inference queries and wall-clock time, we restrict each inference on LLM to a single response.

**Results.** Figure 5 shows our comparison results for the `miniF2F` domain. As we see, COPRA outperforms REPROVER, completing, within just 60 inferences, problems that REPROVER could not solve even after a thousand inferences. This is remarkable given that COPRA is based on a black-box foundation model and REPROVER was fine-tuned for at least a week on a dataset derived from Lean’s Mathlib library. For fairness, we ran REPROVER multiple times with 16, 32, and 64 (default) as the maximum number of inferences per proof step. We obtained success rates of 15.9%, 20.1%, and 22.13% in the respective cases and took the best for comparison.

We find that COPRA is significantly faster than PROVERBOT9001. Since we put a cap of 60 inferences on COPRA, it cannot prove all the theorems that PROVERBOT9001 eventually proves. However, as shown in the figure, COPRA proves many more theorems than PROVERBOT9001 if only 60 inferences as allowed. Specifically, we prove 77.5% of all proofs found by PROVERBOT9001 in less than 60 steps.

Aggregate statistics for the two approaches, as well as a comparison with the one-shot GPT-3.5 and GPT-4 baselines, appear in Table 1. It is clear from this data that the language-agent approach offers a significant advantage over the one-shot approach. For example, COPRA solves more than twice as many problems as the one-shot GPT-4 baseline, which indicates that it does not just rely on GPT-4 recalling the proof from its memory. Also, the use of GPT-4 as opposed to GPT-3.5 seems essential.

We establish the correlation between the number of inferences needed for a proof and wall-clock time in Table 2. Although the average time per inference is higher for COPRA, COPRA still finds proofs almost 9x faster than REPROVER. This can be explained by the fact that our search is more effective as it uses 46x fewer inferences than REPROVER. These inference steps not only contain the average time spent on generating responses from LLM but at times have some contribution corresponding to the execution of the tactic on the Lean environment itself.

```
theorem algebra_sqineq_at2malt1
(a : ℝ) :
a * (2 - a) ≤ 1 :=
begin
  have h : ∀ (x : ℝ), 0 ≤ (1 - x) ^ 2,
  from λ x, pow_two_nonneg (1 - x),
  calc a * (2 - a)
      = 1 - (1 - a) ^ 2 : by ring
      ... ≤ 1 : sub_le_self _ (h a),
end
```

Figure 7: A theorem in the “algebra” category that COPRA could prove but REPROVER could not.

Table 2 also offers data on when the different approaches report failures. Since REPROVER uses a timeout for all theorems, we also use a timeout of 11 minutes while considering failures in Table 2. The data indicates that COPRA is comparatively better at giving up when the problem is too hard to solve. We also note that less time is spent per inference in case of failure for all approaches.

We show the impact of ablating the backtracking feature of COPRA in Table 3. We note that backtracking has a greater positive impact in the `Compcert` domain. We believe this is because the `Compcert` problems are more complex and backtracking helps more when the proofs are longer.

Finally, we offer an analysis of the different categories of `miniF2F` problems solved by COPRA and REPROVER in Figure 8. We see that certain kinds of problems, for example, International Mathematics Olympiad (IMO) problems and theorems that require induction, are difficult for all approaches. However, Figure 8b shows that COPRA takes fewer steps consistently across various categories of problems in `miniF2F`.

From our qualitative analysis, there are certain kinds of problems where the language-agent approach seems especially helpful. For instance, Figure 7 shows a problem

in the ‘algebra’ category that REPROVER could not solve. More examples of interesting Coq and Lean proofs that COPRA found appear in the appendix.



