# Agentic Program Verification

HAOXIN TU, National University of Singapore, Singapore
HUAN ZHAO, National University of Singapore, Singapore
YAHUI SONG, National University of Singapore, Singapore
MEHTAB ZAFAR, National University of Singapore, Singapore
RUIJIE MENG, National University of Singapore, Singapore
ABHIK ROYCHOUDHURY, National University of Singapore, Singapore

Automatically generated code is gaining traction recently, owing to the prevalence of Large Language Models (LLMs). Further, the AlphaProof initiative has demonstrated the possibility of using AI for general mathematical reasoning. Reasoning about computer programs (software) can be accomplished via general mathematical reasoning; however, it tends to be more structured and richer in contexts. This forms an attractive proposition, since then AI agents can be used to reason about voluminous code that gets generated by AI.

In this work, we present a first LLM agent, AUTOROCQ, for conducting program verification. Unlike past works, which rely on extensive training of LLMs on proof examples, our agent learns on-the-fly and improves the proof via an iterative refinement loop. The iterative improvement of the proof is achieved by the proof agent communicating with the Rocq (formerly Coq) theorem prover to get additional context and feedback. The final result of the iteration is a proof derivation checked by the Rocq theorem prover. In this way, our proof construction involves autonomous collaboration between the proof agent and the theorem prover. This autonomy facilitates the search for proofs and decision-making in deciding on the structure of the proof tree.

Experimental evaluation on SV-COMP benchmarks and on Linux kernel modules shows promising efficacy in achieving automated program verification. As automation in code generation becomes more widespread, we posit that our proof agent can be potentially integrated with AI coding agents to achieve a generate and validate loop, thus moving closer to the vision of *trusted automatic programming*.

## 1 Introduction

The widespread adoption of AI-generated code has transformed the landscape of software development. Today, more than 25% of the new code at Google is generated by AI. Yet the increasing amounts of code often carry subtle semantic errors or vulnerabilities [40] that may elude human review and testing, necessitating automated reasoning on program behaviors to engender *trust*. One promising approach to attain strong, machine-checkable evidence about program properties is

Authors' Contact Information: Haoxin Tu, National University of Singapore, Singapore, haoxin.tu@nus.edu.sg; Huan Zhao, National University of Singapore, Singapore, zhaohuan@comp.nus.edu.sg; Yahui Song, National University of Singapore, Singapore, yahuisong123@gmail.com; Mehtab Zafar, National University of Singapore, Singapore, mehtab@nus.edu.sg; Ruijie Meng, National University of Singapore, Singapore, mengrj.cs@gmail.com; Abhik Roychoudhury, National University of Singapore, Singapore, abhik@nus.edu.sg.

formal verification, which has been pivotal in certifying high-assurance systems for decades, such as microkernels [25], compilers [29], and databases [36].

Formal verification techniques commonly (1) reduce programs and specifications to logical formulas as proof obligations, and (2) discharge those obligations to automatic or interactive theorem provers (ITPs). Common ITPs, such as Isabelle/HOL and Rocq/Coq, automatically reduce the proof goals by applying user-supplied proof steps called *tactics*. Unfortunately, formal verification sees limited adoption in production systems, as it is an extremely time- and skill-intensive undertaking. Specifically, (1) capturing program behaviors formally requires significant efforts in manually annotating specifications and crafting loop invariants, and (2) generating a mathematically rigorous proof to validate verification conditions demands a high degree of expertise. For instance, the verification of the seL4 microkernel [25] required 22 person-years of effort, while the proofs for the CompCert compiler took 6 person-years and 100,000 lines of Rocq code—eight times longer than the implementation itself [29]. Most of the specifications and proofs are meticulously crafted by the developers. Recently, advances in Large Language Model (LLM) agents have been transformative across various disciplines. In particular, LLMs have demonstrated remarkable capabilities in both code understanding [5] and general mathematical reasoning [17]—two foundational capabilities required by program verification. In this work, we thus ask the following question: *Can LLM agents be leveraged for automatic and end-to-end verification of real-world programs?*

*State-of-the-art and what is needed.* Since there have been approaches to discuss proof automation by leveraging LLMs [33, 48] and other machine-learning-based techniques [41, 43, 51]—let us try to understand the need for a new approach. First of all, it is reasonable to use machine learning models to learn and predict individual proof steps or tactics in a proof. This has been accomplished in several works, including Proverbot9001 [41]. It is also possible for LLMs to generate whole proofs, albeit potentially incorrect ones, and then have a repair step to correct them, if possible—an approach studied in PALM [33]. Rango [48] goes one step further. Given a sequence of proof steps, it uses an ITP for yes/no validation of the proof steps, thereby allowing the trial of different tactics.

*An agentic approach.* These approaches all use LLMs, but are not agentic. We aim to develop an approach where an LLM agent can understand the proof structure and autonomously seek help from a theorem prover while constructing the proof. Thus, apart from using a theorem prover for validation of proof steps, the agent can actively collaborate with a theorem prover to prove a program property. As such, we develop a proof automation agent AutoRocq which acts as an *interpreter* of a proof derivation, and collaborates with a theorem prover to extend/improve it. Specifically, AutoRocq *autonomously* interacts with the theorem prover ahead of tactic prediction to retrieve relevant lemmas in the context, and generates tactics smartly guided by tree-shaped proof representations. It also incorporates feedback from the interactive theorem prover and proof histories to refine tactic generation.

*Evaluation.* We thoroughly evaluate our approach on existing mathematical lemmas [51], as well as lemmas systematically extracted from SV-COMP programs [7]. The SV-COMP programs' lemmas capture intricate code logic and properties, which are more representative of program verification. Evaluation shows that AutoRocq significantly outperforms state-of-the-art approaches. Specifically, AutoRocq is capable of proving 51.1% mathematical lemmas and 30.9% program lemmas, exceeding baseline approaches by 20.8% to 343.0% on mathematical lemmas, and by 42.4% to 204.6% on program lemmas. Among the successes, 142 lemmas (98 mathematical lemmas and 44 program lemmas) are uniquely proved by AutoRocq, due to its agentic design and access to proof contexts. We also conduct a case study to verify the code in Linux kernel modules (i.e., memory management utilities), where AutoRocq automatically verifies 12 lemmas for the function correctness property, compared to only 2–10 lemmas verified by baseline approaches.

*Contributions.* In summary, the contributions of this paper are as follows.

- We build the first proof automation agent, AUTOROCQ, which is highly effective in automatic proof generation. It acts as an interpreter of proof representations and actively collaborates with the Rocq prover to construct proofs.
- We showcase the feasibility of automatic and end-to-end verification with LLM agents. AUTOROCQ is evaluated on the widely used SV-COMP programs in software verification, as well as Linux kernel modules.
- We conduct an empirical study to demonstrate proof context- and structure-awareness can help LLM agents reason about software *formally*.

## 2 Background

### 2.1 Formal Program Verification

Formal program verification uses mathematically rigorous methods to prove that a program satisfies its specifications. Unlike traditional testing, which can only expose errors for specific inputs or executions, formal verification provides logical guarantees that a program behaves correctly for all possible executions within a given semantic model. Within this landscape, a prominent paradigm is *deductive verification*. In this approach, the program code is annotated with formal specifications, such as preconditions, postconditions, and loop invariants. These annotations form *formal contracts* that describe the program's intended behavior at the module and function boundaries. The annotated code is then compiled into a finite set of logical *proof obligations* under an explicit semantic model. The discharge of all proof obligations implies that the program satisfies its specifications.

To facilitate deductive verification, a variety of tools have been developed, including DAFNY [28, 37], VERUS [27], VIPER [38], and FRAMA-C [24]. These tools support annotating code, generating proof obligations, and attempting to discharge them. Automated solvers [6, 14] can be used to handle simple proof obligations (e.g., quantifier-free linear arithmetic [21]). More complex proof obligations—such as those involving non-linear reasoning, quantifier handling, or sophisticated inductive proofs—require human guidance. In such cases, *interactive theorem provers* (ITPs) are used, where auxiliary *lemmas* are introduced to capture intermediate properties, helping to structure and simplify the proofs of these obligations.

Rocq (formerly Coq) proof assistant [11] is a widely-used interactive theorem prover for semi-automated validation of residual proof obligations. In Rocq, proofs are constructed in a top-down manner from a *proof goal*—typically the statement of a theorem or lemma. Users iteratively apply *proof tactics*, which are simple or parameterized commands that decompose the goal into zero or more subgoals and guide the construction of a proof term. This proof term is type-checked by Rocq's trusted kernel. A proof succeeds when the subgoal list becomes empty, and the output proof term is the formal certificate of validity. This certificate can be reproduced by replaying the exact sequence of tactics, referred to as the *proof script*. Conceptually, successful proof induces a *proof tree* rooted at the original proof goal. Each node represents a subgoal produced during the derivation, and each edge corresponds to a tactic application. The interactive workflow of Rocq supports incremental, exploratory proof development with frequent inspection and backtracking.

### 2.2 Machine Learning for Proof Generation

Proof generation aims to automatically synthesize proof scripts that discharge the given proof obligations, which closes the last mile of program verification. Fueled by recent advances in AI, a growing body of work tackles this task with machine-learning methods. Early works [19, 41, 42, 51] often cast this problem as a sequence generation task. They design smart representations of syntactic constraints [51], proof states [19], and fine-grained proof contexts [42] to enable neural

networks to generate valid tactic sequences. However, applying a tactic successfully does not always advance the proof, as it may yield equivalent or even harder subgoals. This mismatch limits standalone prediction networks, which fixate on next-step tactic generation without a reliable global progress signal. To address this, an orthogonal line of work [8, 43] augments these networks with search strategies to guide tactic generation. For instance, QEDCartographer [43] extends Proverbot9001 [41] with a tactic selection model, trained through reinforcement learning with a fine-grained reward function.