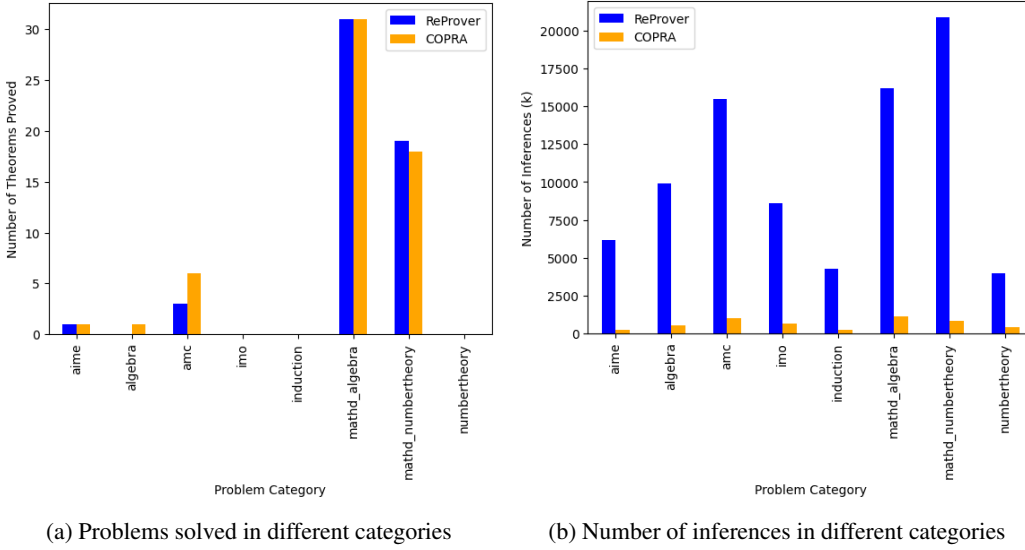


Figure 8: Breakdown of theorems proved in various categories

## 5 RELATED WORK

**Supervised Learning for Theorem-Proving.** There is a sizeable literature on search-based theorem-proving techniques based on supervised learning. These methods train a model to predict the next proof step at each point in a proof. This model is then used to guide a search technique, e.g., best-first or depth-limited search, that synthesizes a proof. Earlier methods of this sort used small-scale neural networks (Yang & Deng, 2019; Sanchez-Stern et al., 2020; Huang et al., 2019) as predictors. More recent methods, such as GPT-f (Polu & Sutskever, 2020), PACT (Han et al., 2021), HyperTree Proof Search (Lample et al., 2022), and REPROVER (Yang et al., 2023), have used LLMs. COPRA has some resemblance with the latter approaches. However, it departs from these prior methods in using execution feedback and a more sophisticated search algorithm.

The recent Draft-Sketch-Proof (Jiang et al., 2022) method relies on informal proofs to generate formal proofs.

Other methods like Baldur (First et al., 2023) generate the whole proof in one shot using an LLM and then *repair* it. The main ideas in these efforts — the use of informal proofs and repair models — are orthogonal to our approach.

**Reinforcement Learning for Theorem-Proving.** Kaliszky et al. (2018) pioneered the use of RL in theorem-proving; subsequently, Wu et al. (2021) gave TacticZero, a deep RL approach to the problem. TacticZero does not use LLMs, thus missing out on a key source of generic mathematical knowledge. Also, COPRA has retrieval capabilities that TacticZero lacks.

**Advanced Prompting Strategies.** Several prompting strategies like Chain-of-Thought (CoT) (Wei et al., 2022), Tree-of-Thought (ToT) (Yao et al., 2023), and Graph-of-Thought (GoT) (Besta et al., 2023) have recently emerged for modeling reasoning using LLMs. However, thought generation is not sufficient to emulate the rigorous verification that a formal proof environment can perform. This is why we use an approach based on language agents.

**Language Agents.** Several distinct LLM agent architectures have been proposed over the last year (Significant-Gravitas, 2023; Yao et al., 2022; Shinn et al., 2023; Wang et al., 2023). These models combine an LLM’s capability to use tools Schick et al. (2023), decompose tasks into subtasks (Wei et al., 2022; Yao et al., 2023), and self-reflect (Shinn et al., 2023). However, we are the first to offer an LLM agent for theorem-proving. We also distinguish ourselves from prior work along these lines by introducing a more efficient stateful search in the policy.