On the other hand, Large Language Models (LLMs), with their remarkable capabilities in mathematical reasoning and high-level understanding, have become a promising alternative. A recent study [33] shows that LLMs can grasp the high-level structure of proofs, though they often stumble on rudimentary errors such as invalid references. Building on this observation, PALM first prompts the LLM to generate a complete proof script, and then applies deterministic repair mechanisms supplemented by the external tool CoqHammer [13] to patch errors. More recently, Rango [48] fine-tuned an LLM as a knowledge base for tactic generation. At every step of the proof, Rango retrieves the most relevant proofs and lemmas for the current proof state based on the predefined similarity function, and then prompts its knowledge base to generate the next tactic.

These existing works have made significant advances in automated proof generation. However, most approaches rely on models that are trained or fine-tuned on large corpora of existing proofs to guide tactic generation. A widely-used large-scale training dataset is CoqGym [51], which contains 71K human-written proofs that cover mathematics, computer hardware, and programming language. Models trained or fine-tuned on this dataset are effective in automatically generating proof scripts for mathematical theorems and libraries. Unfortunately, their effectiveness is reduced when applied to proofs about computer programs, which are typically more complex in both structure and semantics. Moreover, existing works adopt deterministic mechanisms to patch proof errors, retrieve relevant contexts, and generate subsequent tactics, preventing them from making adaptive, demand-driven decisions. For example, context retrieval is usually based solely on similarity functions, which may not provide sufficient information for generating the next tactic. Such a *non-agentic* context search strategy further limits the ability of these approaches.

## 3 Motivating Example

Although many automatic proving approaches have been proposed, to our knowledge, none of them prove lemmas in an *agentic* fashion. To motivate our agentic approach, AutoRocq, we use a concrete example shown in Figure 1. Figure 1(a) shows the proof obligation (wp_goal, lines 1-7) extracted from verifying the *function correctness* property of a real-world program cggmp2005b in SV-COMP [7], along with the complete proof (lines 9-23) generated by AutoRocq. The proof tree representation that guides AutoRocq's proving process is visualized in Figure 1(b). This example proof goal requires proving that "*i1 = 10%Z*" (line 7) holds for integer-typed (Z) variables i1 and i under a complex set of hypotheses involving modular arithmetic, inequalities, and type constraints (lines 2-6). The hypotheses capture program semantics including path conditions (e.g., "*i <= 10%Z*") and non-overflow constraints ("*(-2147483648%Z)%Z <= x*"), all discharged by Frama-C during automated verification (see Section 5.1.2 for more details). This example represents a typical lemma derived from real-world verification tasks. It captures intricate program logic and thus requires sophisticated reasoning about nested quantifiers, integer arithmetic, and logical constraints to handle. The proof tree in Figure 1(b) illustrates how AutoRocq systematically decomposes the problem through strategic tactic application, including case analysis with destruct that creates multiple proof branches, and maintains rich contextual information throughout the proving process. Such complex lemmas pose significant challenges to existing approaches (e.g., [33, 41, 43, 48]),
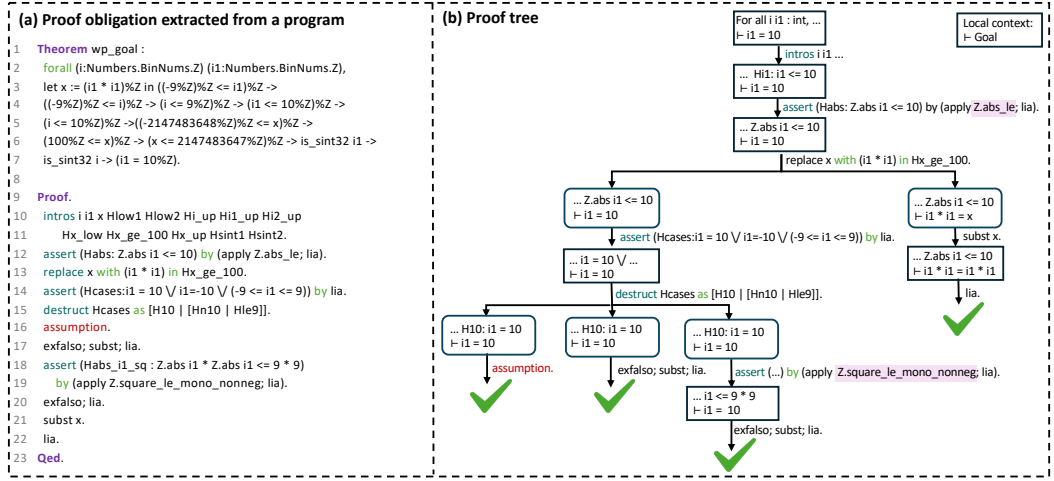
Fig. 1. (a) Proof obligation `wp_goal` extracted from `benchmark52_polynomial` in SV-COMP [7] and its proof generated by AutoRocq, with (b) the proof tree constructed during proving, where lemmas highlighted are retrieved from the global context autonomously by AutoRocq.

and none of them can prove this lemma during our experiments. We highlight these challenges to motivate AutoRocq's agentic design in the following.

**Verbose and crowded context.** The local context of `wp_goal` carries rich semantic information—six hypotheses are present in the statement simultaneously. These hypotheses include conjunctions, disjunctions, inequalities, and type constraints. Conventional approaches, which mostly employ naive or similarity-based retrieval, could easily get a surfeit or an insufficiency of contextual information. For example, before proof generation, PALM [33] simply retrieves *all* relevant premises from the global environment. In this case, it would be overwhelmed by hundreds of existing premises on integer arithmetic, most of which are irrelevant in the context. Conversely, Rango [48] retrieves lemmas based on lexical similarities, an imprecise metric that often fails to identify semantically relevant lemmas. This is why Rango fails to prove the example: it cannot pinpoint the specific contextual assumptions on integer bounds and arithmetic properties required for the proof.

In contrast, AutoRocq employs an agentic context search that *autonomously* decides, based on the current proof state, if it should make a tactic prediction, or if it should gather more context. In the latter case, we facilitate its requests for context with fine-grained query commands (see Table 1). For instance, it can leverage "*Search (Z.abs _ <= _).*" to fetch all lemmas and definitions that involve absolute values (`Z.abs`) and match the specified wildcard pattern. As a result, the lemma `Z.abs_le` (line 12 in Figure 1(a)) is retrieved, which is crucial for reasoning about absolute values in the context. This autonomy enables the agent to retrieve additional context from the theorem prover *on demand*, supporting it to reason about lemmas with rich contexts. It is worth noting that our contribution is not on how to extract certain terms or patterns in the theorem prover; rather, it is to build an agentic system that knows *when* to search and *what* patterns to search intelligently.

**Understanding of proving progress.** A complicated lemma such as `wp_goal` requires a delicate sequence of tactics to prove. When generating such sequences, however, existing LLM-based methods lack structured representations of the proof process. They typically rely on simple goal lists rather than structured proof trees. This textual and linear representation of the proof state makes them overly focus on predicting a single tactic without a holistic view of the proving

progress. As a result, they often fail to adapt to the evolving state of the partially generated proof. AutoRocq explicitly maintains a well-structured proof tree representation, as shown in Figure 1(b). It captures the hierarchical structure and dependencies of the proof derivation. The proof tree is built gradually as the proving advances. For instance, AutoRocq carefully decomposes the original goal into multiple subgoals with the `destruct` tactic, resulting in the branching in Figure 1. This structure-aware representation enables strategic reasoning about proof progress, tactical decisions, and effective tracking of proof dependencies. As such, our proof agent is able to *interpret* the proof derivation from a higher level, thus achieving more effective proof generation.

**Harnessing the feedback.** Developing proofs, especially for intricate lemmas, is a trial-and-error process, and interactive theorem provers are designed to provide timely feedback to guide the proving process. Nonetheless, existing approaches hardly leverage this opportunity to refine their tactics. For example, PALM [33] only performs deterministic repairs to failed tactics, whereas Rango [48] takes only binary signals from the proof assistant and retries the prediction again. AutoRocq implements an effective feedback mechanism that incorporates the feedback from the prover to refine individual tactics. In addition, it also recognizes persistent errors from the history, and conducts autonomous context searches to collect additional context progressively. This process allows our agent to learn continuously from both failures and successes. In this example, despite 48 failed tactic applications (59 attempts in total) during proving, AutoRocq eventually recovers and synthesizes a complete proof in a short time (i.e., 156.7 seconds).

Collectively, we have an agentic proof system that can autonomously decide when and how to incorporate additional contexts and adapt strategies based on the feedback. A high-level interpretation of the proof derivation from the expressive proof tree supports such an agency. In principle, this process is analogous to how an expert human prover would approach the task. In Section 4, we will explain each of these components in detail.

## 4 Proof Automation with AutoRocq

**Overview.** Figure 2 shows the overall workflow of AutoRocq. At a high level, it is an LLM agent that takes autonomous actions to carry out a human-like proving process. Specifically, given a lemma and the initial context (i.e., proof state), AutoRocq performs an agentic context analysis to generate a tactic if the context is sufficient, or a query command if additional context is needed (Section 4.1). Then, AutoRocq generates a sequence of tactics by interpreting formal proof representations (i.e., proof tree), analyzing them, and communicating with the Rocq proof assistant autonomously (Section 4.2). During communication, a context-assisted feedback loop is established to either refine an incorrect tactic application or provide additional context to the LLM if the errors persist (Section 4.3). The output of a successful proof is a certificate consisting of a sequence of tactics that can be replayed in the interactive theorem prover to verify the proof.