- Stanislas Polu and Ilya Sutskever. Generative language modeling for automated theorem proving. *arXiv preprint arXiv:2009.03393*, 2020.
- Alex Sanchez-Stern, Yousef Alhessi, Lawrence Saul, and Sorin Lerner. Generating correctness proofs with neural networks. In *Proceedings of the 4th ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*, pp. 1–10, 2020.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*, 2023.
- Noah Shinn, Federico Cassano, Beck Labash, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *arXiv preprint arXiv:2303.11366*, 2023.
- Significant-Gravitas. Autogpt. <https://github.com/Significant-Gravitas/Auto-GPT>, 2023.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.
- Minchao Wu, Michael Norrish, Christian Walder, and Amir Dezfouli. Tacticzero: Learning to prove theorems from scratch with deep reinforcement learning. *Advances in Neural Information Processing Systems*, 34:9330–9342, 2021.
- Kaiyu Yang and Jia Deng. Learning to prove theorems via interacting with proof assistants. In *International Conference on Machine Learning*, pp. 6984–6994. PMLR, 2019.
- Kaiyu Yang, Aidan M Swope, Alex Gu, Rahul Chalamala, Peiyang Song, Shixing Yu, Saad Godil, Ryan Prenger, and Anima Anandkumar. Leandojo: Theorem proving with retrieval-augmented language models. *arXiv preprint arXiv:2306.15626*, 2023.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*, 2023.
- Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. Minif2f: a cross-system benchmark for formal olympiad-level mathematics. *arXiv preprint arXiv:2109.00110*, 2021.

## A APPENDIX

### A.1 EXAMPLE PROOFS GENERATED FOR MINIF2F

Fig. 9 shows some other interesting proofs generated by our approach on miniF2F dataset.

### A.2 SYSTEM PROMPTS FOR MINIF2F

Parts of the ‘system prompt’ used by COPRA for theorem proving in Lean are shown in Fig. 10.

### A.3 SYSTEM PROMPTS FOR COMPCERT

Parts of the ‘system prompt’ used by COPRA for theorem proving in Coq are shown in Fig. 11.