### 4.1 Context-Aware Tactic Generation

Previous works [30, 33, 48, 52] have demonstrated that additional information helps LLMs generate valid tactics. Such additional information, which we refer to as the *context*, typically include lemmas, definitions, and existing proof steps that are relevant in constructing the proof. Existing approaches, however, tend to augment LLMs with blindly selected context – typically a long list of premises [33] or existing proofs [48] ordered by a predefined similarity metric – without a good understanding of the current proof state or the progress of the proof. We argue that, before suggesting a new tactic, an intelligent proof system should be able to analyze the current proof state and the progress of the proof to decide if more context is needed. This system should carefully determine *what* context information is needed and *when* to add it, so that unnecessary noise is reduced to a minimum. To
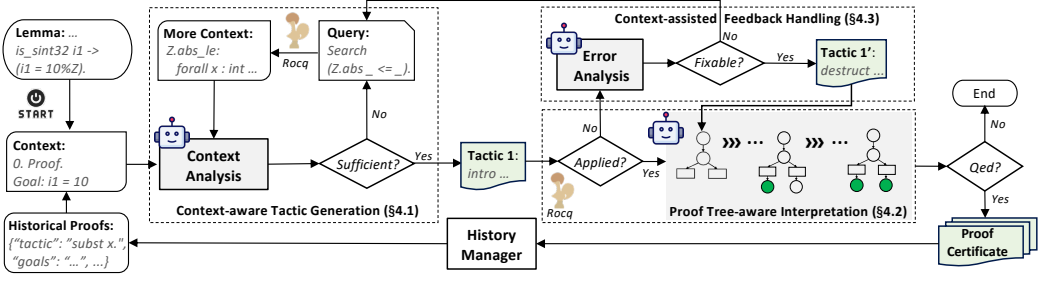
Fig. 2. Overview of AᴜᴛᴏRᴏᴄǫ, where components involving decision-making by LLMs are highlighted.

Table 1. List of Supported Context Query Commands in AᴜᴛᴏRᴏᴄǫ

| Command | Description | Output |
|---|---|---|
| Search <pattern> | Search for a pattern (e.g., a term). | A list of matched declarations. |
| Print <identifier> | Print the definition of an identifier. | Full definition of the identifier. |
| Locate <identifier> | Locate where an identifier is defined. | Path or location of the identifier. |
| About <identifier> | Show information about an identifier. | Type and summary of the identifier. |
| Check <term> | Type-check a term or expression. | Type of the term or expression. |

support this autonomy, an agentic context-aware search mechanism is needed. Such a context-aware approach can significantly reduce the amount of context added to the LLM, which could help the LLM focus on the current proof state and generate a more accurate tactic.

To achieve agentic context search, we couple a context analysis module with an agent that autonomously decides when and how to query the proof assistant's database. The agent performs continuous context analysis: if the current context is deemed sufficient, the agent directly generates a tactic to advance the proof; if a key piece of information (e.g., the definition of a term is_sint32 or a proved lemma in an imported library) is missing, the agent can submit a context query to retrieve it. To retrieve the context, it constructs a precise *Search* command, such as "*Search is_sint32*". This command is forwarded to the Rocq proof assistant, which searches its currently loaded libraries for matching identifiers, lemmas, and definitions. The results of this query are then incorporated into the LLM's context, enriching it for subsequent attempts at tactic generation.

---

**Agentic Context Analysis**

**Prompt:** "Analyze the current {proof tree} and Top-5 {historical tactics}. If sufficient information is available, generate a tactic to proceed. Otherwise, output a query command to retrieve additional context."
**Response:** Search (Z.abs _ <= _).

---

The supported context query commands are listed in Table 1. By default, AᴜᴛᴏRᴏᴄǫ supports the standard query commands available in the Rocq GUI (CoqIDE), including *Search*, *Print*, *Locate*, *About*, and *Check*. These commands enable AᴜᴛᴏRᴏᴄǫ to retrieve critical context information on demand, significantly aiding the proof process. Among these, the *Search* command is the most frequently used (see Section 6.2). Its versatility is the key: it allows for discovery based on name patterns, type signatures, and existing lemma names [47]. This makes *Search* the primary tool for pinpointing relevant declarations without knowing their exact names or forms.

## 4.2 Proof Tree-Aware Interpretation

The key component empowering AutoRocq's agency is an expressive proof tree representation that enables high-level interpretation of the proof derivation. Compared to the textual or linear representation of proof states (e.g., a list of goals), the explicit tree structure conduces understanding of the proving process. Formally, we define a proof tree as follows [22, 26].

**Definition** (*Proof Tree*). Let $\Sigma$ be the proving context[1], $G$ be the initial proof goal, and $S$ be a list of tactics that proves $G$ under $\Sigma$. A proof tree $T(\Sigma, G, S)$ is inductively defined as follows:
- Each node $N_A$ corresponds to a goal $\Sigma \vdash A$, i.e., a statement $A$ to be proved under $\Sigma$.
- The root node is $N_G := \Sigma \vdash G$.
- When a tactic $\tau \in S$ is applied to a node $N_A := \Sigma \vdash A$, it either (1) generates zero subgoals, in which case $N_A$ is a leaf, or (2) generates $n$ subgoals/nodes $\{N_{A_i} := \Sigma \vdash A_i \mid i = 1, \ldots, n\}$, and creates edges $N_A \rightarrow N_{A_i}$.
- A node is a leaf if it is directly derivable under $\Sigma$, i.e., no further subgoals are produced.

Intuitively, a proof script $S$ induces a proof tree $T$ in its enclosing context $\Sigma$, rooted at the original proof goal $G$. Each node represents a subgoal to be proved, produced a tactic application is represented as an edge. Leaves are subgoals that are trivially true or could be easily proved with a single tactic. The goal is complete when all the sub-goals in the leaves are proved.

AutoRocq explicitly maintains a (partial) proof tree as it progresses in tactic generation. After a successful tactic application, it examines the open subgoals from the Rocq proof assistant. If the tactic successfully proves the current subgoal (i.e. a leaf node), AutoRocq marks it as closed, and shifts its focus to another residual subgoal. When the tactic creates multiple subgoals, it creates these nodes in the tree and determines one of the new subgoals to work on. In this way, the proof tree is updated by extending the current node. On the other hand, if the tactic fails, our agent remains at the current node on the proof tree for further attempts. During the entire proving process, the proof states in subgoals and tactic applications are continuously updated until a complete tree is constructed. Figure 1(b) shows an example of a proof tree after all tactics were applied in the motivating example presented in Figure 1(a). From the tree, we could see that four subgoals were generated in total, and the lemmas with a violet background were retrieved from our agentic context search. The proof is completed when all leaves are closed.

Whenever the proof tree is updated, AutoRocq *interprets* the tree structure comprehensively to reevaluate the current progress. It does so by traversing the (partial) proof tree to identify the open subgoals, and how they relate to the existing proof structure. With a holistic view of the entire derivation process, it then proceeds to generate a sequence of tactics iteratively.

We note that our key contribution is not about formulating the proof tree structure. Instead, our insight lies in the realization that high-level interpretation enables LLM agents to carry out autonomous decision-making and actions. The benefits provided by our proof tree representation are twofold. First, it captures the hierarchical structure of the proof, including the relationships between tactics and subgoals, which provides a richer context for interpreting the proof process. Second, it allows AutoRocq to systematically break down the proof into smaller subgoals and tackle them one by one, which is more aligned with how humans approach proving.

## 4.3 Context-Assisted Feedback Handling

To emulate a trial-and-error feedback loop, we augment AutoRocq with feedback handling mechanisms to collaborate with the Rocq ITP. Based on the error that occurs, AutoRocq employs two strategies during tactic generation: (1) revising the tactic based on the error message when a single

---

[1]Here, $\Sigma$ captures both the global environment and the local context, as they are not differentiated in our context search.

error occurs, and (2) initiating a context search to retrieve additional context when persistent errors occur. In addition, AUTOROCQ takes positive feedback from past successes. We delegate each mechanism to a subsection as below.

*4.3.1  Handling a Single Error with Strategic Fixing.* When a tactic fails to apply to a subgoal, the ITP provides an error message indicating the reason for the failure. As shown in the following example, AUTOROCQ captures this error message and feeds it back to the LLM, along with the current proof tree. The LLM then analyzes the error message and the context to generate a revised tactic to fix the issue. Subsequently, the revised tactic is applied again, and the proof tree is updated accordingly.

---

**Agentic Error Handling of a Single Error**

---

***Prompt:*** "The previous {tactic} failed to apply to the current subgoal with the following {error message} from Rocq: {error message}. Analyze the error and generate a corrected tactic."
***Response:*** `destruct` H`cases as` [H10 | [Hn10 | H1e9]].

---

*4.3.2  Handling Persistent Errors with Context Search.* If the same error persists after several repair attempts (which is treated as *not fixable* as shown in Figure 2), AUTOROCQ recognizes that the current context may be insufficient to resolve the subgoal. In such cases, AUTOROCQ initiates a context search by issuing a query command to the ITP to retrieve additional context information, similar to the process described in Section 4.1. The retrieved context is then incorporated into the LLM's input, and the LLM generates a new tactic based on the enriched context. This context-assisted feedback loop continues until the subgoal is successfully solved or a predefined termination condition is met (e.g., maximum number of attempts). For example, the following shows how AUTOROCQ handles persistent errors. Given the failed tactics and the proof tree, our agent returns a query command to search for what is defined in the integer quotient from the `Z.quot` module.

---

**Agentic Error Handling of Persistent Errors**

---

***Prompt:*** "The agent has repeatedly generated {failed tactics} for the current subgoal multiple times. Analyze the current {proof tree} and determine what additional context is needed to proceed. Output a query command to retrieve the necessary context information."
***Response:*** `Search` (Z.quot _ _).

---

*4.3.3  Managing Historical Proofs.* When a lemma is successfully proved, AUTOROCQ stores comprehensive tactic history records that capture the complete context and evolution of each proof step. Each record contains the applied tactic, the proof goals before and after tactic application, the available hypotheses and their changes, along with metadata including the theorem name, and the tactic ID within the proof sequence. This rich historical information serves multiple purposes: (1) it enables AUTOROCQ to learn from successful proof patterns and reuse effective tactic sequences in similar contexts, (2) it provides detailed examples for context-aware tactic generation by showing how specific goals were transformed, and (3) it supports the feedback handling mechanism by offering concrete instances of successful problem-solving strategies. By maintaining this detailed provenance of proof construction, AUTOROCQ can build a knowledge base of proven tactics that enhances its capability to tackle new lemmas with similar structural patterns or mathematical properties (the experiment results presented in Section 6.2 also support our claim).

**Implementation.** We implemented AUTOROCQ in approximately 8k lines of Python code with a modular architecture that separates the core proving logic, theorem prover interface, and supporting utilities. The main Rocq backend is implemented using the CoqPyt (v1.0.0) library [9] to interact with the Rocq proof assistant. The implementation follows an iterative workflow as shown in

Figure 2. The modular design enables easy extension and maintenance while supporting the three key components of our approach: context-aware tactic generation, proof tree management, and feedback handling. Historical proof data is stored in JSON format to enable learning from successful proof patterns for future lemmas. For the LLM backend, we use GPT-4.1 as the backbone model for all LLM interactions. We set the temperature to 0 for reproducibility.

## 5  Experimental Setup

To evaluate the effectiveness of AutoRocq on program verification tasks, we seek to answer the following research questions (RQs):

**RQ.1:** Is AutoRocq more effective at proving lemmas than state-of-the-art approaches?
**RQ.2:** How does each component of AutoRocq contribute to its effectiveness?
**RQ.3:** How do the proofs generated by AutoRocq compare with human written proofs?

In this section, we first present our experimental setup. We then analyze the detailed results in Section 6. We also conduct case studies on verifying individual Linux kernel modules in Section 7.

### 5.1  Preparation of Lemmas from Programs

*5.1.1  Lemmas for Program Verification Tasks.* While benchmarks like CoqGym [51] provide a valuable corpus of human-written theorems, they are drawn primarily from mathematical libraries and do not capture the complexity of program verification lemmas found in real-world software. Since AutoRocq is designed to automate program verification, its evaluation requires a benchmark of lemmas extracted from real programs. To our knowledge, no such dataset currently exists, leaving a critical gap in the program verification field.

To address this gap, we construct a new benchmark by systematically extracting lemmas from real-world C programs. This benchmark enables the evaluation of AutoRocq and facilitates future research in automatic program verification. We source our subject programs from SV-COMP [7] for two key reasons: their widespread adoption in the formal verification community ensures they represent a diverse set of C language constructs, and their inclusion of ground-truth specifications provides meaningful proof obligations. Our selection focused on the most common property types in SV-COMP: functional correctness, defined by the unreachability of error calls, and non-overflow. The final criteria mandated that a program must (1) be verified for both properties and (2) exhibit deterministic behavior (e.g., no multi-threading). Applying these criteria yielded a final benchmark of 131 C programs from SV-COMP. The programs have an average of 43.06 lines of code, with the largest (the BusyBox [1] utility basename) comprising 428 lines.

*5.1.2  Lemma Extraction Methodology.* We designed a systematic method to automatically extract lemmas from the selected SV-COMP programs. We focus on generating non-trivial lemmas that necessitate reasoning about program semantics, as opposed to simple lemmas that can be discharged with basic tactics (e.g., auto.). Our approach utilizes Frama-C [12] to generate the necessary proof obligations from the program code.

Technically, we first use the RTE plug-in [3] to annotate the source code with formal contracts in the ANSI/ISO C Specification Language (ACSL), including preconditions and postconditions. As Frama-C cannot automatically infer loop invariants, we address this by employing LLMs to generate candidate invariants based on the loop's context and structure. Each candidate is validated through property testing with the Eva plug-in [2]; unsuccessful candidates are refined iteratively until a verifiably correct invariant is established. Note that while property testing can falsify incorrect invariants by discovering counterexamples, it cannot formally prove the correctness of an invariant – that it holds for all possible program executions. We justify this design choice as loop invariant inference is a long-standing, challenging problem [20, 35] that is orthogonal to our core
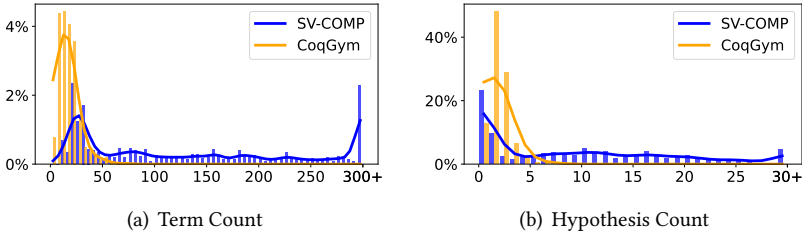
(a) Term Count             (b) Hypothesis Count

Fig. 3. Histogram of different complexity metrics: lemmas from SV-COMP programs (blue) vs. CoqGym (orange). Not a single lemma from CoqGym involves more than 100 terms or 7 hypotheses.

focus on proof automation. Once the code is fully annotated and verified, we employ Frama-C's WP plug-in [4] to generate proof obligations. These obligations are automatically translated into lemmas, which can be discharged in the Rocq proof assistant. The overall framework completes an end-to-end automated workflow from C programs to verifiable Rocq lemmas.

As a result, we collected 641 lemmas extracted from the verification of two types of target properties: non-overflow and functional correctness. Since both property types often require reasoning about loops, a significant portion of the generated proof obligations are related to loop invariants. These are essential intermediate lemmas needed to conclude that either a non-overflow or a functional correctness property holds for the entire program. Consequently, the 641 lemmas comprise three categories, reflecting their role in the proof: i. non-overflow lemmas (203, 31.7%): direct proof obligations for non-overflow properties; ii. functional correctness lemmas (153, 23.9%): direct proof obligations for functional properties; and iii. loop invariant lemmas (285, 44.5%): Supporting lemmas required to prove the loop invariants necessary for establishing both non-overflow and functional correctness properties. Together, these lemmas represent a diverse set of challenging proof obligations that arise in real-world program verification tasks.

*5.1.3 Comparison of Extracted Lemmas and Existing Lemmas.* To quantify the complexity and difficulty of lemmas extracted from SV-COMP programs, we compare them against existing lemmas from CoqGym. To do this, one commonly used proxy is the difficulty of proving a lemma, since more sophisticated lemmas generally necessitate longer proofs involving more cases and derivation steps. Unfortunately, this is infeasible for our purpose, since no ground-truth proofs are readily available for our extracted lemmas. As such, we directly examine these lemmas themselves from the extracted 641 lemmas from Section 5.1 and 625 lemmas across four projects (i.e., dblib, zfc, hoare-tut, and huffman) from CoqGym's testing partition [51]. We choose these four projects as they represent the common type of existing human-written lemmas, such as logical tautology, coding algorithms, and theory formalizations, and they have been extensively studied in the research community [8, 33, 48]. Specifically, we propose two metrics: (a) term count: the number of terms (variables, quantifiers, operators) to gauge structural complexity, and (b) hypothesis count: the number of assumptions to measure the richness of the semantic context.

In Figure 3, we report the complexity distribution of the lemmas by measuring the percentage of lemmas (y-axis) with a given complexity metric (term or hypothesis count, x-axis). For both metrics, the lemmas from CoqGym (orange) densely concentrate on the left end of the axis, suggesting uniformly lower complexity. In contrast, lemmas extracted from SV-COMP programs (blue) exhibit a more even distribution, with a lower peak and a long right tail, reflecting generally higher complexity. Concretely, 48% of these lemmas from SV-COMP programs consist of >100 terms, and 52% of them are associated with >7 hypotheses; *no* lemma from CoqGym satisfies either constraint. The stark contrast across both metrics suggests that the obligations extracted from real-world

programs, which often capture intricate logic in the code, tend to have much higher complexity than their benchmark counterparts. Concretely, below shows an example lemma from hoare_tut, a project in CoqGym. It states that the non-equality function Zneq_bool correctly returns false only if x = y. It is much simpler than lemmas from program verification tasks, as shown in Figure 1(a).

```
Lemma Zneq_bool_false: forall x y, Zneq_bool x y=false → x=y.
```

## 5.2 Baselines and Benchmarks

**Baselines.** We compare AutoRocq with state-of-the-art tools developed for proof automation in Rocq. We choose four baselines, namely Rango [48], PALM [34], QEDCartographer [43] (or QEDC for short), and Proverbot9001 [41] (or P9001 for short). These tools have shown exceptional results on proof automation tasks, and they employ different machine-learning techniques: standard LLMs, fine-tuned LLMs, reinforcement learning, and recurrent neural networks (RNNs), respectively. AutoRocq and PALM, which invoke closed-source LLM through APIs, use GPT-4.1 as the backend with temperature 0. For the other three baselines, we use the pre-trained weights made available in their artifacts. These include a fine-tuned version of DeepSeek-Coder used by Rango, a reinforcement-learning-based tactic selection agent employed by QEDC, and a custom tactic prediction network in P9001. They also have access to a single Nvidia A40 GPU with 48GB of VRAM if needed. We use the default settings for all the tools. Rango times out after 10 minutes, and QEDC is limited to 512 generation steps. PALM and P9001 are executed until completion. We set up all the baselines in Docker on Ubuntu 22.04. Other than PALM, which uses older Rocq 8.12 for compatibility reasons, all the other tools are configured with Rocq 8.18.

**Benchmarks used in Evaluation.** To conduct a comprehensive evaluation, we use the two benchmark sets mentioned in Section 5.1.3, including 625 lemmas from CoqGym [51] and 641 lemmas extracted from real-world programs in SV-COMP [7]. At the end, we have 1,266 lemmas evenly distributed across CoqGym and SV-COMP programs for our evaluation.

## 6 Experimental Results
### 6.1 RQ.1: Effectiveness of AutoRocq

For this research question, we investigate AutoRocq's effectiveness in proof generation on both (a) a non-program-related subset of lemmas from CoqGym, and (b) new lemmas extracted from SV-COMP programs. We report the number of proved lemmas from each tool in Figure 4. Figure 4(a) shows the total number of lemmas successfully proved by each tool. Figure 4(b) presents the detailed breakdown of the lemmas proven by each tool, on CoqGym (left) and SV-COMP programs (right). We note that PALM's results (in green) may not reflect its full capability, as several lemmas in the benchmarks are incompatible with Rocq 8.12.