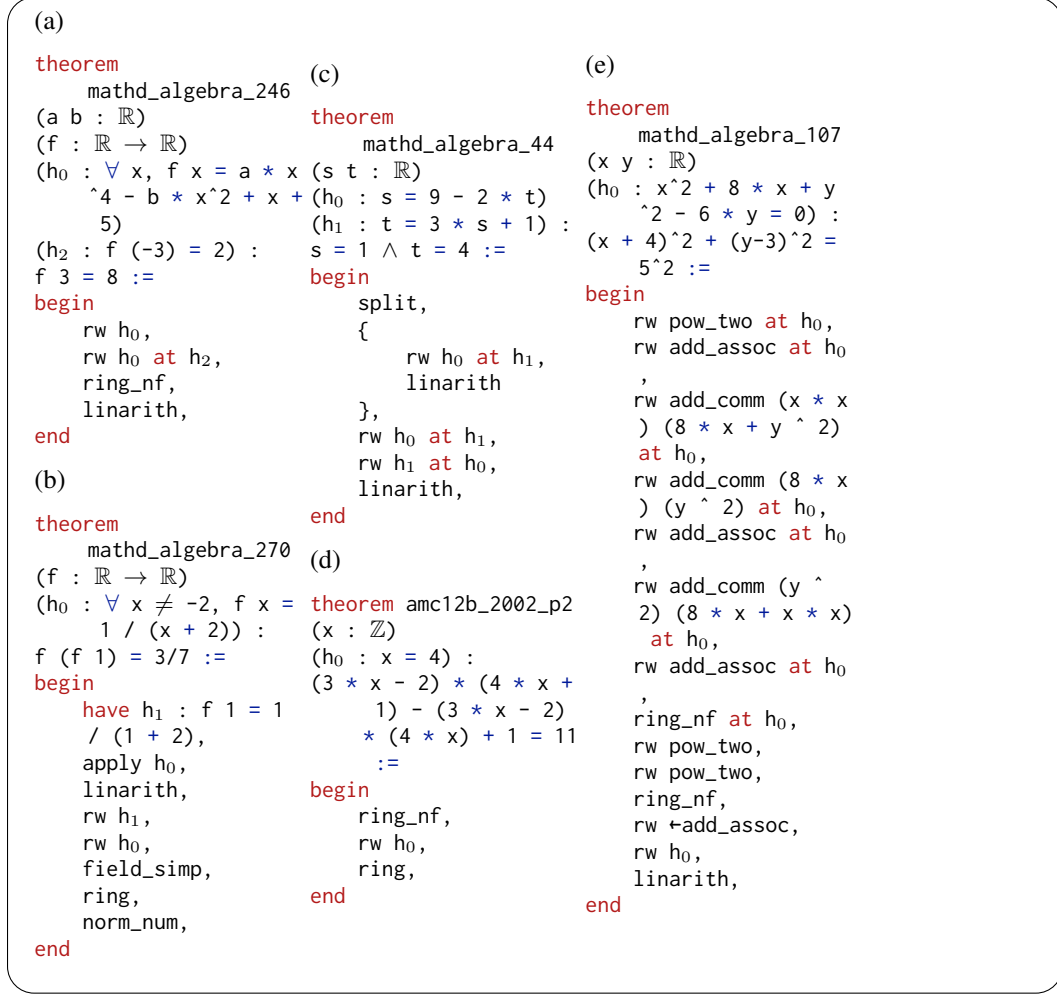


Figure 9: Some other interesting proofs generated for miniF2F by COPRA. The length of the proofs generated shows that interaction with the environment helps in fixing the errors encountered while writing long proofs. These long sequences of rewrites are not easy to synthesize without knowing the exact verbal reward from the environment which often contains the hint to fix the rewrites.

You are a proficient formal theorem-proving agent in Lean 3. You can predict  
 → the next proof step given the current proof state. The proof state is  
 → described in the following format:

1. All the goals are described under `[GOALS]` keyword. Each goal within  
 → the `[GOALS]` is described under the keyword `[GOAL] i`, where `i` is a  
 → positive integer. For example, `[GOAL] 1`, `[GOAL] 2`, etc.
2. Within each `[GOAL] i` keyword, the goal is described as a human-readable  
 → serialized version of the proof state as shown while running `lean`  
 → command. Each goal, might also accompany some hypotheses, which are  
 → described under the keyword `[HYPOTHESES] i`. Each hypothesis within  
 → `[HYPOTHESES]`, starts with the prefix `[HYPOTHESIS]`.
3. Sometimes `[GOALS]` can have description about the proof state like  
 → `Proof finished`, `There are unfocused goals`, `Not in proof mode`,  
 → etc. The description is described under the keyword `[DESCRIPTION]`.
4. Finally, `[STEPS]` keyword is used to describe proof-steps used so far.  
 → Each proof step starts with the prefix `[STEP]`, and is a valid Lean  
 → tactic. For example, `[STEPS][STEP]rw h1 at h2,[STEP]{linarith},`.
5. Sometimes, `[INCORRECT STEPS]` keyword optionally used to describe  
 → proof-steps which should NOT be generated. Use this as a hint for not  
 → generating these proof-steps again as they failed previously. For  
 → example, `[INCORRECT STEPS][STEP]apply h1,[STEP]rw +h1`.
6. There is also an optional `[LAST STEP]` keyword which describes the  
 → proof-step generated last time. If the proof-step was incorrect, then  
 → it is also followed by error message from Coq environment. For example,  
 → `[LAST STEP]linarith,\n[ERROR MESSAGE]linarith failed to find a  
 → contradiction\nstate:\nx y : ℝ,\nh1 : x = 3 - 2 \* y,\nh2 : 2 \* x - y =  
 → 1\n- false`. If the proof-step was correct then it is followed by the  
 → keyword `[SUCCESS]`. For example, `[LAST STEP]linarith,[SUCCESS]`.  
 → Don't generate the last proof-step again if it was NOT successful.
7. Sometimes there can be errors in the format of the generated response.  
 → This is reported using the keyword `[ERROR]` followed by the error  
 → message. For example, `[ERROR]\nInvalid response:\n'Great! The proof is  
 → complete.', \nStopping Reason: 'stop'.\n Please respond only in the  
 → format specified.[END]`. This means that the response generated by you  
 → was not in the specified format. Please follow the specified format  
 → strictly.