**Lemmas from CoqGym.** From Figure 4(a), we observe that AutoRocq successfully proves 319 of them (51.1%), outperforming other tools by a large margin of 20.8% (compared with Rango [48]) to 343.0% (compared with PALM [33]). Among the baseline tools, Rango performs the best on this dataset, synthesizing proofs for 264 (42.3%) of the lemmas on CoqGym. The Venn diagram (Figure 4(b), left) also demonstrates that AutoRocq proves the most lemmas (98) uniquely. Our tool is followed by Rango and PALM, which account for 30 and 18 uniquely proved lemmas, respectively. The results suggest that AutoRocq is highly effective in automatic proof generation.

**Lemmas from SV-COMP programs.** More importantly, we study complex proof obligations that represent real verification workflows. These lemmas from SV-COMP programs indeed pose a greater challenge to all LLM-based tools, namely AutoRocq, Rango, and PALM, as evidenced by their lower success rates overall. AutoRocq manages to prove 198 lemmas (30.9%) from the

(a) # of lemmas proved by each tool.

(b) Venn diagram of successfully proved lemmas by each tool, from CoqGym (left) and SV-COMP programs (right).
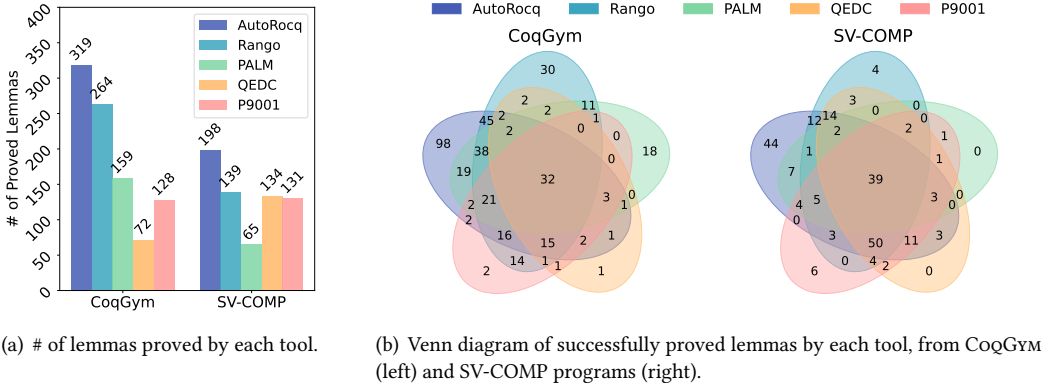
Fig. 4. [RQ.1] Lemmas proved by each tool on CoqGym and SV-COMP programs. On both benchmarks, AutoRocq is able to prove more lemmas, and has the most number of uniquely proved lemmas.



(a) Breakdown by term count.
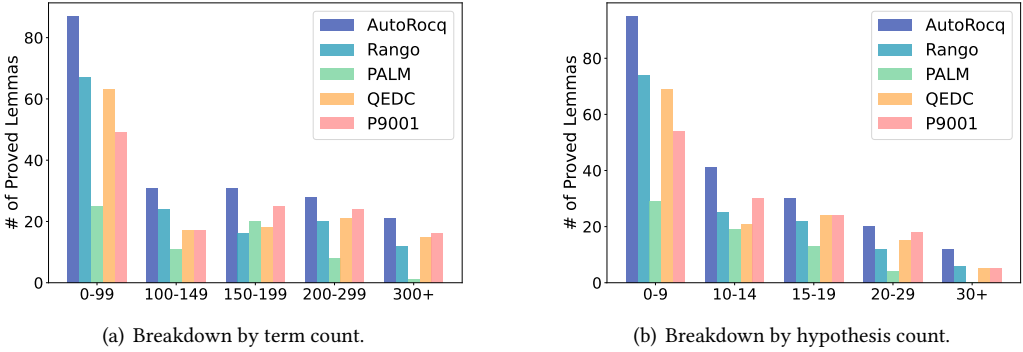
(b) Breakdown by hypothesis count.

Fig. 5. [RQ.1] # of lemmas proved from SV-COMP programs: a breakdown by lemmas' complexity.

benchmark, which is significantly higher than all other baseline tools. Results produced by Rango, QEDC, and P9001 lie in close proximity with one another, with each producing proofs for around 135 lemmas. AutoRocq outperforms them by 42.4%, 47.8%, and 51.1%, respectively. The number of lemmas proved by AutoRocq is also 204.6% higher than that of PALM.

Figure 4(b) (right) shows that AutoRocq generates proof scripts for 44 unique lemmas (owing to its agentic workflow) that *no* other tool can prove. In contrast, no other tool is able to generate proofs uniquely for more than 6 lemmas. In fact, only 23 lemmas proved by *any* other baselines elude AutoRocq! We note that there is a large overlap among the successes of baseline tools. There are 106 lemmas that are proved commonly by at least three of the baselines, suggesting a strong convergence in capabilities among them.

To better understand the high effectiveness of AutoRocq, we further correlate the number of successfully proved SV-COMP lemmas with their complexity. Figure 5 details the breakdown of successful attempts from different tools, grouped into five complexity buckets. We present the breakdown for both the term count (left) and the hypothesis count (right) in the original goal. Buckets on the right correspond to lemmas containing more terms/hypotheses, and thus are typically more structurally and contextually complicated. Results reveal an emerging pattern: as the original goal becomes more verbose and richer in context, the number of successfully proved lemmas decreases in general. However, AutoRocq retains its relative effectiveness in proving longer, more complex goals, resulting in a uniform lead across different complexity levels.

Additionally, we report the breakdown by the categories of lemmas, as shown in Figure 6. Results show that AUTOROCQ is particularly effective in proving non-overflow-related lemmas and loop invariants in the program, outperforming the best baselines by 53% and 22%, respectively. On lemmas related to functional correctness, AUTOROCQ also performs comparably with three baseline approaches. Results indicate AUTOROCQ's performance is consistent for verifying different types of program properties, showing its versatility and generality.

We also compare the efficiency of different tools in terms of the time taken to generate a successful proof on the SV-COMP lemmas. On average, AUTOROCQ takes 21.3 seconds to generate a successful proof, expending less time than other LLM-based approaches, namely RANGO (105.5 seconds) and PALM (45.4 seconds). AUTOROCQ's result is also comparable to P9001 (22.6 seconds). Notably, QEDC finishes extremely fast (5.1 seconds) when it does succeed, thanks to its well-optimized tactic-prediction model for the tactic generation.

We hypothesize that AUTOROCQ's remarkable efficacy and efficiency stem from its agentic access to the proof context, strategic tactic fixing, and effective communication between the LLM and ITPs. These collectively enable the agent to discover and apply relevant facts and lemmas dynamically during the proof search. These capabilities could be particularly beneficial for complex proof obligations, which often require non-trivial reasoning steps and the application of specific lemmas or theorems. In the next research question (Section 6.2), we evaluate the design decisions of AUTOROCQ to gain a deeper understanding of the system.

---

**Answer to RQ.1**: Overall, AUTOROCQ is significantly more effective in proving obligations extracted from real C programs (SV-COMP), outperforming baselines by 42%–205%. Notably, AUTOROCQ proves 44 lemmas that none of the other approaches can prove. It also outperforms the baselines by a similar margin on mathematical lemmas from COQGYM.
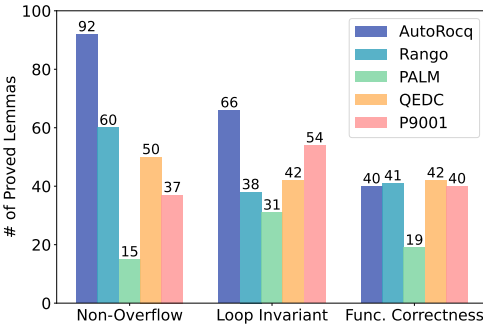
---



Fig. 6. [RQ.1] # of lemmas proved from SV-COMP programs: a breakdown by the category of lemmas.
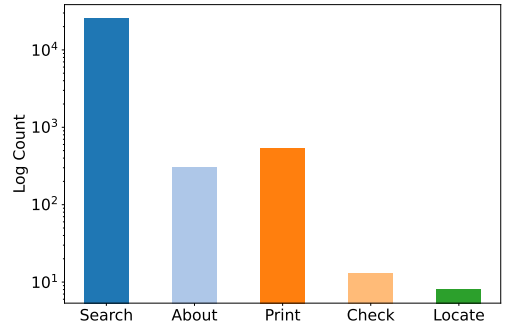


Fig. 7. [RQ.2] The log frequency of invocation for context search commands on SV-COMP programs.

## 6.2 RQ.2: Design Choices of AUTOROCQ

To better understand the interplay between AUTOROCQ's components, we evaluate how different design choices affect the overall effectiveness by implementing 6 variants of AUTOROCQ. We conduct the evaluation on 70 randomly sampled lemmas from the SV-COMP programs. Specifically, we examine the following variant approaches:

- ¬*CS*: No context search (Section 4.1). Instead, tactics are directly generated through prompting.
- ¬*PT*: No proof tree awareness (Section 4.2). Instead, proof states are encoded as plain texts.
- ¬*EF*, ¬*HM*: No error feedback or history manager (Section 4.3).

Table 2. [RQ.2] Success rates of different variants of AUTOROCQ on verifying sampled SV-COMP programs. All components contribute to the overall effectiveness, with context search ($\neg CS$) being the most significant.

|  | AUTOROCQ | $\neg CS$ | $\neg PT$ | $\neg EF$ | $\neg HM$ | $Err_1$ | $Err_5$ |
|---|---|---|---|---|---|---|---|
| Success Rate | 44.3% | 30.0% | 40.0% | 38.6% | 41.4% | 37.1% | 38.6% |
| Relative Improv. | – | +47.6% | +10.7% | +14.8% | +6.9% | +19.2% | +14.8% |

- $Err_1$, $Err_3$, and $Err_5$: The maximum number of errors encountered before context search. We evaluate three settings of the maximum number of errors encountered before seeking context search. We found that the value 3 is the most effective setting, and is used in AUTOROCQ. This is expected, as more frequent context queries ($Err_1$) may confuse the agent with unnecessary information, whereas fewer queries ($Err_5$) may not supply sufficient context.

We report the success rates, and AUTOROCQ's relative improvement, of different variants in Table 2. Results show that all variants only prove a fraction of lemmas compared to full-fledged AUTOROCQ. In particular, AUTOROCQ sees the most significant improvement from the ablative version $\neg CS$, which replaces the context search component with a direct request of tactic through prompting. This indicates that agentic access to the proof context contributes the most to our approach's effectiveness. Intriguingly, querying for contexts much more frequently ($Err_1$) or less frequently ($Err_5$) also affects the results negatively, suggesting that the amount of contexts supplied to the agent is paramount in practice. Removing other components in our approach, namely the proof-tree representation ($\neg PT$) or the feedback mechanisms ($\neg EF$ and $\neg HM$), each reduces the efficacy moderately by $\sim 10\% - 15\%$. We conclude that all components of AUTOROCQ are useful.

To visualize how AUTOROCQ conducts the context search, we summarize the frequency of invocation for different search commands provided by our framework (Table 1). We plot these frequencies in log scale, as shown in Figure 7. Statistics show that *Search* is the most frequently used command. Its primacy stems from its unique ability to solve the fundamental problem of discovery: it allows users to find relevant lemmas and theorems without prior knowledge of their names by searching for patterns. In essence, *Search* directly facilitates the core intellectual challenge of finding the right fact at the right time. It also reflects the fact that many proof obligations in our benchmarks require non-trivial reasoning steps, which often hinge on the application of specific lemmas or theorems. We dig into the specific patterns searched by AUTOROCQ, and find that they often involve complex expressions with multiple operators and operands, e.g., "$(\_ * \_ <= \_ * \_)$", "$Z.abs\ (\_ + \_)$", and "$(\_ + \_ <= \_)$", which are challenging to write for inexperienced Rocq users. In short, the above facts highlight that AUTOROCQ's agentic access to the proof context through invocation of *Search* and other queries is pivotal in its efficacy.

> **Answer to RQ.2**: Each component of AUTOROCQ contributes to its overall effectiveness. AUTOROCQ's agentic context search enhances its efficacy most significantly.

### 6.3 RQ.3: Comparison to Human-written Proofs

Since there is no ground truth for the lemmas extracted from SV-COMP programs, we cannot directly compare the proofs synthesized by AUTOROCQ with human-written proofs. In this subsection, we closely examine a proof synthesized by AUTOROCQ, and compare it to a manually crafted proof by a Rocq expert. The expert has more than six years of formal verification experience and one year of practical experience with Rocq. We use the same theorem `wp_goal` as presented in Figure 1(a), and report the human-written proof in Figure 8.

```
1  Theorem wp_goal :
2    forall (i:Numbers.BinNums.Z) (i1:Numbers.BinNums.Z),
3    let x := (i1 * i1)%Z in ((-9%Z)%Z <= i)%Z →
4    ((-9%Z)%Z <= i)%Z → (i <= 9%Z)%Z → (i1 <= 10%Z)%Z →
5    (i <= 10%Z)%Z → ((-2147483648%Z)%Z <= x)%Z →
6    (100%Z <= x)%Z → (x <= 2147483647%Z)%Z → is_sint32 i1 →
7    is_sint32 i → (i1 = 10%Z).
8
9  (* Helper lemmas written by the human prover*)
10 Lemma square_le_mono_nonneg :
11   forall a b, 0 <= a → 0 <= b → a * a <= b * b → a <= b.
12 Proof. intros a b Ha Hb Hsq. nia. Qed.
13 Lemma square_le_mono_neg :
14   forall a b, a < 0 → b < 0 → a * a <= b * b → b <= a.
15 Proof. intros a b Ha Hb Hsq. nia. Qed.
```

```
1  Proof.
2    intros. subst x. assert (Hcases: 0 <= i1 ∨ i1 < 0) by lia.
3    destruct Hcases as [Hge0 | Hlt0].
4    - assert (H10: 0 <= 10) by lia.
5      pose square_le_mono_nonneg.
6      specialize (l 10 i1 H10 Hge0).
7      simpl in l.
8      specialize (l H5). lia.
9    - pose square_le_mono_neg.
10     specialize (l (-10) i1).
11     assert (Hn10: -10 < 0) by lia.
12     specialize (l Hn10 Hlt0).
13     simpl in l.
14     specialize (l H5). lia.
15 Qed.
```

(a) `wp_goal` in Figure 1(a) with human helper lemmas.   (b) Main proof script written by a human Rocq expert.

Fig. 8. [RQ.3] Proof for `wp_goal` in Figure 1(a), written by a human expert prover in 20 minutes.

To construct the proof in Figure 8(b), the expert had first to recognize that the problem's core reduces to the mathematical principle that, for any integer $a$, $a^2 \geq 100$ implies $|a| \geq 10$. This insight is non-trivial, as it requires moving beyond a direct case-based enumeration of values to grasp the underlying structural inequality. Only after achieving this conceptual leap could the expert elegantly bifurcate the problem based on the sign of `i1` and construct the two custom, symmetric helper lemmas (lines 9–15 in Figure 8(a)) to complete the solution. Crafting this proof takes the Rocq expert 20 minutes of time.

In contrast, let us consider again the proof generated automatically by AutoRocq in lines 9–23 in Figure 1(a). It took AutoRocq only 156.7 seconds to synthesize a proof for the same lemma. At a high level, it achieves the same realization as the human prover, and focuses its proof on an explicit, enumerative case analysis for the value of `i1`. It also quickly connects the goal to properties of absolute values (`Z.abs`). Crucially, AutoRocq leverages two key lemmas from the global context to assist its reasoning on absolute values, namely `Z.abs_le` at line 12 and `Z.square_le_mono_nonneg` at line 19. This not only allows it to prove the goal successfully, but also proves it in noticeably fewer steps than the human expert. Our log reveals that AutoRocq autonomously invokes multiple context searches during its proving. Specifically, searching for the pattern "`(_ * _ <= _ * _)`" surfaces the crucial lemma used at line 19, which eventually helps AutoRocq to succeed. Unsurprisingly, *no* other baseline tools are able to prove this lemma due to their lack of agency.

---

**Answer to RQ.3**: With the help of context queries, AutoRocq is able to prove challenging lemmas, sometimes in even fewer steps than human experts. It also takes much less time.

---

### 6.4 Threats to Validity

**Trustworthiness and Cost of LLMs.** LLMs may give incomplete or factually wrong responses to questions. As a result, verifying programs with untrusted LLMs may seem off-putting at first glance. However, we note that the certificate essentially comes from the trusted kernel of Rocq, which ensures that only derivations conforming to the formal rules of the proof assistant are accepted [46]. Therefore, all the generated proof scripts are *true positives* and are trustworthy, as any erroneous or misleading outputs from the LLMs will simply fail to be verified. Throughout our experiments, it takes $1.22 on average to prove a single lemma without token caching. We deem this cost acceptable, compared to the costly human involvement in lemma-proving activities.

Table 3. [Cast Study] Verifying Linux kernel modules: # of lemmas proved, and the average time/steps (# of tactics) expended on proved lemmas. Across all tools, proofs are constructed for 13 unique lemmas. AutoRocq proves 12 (92%) of them, or 20% of all 60 lemmas.

| Tools | AutoRocq | Rango | PALM | QEDC | P9001 | Combined |
|---|---|---|---|---|---|---|
| Proved? (max. 60) | **12** | 3 | 10 | 3 | 2 | 13 |
| Avg. Time (sec) | 29.9 | 96.2 | 39.4 | **2.0** | 160.8 | – |
| Avg. Steps | **8.3** | 32.0 | – | 32.0 | 979.0 | – |

**Data Leakage.** LLMs are trained on a large corpus of data; thus, they may have been exposed to open-sourced Rocq projects. As such, LLM-based approaches, including AutoRocq, may report inflated results on CoqGym, which comprises only lemmas and proofs in the public domain. In our evaluation, we mitigate this risk by including obligations from SV-COMP programs. These lemmas are extracted by us automatically (see more details in Section 5.1.2), and their ground-truth proofs are not available. In fact, it is unclear if these lemmas are provable at all until we find a correct proof, so we believe the risk of data leakage is minimal.

**Lemma Selection.** In our evaluation, we include established benchmarks and verification targets, namely CoqGym and SV-COMP programs. The lemmas are selected systematically as detailed in Section 5.1.1. We study (un)reachability and overflow freedom as examples of correctness and safety properties, respectively, due to their universality. We also study the loop invariants proposed by LLMs. However, certain types of programs or properties may pose systematic challenges to our approach that we are unaware of. We plan to investigate more types of programs (e.g., multi-threaded) or properties (e.g., memory safety) in future work.

## 7   Case Study: Verifying Linux Kernel Modules

AutoRocq has demonstrated remarkable efficacy in automatically synthesizing rigorous proofs for lemmas related to smaller programs or a mathematical context. In this case study, we assess if AutoRocq can handle the complexity of real-world software. To this end, we build on earlier efforts [18, 49] to reason about properties in the Linux kernel. Specifically, we examine 60 lemmas that formalize functional correctness in Linux kernel code. These lemmas reason about the correctness of intricate program behaviors, and come from various source files, such as memory management (e.g. memmove) and utilities such as string operations (e.g. strcpy) and type conversion (e.g. hex2bin). These lemmas tend to be significantly more involved and may contain up to hundreds of terms.