If you think you know the next proof step, then start your response with  
 → `[RUN TACTIC]` followed by the next proof-step which will help in  
 → simplifying the current proof state. For example, `[RUN  
 → TACTIC]induction c,[END]`. Generate exactly ONE proof-step. Multiple  
 → proof steps are more error prone, because you will not get a chance to  
 → see intermediate proof state descriptions. Make sure that the proof  
 → step is valid and compiles correctly in Lean 3.

You can refer to the example conversation to understand the response format  
 → better. It might also contain some similar proof states and their  
 → corresponding proof-steps.

Please take a note of the following:

1. Make sure to end all your responses with the keyword `[END]`. Follow the  
 → specified format strictly.
2. While generating `[RUN TACTIC]` keyword, do NOT generate the tactics  
 → mentioned under `[INCORRECT STEPS]`.....  
 .....

Figure 10: Parts of ‘system prompt’ used by COPRA for Lean

You are a proficient formal theorem-proving agent in Coq. You can predict

- the next proof step given the current proof state, relevant
- definitions, and some possible useful lemmas/theorems. The proof state
- is described in the following format:

1. All the goals are described under ``[GOALS]`` keyword. Each goal within
  - the ``[GOALS]`` is described under the keyword ``[GOAL] i``, where ``i`` is a
  - positive integer. For example, ``[GOAL] 1``, ``[GOAL] 2``, etc.
2. Within each ``[GOAL] i`` keyword, the goal is described as a human-readable
  - serialized version of the proof state as shown while running ``coqtop``
  - command. Each goal, might also accompany some hypotheses, which are
  - described under the keyword ``[HYPOTHESES] i``. Each hypothesis within
  - ``[HYPOTHESES]``, starts with the prefix ``[HYPOTHESIS]``. Apart from the
  - goal and hypothesis, some OPTIONAL keywords like ``[DEFINITIONS] i`` and
  - ``[THEOREMS] i`` are also present which describe the relevant definitions
  - of symbols used in that goal, and some possible useful theorems or
  - lemmas which might help in simplifying the goal. Each definition within
  - ``[DEFINITIONS]`` starts with the prefix ``[DEFINITION]``. Similarly, each
  - theorem/lemma under ``[THEOREMS]`` keyword starts with the prefix
  - ``[THEOREM]``. These definitions and theorems can be used to simplify the
  - goal using the tactics like `rewrite`, `apply`, etc. However, it is also
  - possible that these definitions and theorems are not used at all.
3. Sometimes ``[GOALS]`` can have description about the proof state like
  - ``Proof finished``, ``There are unfocused goals``, ``Not in proof mode``,
  - etc. The description is described under the keyword ``[DESCRIPTION]``.
4. Finally, ``[STEPS]`` keyword is used to describe proof-steps used so far.
  - Each proof step starts with the prefix ``[STEP]``, and is a valid Coq
  - tactic ending with a ``.``. For example, ``[STEPS][STEP]intros`
  - `a.[STEP]induction a.``.
5. Sometimes, ``[INCORRECT STEPS]`` keyword optionally used to describe
  - proof-steps which should NOT be generated. Use this as a hint for not
  - generating these proof-steps again as they failed previously. For
  - example, ``[INCORRECT STEPS][STEP]apply mul_assoc.[STEP]rewrite <- H.``.
6. There is also an optional ``[LAST STEP]`` keyword which describes the
  - proof-step generated last time. If the proof-step was incorrect, then
  - it is also followed by error message from Coq environment. For example,
  - ``[LAST STEP]reflexivity.[ERROR MESSAGE]Error: In environment\nn :`
  - `nat\nUnable to unify "n" with "n + 0".``. If the proof-step was correct
  - then it is followed by the keyword ``[SUCCESS]``. For example, ``[LAST`
  - `STEP]reflexivity.[SUCCESS]``. Don't generate the last proof-step again
  - if it was NOT successful.
7. Sometimes there can be errors in the format of the generated response.
  - This is reported using the keyword ``[ERROR]`` followed by the error
  - message. For example, ``[ERROR]\nInvalid response:\nGreat! The proof is`
  - `complete.', \nStopping Reason: 'stop'.\n Please respond only in the`
  - `format specified.[END]``. This means that the response generated by you
  - was not in the specified format. Please follow the specified format
  - strictly.