We report the statistics in Table 3. Results show that AutoRocq is the most effective, being able to synthesize the proofs for 12/60 (20%) lemmas. AutoRocq fails to construct a proof for only one lemma that is proved by *any* other tools (in this case PALM and P9001). In contrast, Rango, QEDC, and P9001 perform poorly on these tasks, which only manage to succeed in 3, 3, 2 lemmas, respectively. We also note that QEDC runs exceptionally fast when it is able to generate a proof, being 15x faster than AutoRocq, which is the second fastest. This suggests that QEDC's search heuristic performs extremely well on certain cases. Nonetheless, such cases seem rare when verifying real, complex software, as it only succeeds on 3 lemmas.

Notably, PALM also achieves remarkable effectiveness and manages to prove 10/60 lemmas. Upon manual examination, we note that *all* successful proofs generated by PALM rely on heavy use of CoqHammer [13], a plug-in that directly discharges remaining subgoals to automated theorem provers (ATPs). CoqHammer is invoked by PALM as a last resort when *none* of its deterministic tactics repairs or resolves the encountered errors. AutoRocq currently does not rely on any

external ATP to boost its performance. We note that more tools, such as with ATPs, could be easily incorporated into our agentic workflow, which we leave as future work.

> AutoRocq can be applied to verify properties in complex software such as the Linux kernel. Access to more tools, e.g., automated theorem provers, may further improve our proof agent.

## 8  Related Work

### 8.1  Formal Verification Tools and Practices

Our work continues the advances on forging an end-to-end verification workflow for software, and thus is closely related to pioneering efforts on verifying critical software systems [10, 25, 29, 36], network protocols [53], and microprocessor designs [23]. Automation in formal verification is largely underpinned by verification-oriented languages including Frama-C [24], Dafny [28], Verus [27], and Viper [38], which integrate specification and proof directly into the programming workflow through the use of annotations. These tool chains generate proof obligations automatically and then attempt to discharge them to automated solvers [6, 16]. When such automation is insufficient, interactive theorem provers (ITPs), such as Isabelle/HOL [39], Lean [15], and Rocq [11], are used as the back end of the workflow. Our approach critically depends on the verification infrastructure.

### 8.2  Machine Learning for Program Verification

Our work builds on a rapidly growing body of work [31] that applies machine learning to program verification. Related approaches include loop invariants inference [44, 45], automatic insertion of annotations [32], and direct synthesis of verified methods [37]. Among these directions, our work is most closely related to automatic approaches in theorem proving. We introduce a new paradigm in this space: *agentic* approaches, in which agents interpret partial proofs and collaborate with a theorem prover to construct complete ones. Our method extends and unifies three broad paradigms of prior work, as detailed below.

**Neural Proof Generation.** A significant body of work treats theorem proving as a sequence-to-sequence (seq2seq) translation task, and thus employs neural networks to learn from large corpora of human-written proofs [19, 41, 51]. While effective at capturing common proof patterns, these networks generate tactics purely based on their parametric knowledge, without a fine-grained, dynamic understanding of the evolving proof state. This limitation is particularly acute in software verification, where proofs are long, the state space is enormous, and the correct next step is highly sensitive to the precise context.

**Retrieval-Augmented Proving.** To augment parametric models with more context, recent works have integrated retrieval mechanisms to fetch relevant premises from a large library of existing theorems during proof generation [33, 48]. This approach is promising because it grounds the generation process in established facts. However, these approaches perform retrievals *statically*, lacking a tight, iterative feedback loop between the proof state, the retriever, and the generator. Consequently, the retriever may fetch axioms that are relevant to the overall theorem but useless for the specific sub-goal the prover is currently tackling.

**Reinforcement Learning.** To better understand the proving progress, approaches like [50] employ reinforcement learning, where each tactic application is awarded with immediate feedback to facilitate incremental progress. The inherent sparsity of positive feedback prompts clever design of reward functions that estimate the value of intermediate states [43]. However, such reward functions often provide noisy signals or overfit to short-term goals, thereby hindering strategic, long-horizon reasoning.

## 9 Perspectives

We present an agentic proving system that automatically generates machine-checkable certificates for proof obligations. At its core, our agent is an autonomous process that retrieves relevant contexts on demand, incorporates feedback from the proof assistant, and generates tactics adaptively – all guided by interpreting the proof derivation tree. Together with a lemma extraction process, our agent can achieve effective push-button verification that takes source code in C (i.e., SV-COMP programs and Linux kernel modules) and automatically generates proofs without any human effort.

As AI-generated code becomes increasingly prevalent in software development, the need for automated verification becomes more pressing. Our work demonstrates that LLM agents can bridge the gap between automatic code generation and formal verification, moving closer to the vision of trusted automatic programming. The agentic nature of our approach suggests a paradigm shift where verification tools act as intelligent collaborators rather than passive validators, capable of autonomously navigating complex proof spaces and adapting to diverse verification challenges. This represents a crucial step toward making formal verification accessible for real-world software systems, where the combination of automated reasoning and agentic intelligence can provide the rigor and scalability needed for next-generation software assurance.

## References

[1] [n. d.]. BusyBox. https://busybox.net/
[2] [n. d.]. Evolved Value Analysis (Eva) Plug-in in Frama-C. https://frama-c.com/fc-plugins/eva.html
[3] [n. d.]. Runtime Error (RTE) Plug-in in Frama-C. https://www.frama-c.com/fc-plugins/rte.html
[4] [n. d.]. Weakest Precondition (WP) Plug-in in Frama-C. https://www.frama-c.com/fc-plugins/wp.html
[5] Wasi Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2021. Unified Pre-training for Program Understanding and Generation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics.
[6] Haniel Barbosa, Clark Barrett, Martin Brain, Gereon Kremer, Hanna Lachnitt, Makai Mann, Abdalrhman Mohamed, Mudathir Mohamed, Aina Niemetz, Andres Nötzli, et al. 2022. cvc5: A versatile and industrial-strength SMT solver. In *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*. Springer, 415–442.
[7] Dirk Beyer and Jan Strejček. 2025. Improvements in Software Verification and Witness Validation: SV-COMP 2025. In *Tools and Algorithms for the Construction and Analysis of Systems*, Arie Gurfinkel and Marijn Heule (Eds.). Springer Nature Switzerland, Cham, 151–186.
[8] Lasse Blaauwbroek, Josef Urban, and Herman Geuvers. 2020. The tactician: A seamless, interactive tactic learner and prover for coq. In *International Conference on Intelligent Computer Mathematics*. Springer, 271–277.
[9] Pedro Carrott, Nuno Saavedra, Kyle Thompson, Sorin Lerner, João F Ferreira, and Emily First. 2024. CoqPyt: Proof Navigation in Python in the Era of LLMs. In *Companion Proceedings of the 32nd ACM International Conference on the Foundations of Software Engineering*. 637–641.
[10] Tej Chajed, Joseph Tassarotti, Mark Theng, M. Frans Kaashoek, and Nickolai Zeldovich. 2022. Verifying the DaisyNFS concurrent and crash-safe file system with sequential reasoning. In *16th USENIX Symposium on Operating Systems Design and Implementation, OSDI 2022, Carlsbad, CA, USA, July 11-13, 2022*, Marcos K. Aguilera and Hakim Weatherspoon (Eds.). USENIX Association, 447–463. https://www.usenix.org/conference/osdi22/presentation/chajed
[11] Projet Coq. 1996. The coq proof assistant-reference manual. *INRIA Rocquencourt and ENS Lyon, version* 5 (1996), 7–1.
[12] Pascal Cuoq, Florent Kirchner, Nikolai Kosmatov, Virgile Prevosto, Julien Signoles, and Boris Yakobowski. 2012. Frama-C: A software analysis perspective. In *International Conference on Software Engineering and Formal Methods (SEFM)*. Springer, 233–247.
[13] Łukasz Czajka and Cezary Kaliszyk. 2018. Hammer for Coq: Automation for dependent type theory. *Journal of automated reasoning* 61, 1 (2018), 423–453.
[14] Leonardo De Moura and Nikolaj Bjørner. 2008. Z3: An efficient SMT solver. In *International conference on Tools and Algorithms for the Construction and Analysis of Systems*. Springer, 337–340.
[15] Leonardo de Moura and Sebastian Ullrich. 2021. The Lean 4 Theorem Prover and Programming Language. In *Automated Deduction - CADE 28 - 28th International Conference on Automated Deduction, Virtual Event, July 12-15, 2021, Proceedings (Lecture Notes in Computer Science, Vol. 12699)*, André Platzer and Geoff Sutcliffe (Eds.). Springer, 625–635. doi:10.1007/978-3-030-79876-5_37

[16] Leonardo Mendonça de Moura and Nikolaj S. Bjørner. 2008. Z3: An Efficient SMT Solver. In *Tools and Algorithms for the Construction and Analysis of Systems, 14th International Conference, TACAS 2008, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2008, Budapest, Hungary, March 29-April 6, 2008. Proceedings (Lecture Notes in Computer Science, Vol. 4963)*, C. R. Ramakrishnan and Jakob Rehof (Eds.). Springer, 337–340. doi:10.1007/978-3-540-78800-3_24

[17] Google DeepMind. [n. d.]. AlphaProof. https://deepmind.google/discover/blog/ai-solves-imo-problems-at-silver-medal-level/

[18] Denis Efremov, Mikhail Mandrykin, and Alexey Khoroshilov. 2018. Deductive verification of unmodified Linux kernel library functions. In *International Symposium on Leveraging Applications of Formal Methods*. Springer, 216–234.

[19] Emily First, Yuriy Brun, and Arjun Guha. 2020. TacTok: Semantics-aware proof synthesis. *Proceedings of the ACM on Programming Languages* 4, OOPSLA (2020), 1–31.