If you think you know the next proof step, then start your response with

- ``[RUN TACTIC]`` followed by the next proof-step which will help in
- simplifying the current proof state. For example, ``[RUN TACTIC]destruct`
- `c.[END]``. Generate exactly ONE proof-step. Multiple proof steps are
- more error prone, because you will not get a chance to see intermediate
- proof state descriptions. Make sure that the proof step is valid and
- compiles correctly with Coq.

.....

Figure 11: Parts of 'system prompt' used by COPRA for Coq

```

(a)
gss :
forall l v m,
(set l v m) l = match l with R (c)
  ⇨ r => disjoint_cons_right
    v | S sl ofs ty => ⇨ : (e)
    Val.load_result    forall a l1 l2, set_locals_lessedef
    ⇨ (chunk_of_type ty) v disjoint l1 (a :: ⇨ : forall e1
    ⇨ end.                ⇨ l2) -> disjoint ⇨ e2,
Proof.                  ⇨ l1 l2.          ⇨ env_lessedef e1
  intros l v m.          Proof.          ⇨ e2 -> forall
  destruct l as [r | s o t].  ⇨ H.          ⇨ il,
  - unfold set.           ⇨ unfold          ⇨ env_lessedef
  destruct (Loc.eq (R r) (R ⇨ disjoint.    ⇨ (set_locals il
    ⇨ r)); [reflexivity |   ⇨ e1)
    ⇨ contradiction].    ⇨ (set_locals il
  - unfold set.           ⇨ H2.          ⇨ e2).
  destruct (Loc.eq (S s o t) ⇨ apply H.    Proof.
    ⇨ (S s o t));         assumption.    intros e1 e2 H.
    ⇨ [reflexivity |      right.          induction il as
    ⇨ contradiction].    assumption.    ⇨ [| a il'].
Qed.                    Qed.          - apply H.
(b)                    (d)          - intros.
eq : forall (p q: loc), {p = apply
  ⇨ q} + {p <> q}.        eq_int_type :    ⇨ set_var_lessedef.
Proof.                  ⇨ forall (x y:    apply IHil'.
  decide equality.       ⇨ int_type),    apply
  - apply mreg_eq.       ⇨ {x=y} + {x<>y}.    ⇨ Val.lessedef_refl.
  - decide equality.    Proof.          Qed.
  - decide equality.    decide
  apply Pos.eq_dec.     ⇨ equality.
  decide equality.      Qed.
  - decide equality.
Qed.

```

Figure 12: Some other interesting proofs generated for CompCert by COPRA. We can see that these proofs are long, and often use ‘apply’ tactic which shows that COPRA can effectively use the retrieved information to discharge the current proof state.

#### A.4 EXAMPLE PROOFS GENERATED FOR COMPCERT

Fig. 12 shows some interesting proofs generated by our approach on the CompCert dataset.