[20] Pranav Garg, Daniel Neider, P. Madhusudan, and Dan Roth. 2016. Learning invariants using decision trees and implication counterexamples. In *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (St. Petersburg, FL, USA) *(POPL '16)*. Association for Computing Machinery, New York, NY, USA, 499–512. doi:10.1145/2837614.2837664

[21] Yeting Ge and Leonardo De Moura. 2009. Complete instantiation for quantified formulas in satisfiabiliby modulo theories. In *Computer Aided Verification: 21st International Conference, CAV 2009, Grenoble, France, June 26-July 2, 2009. Proceedings 21*. Springer, 306–320.

[22] Hendrik. 2025. *Proof tree visualization for Proof General.* https://askra.de/software/prooftree/

[23] Robert B. Jones, John W. O'Leary, Carl-Johan H. Seger, Mark D. Aagaard, and Thomas F. Melham. 2001. Practical Formal Verification in Microprocessor Design. *IEEE Des. Test Comput.* 18, 4 (2001), 16–25. doi:10.1109/54.936245

[24] Florent Kirchner, Nikolai Kosmatov, Virgile Prevosto, Julien Signoles, and Boris Yakobowski. 2015. Frama-C: A software analysis perspective. *Formal Aspects Comput.* 27, 3 (2015), 573–609. doi:10.1007/S00165-014-0326-7

[25] Gerwin Klein, Kevin Elphinstone, Gernot Heiser, June Andronick, David A. Cock, Philip Derrin, Dhammika Elkaduwe, Kai Engelhardt, Rafal Kolanski, Michael Norrish, Thomas Sewell, Harvey Tuch, and Simon Winwood. 2009. seL4: formal verification of an OS kernel.. In *Proceedings of the 22nd ACM Symposium on Operating Systems Principles 2009, SOSP 2009, Big Sky, Montana, USA, October 11-14, 2009*, Jeanna Neefe Matthews and Thomas E. Anderson (Eds.). ACM, 207–220. doi:10.1145/1629575.1629596

[26] Guillaume Lample, Timothee Lacroix, Marie-Anne Lachaux, Aurelien Rodriguez, Amaury Hayat, Thibaut Lavril, Gabriel Ebner, and Xavier Martinet. 2022. Hypertree proof search for neural theorem proving. *Advances in neural information processing systems* 35 (2022), 26337–26349.

[27] Andrea Lattuada, Travis Hance, Chanhee Cho, Matthias Brun, Isitha Subasinghe, Yi Zhou, Jon Howell, Bryan Parno, and Chris Hawblitzel. 2023. Verus: Verifying Rust Programs using Linear Ghost Types (extended version). *CoRR* abs/2303.05491 (2023). arXiv:2303.05491 doi:10.48550/ARXIV.2303.05491

[28] K. Rustan M. Leino. 2010. Dafny: An Automatic Program Verifier for Functional Correctness. In *Logic for Programming, Artificial Intelligence, and Reasoning - 16th International Conference, LPAR-16, Dakar, Senegal, April 25-May 1, 2010, Revised Selected Papers (Lecture Notes in Computer Science, Vol. 6355)*, Edmund M. Clarke and Andrei Voronkov (Eds.). Springer, 348–370. doi:10.1007/978-3-642-17511-4_20

[29] Xavier Leroy, Sandrine Blazy, Daniel Kästner, Bernhard Schommer, Markus Pister, and Christian Ferdinand. 2016. CompCert-a formally verified optimizing compiler. In *ERTS 2016: Embedded Real Time Software and Systems, 8th European Congress*.

[30] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems* 33 (2020), 9459–9474.

[31] Zhaoyu Li, Jialiang Sun, Logan Murphy, Qidong Su, Zenan Li, Xian Zhang, Kaiyu Yang, and Xujie Si. [n. d.]. A Survey on Deep Learning for Theorem Proving. In *First Conference on Language Modeling*.

[32] Chloe Loughridge, Qinyi Sun, Seth Ahrenbach, Federico Cassano, Chuyue Sun, Ying Sheng, Anish Mudide, Md Rakib Hossain Misu, Nada Amin, and Max Tegmark. 2025. DafnyBench: A Benchmark for Formal Software Verification. *Trans. Mach. Learn. Res.* 2025 (2025). https://openreview.net/forum?id=yBgTVWccIx

[33] Minghai Lu, Benjamin Delaware, and Tianyi Zhang. 2024. Proof automation with large language models. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*. 1509–1520.

[34] Minghai Lu, Benjamin Delaware, and Tianyi Zhang. 2024. Proof automation with large language models. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*. 1509–1520.

[35] Lezhi Ma, Shangqing Liu, Yi Li, Xiaofei Xie, and Lei Bu. 2025. SpecGen: Automated Generation of Formal Program Specifications via Large Language Models . In *2025 IEEE/ACM 47th International Conference on Software Engineering (ICSE)*. IEEE Computer Society, Los Alamitos, CA, USA, 16–28. doi:10.1109/ICSE55347.2025.00129

[36] Gregory Malecha, Greg Morrisett, Avraham Shinnar, and Ryan Wisnesky. 2010. Toward a verified relational database management system. In *Proceedings of the 37th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (Madrid, Spain) *(POPL '10)*. Association for Computing Machinery, New York, NY, USA, 237–248. doi:10.1145/1706299.1706329

[37] Md Rakib Hossain Misu, Cristina V. Lopes, Iris Ma, and James Noble. 2024. Towards AI-Assisted Synthesis of Verified Dafny Methods. *Proc. ACM Softw. Eng.* 1, FSE (2024), 812–835. doi:10.1145/3643763

[38] Peter Müller, Malte Schwerhoff, and Alexander J. Summers. 2016. Viper: A Verification Infrastructure for Permission-Based Reasoning. In *Verification, Model Checking, and Abstract Interpretation - 17th International Conference, VMCAI 2016, St. Petersburg, FL, USA, January 17-19, 2016. Proceedings (Lecture Notes in Computer Science, Vol. 9583)*, Barbara Jobstmann and K. Rustan M. Leino (Eds.). Springer, 41–62. doi:10.1007/978-3-662-49122-5_2

[39] Lawrence C Paulson. 1994. *Isabelle: A generic theorem prover*. Springer.

[40] Hammond Pearce, Baleegh Ahmad, Benjamin Tan, Brendan Dolan-Gavitt, and Ramesh Karri. 2025. Asleep at the keyboard? assessing the security of github copilot's code contributions. *Commun. ACM* 68, 2 (2025), 96–105.

[41] Alex Sanchez-Stern, Yousef Alhessi, Lawrence Saul, and Sorin Lerner. 2020. Generating correctness proofs with neural networks. In *Proceedings of the 4th ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*. 1–10.

[42] Alex Sanchez-Stern, Emily First, Timothy Zhou, Zhanna Kaufman, Yuriy Brun, and Talia Ringer. 2023. Passport: Improving automated formal verification using identifiers. *ACM Transactions on Programming Languages and Systems* 45, 2 (2023), 1–30.

[43] Alex Sanchez-Stern, Abhishek Varghese, Zhanna Kaufman, Dylan Zhang, Talia Ringer, and Yuriy Brun. 2024. QED-Cartographer: Automating formal verification using reward-free reinforcement learning. In *2025 IEEE/ACM 47th International Conference on Software Engineering (ICSE)*. IEEE Computer Society, 405–418.

[44] Xujie Si, Hanjun Dai, Mukund Raghothaman, Mayur Naik, and Le Song. 2018. Learning loop invariants for program verification. *Advances in Neural Information Processing Systems* 31 (2018).

[45] Xujie Si, Aaditya Naik, Hanjun Dai, Mayur Naik, and Le Song. 2020. Code2Inv: A Deep Learning Framework for Program Verification. In *Computer Aided Verification - 32nd International Conference, CAV 2020, Los Angeles, CA, USA, July 21-24, 2020, Proceedings, Part II (Lecture Notes in Computer Science, Vol. 12225)*, Shuvendu K. Lahiri and Chao Wang (Eds.). Springer, 151–164. doi:10.1007/978-3-030-53291-8_9

[46] Matthieu Sozeau, Yannick Forster, Meven Lennon-Bertrand, Jakob Nielsen, Nicolas Tabareau, and Théo Winterhalter. 2025. Correct and Complete Type Checking and Certified Erasure for Coq, in Coq. *J. ACM* 72, 1, Article 8 (Jan. 2025), 74 pages. doi:10.1145/3706056

[47] Rocq Development Team. 2025. *Search Comand in Rocq Documentation.* https://rocq-prover.org/doc/V9.0.0/refman/proof-engine/vernacular-commands.html#coq:cmd.Search

[48] Kyle Thompson, Nuno Saavedra, Pedro Carrott, Kevin Fisher, Alex Sanchez-Stern, Yuriy Brun, Joao F. Ferreira, Sorin Lerner, and Emily First. 2025. Rango: Adaptive Retrieval-Augmented Proving for Automated Software Verification . In *2025 IEEE/ACM 47th International Conference on Software Engineering (ICSE)*. IEEE Computer Society, Los Alamitos, CA, USA, 347–359. doi:10.1109/ICSE55347.2025.00161

[49] Grigoriy Volkov, Mikhail Mandrykin, and Denis Efremov. 2018. Lemma functions for Frama-C: C programs as proofs. In *2018 Ivannikov Ispras Open Conference (ISPRAS)*. IEEE, 31–38.

[50] Minchao Wu, Michael Norrish, Christian Walder, and Amir Dezfouli. 2021. TacticZero: Learning to Prove Theorems from Scratch with Deep Reinforcement Learning. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (Eds.). 9330–9342. https://proceedings.neurips.cc/paper/2021/hash/4dea382d82666332fb564f2e711cbc71-Abstract.html

[51] Kaiyu Yang and Jia Deng. 2019. Learning to prove theorems via interacting with proof assistants. In *International Conference on Machine Learning*. PMLR, 6984–6994.

[52] Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. 2024. Autocoderover: Autonomous program improvement. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*. 1592–1604.

[53] Naiqian Zheng, Mengqi Liu, Yuxing Xiang, Linjian Song, Dong Li, Feng Han, Nan Wang, Yong Ma, Zhuo Liang, Dennis Cai, Ennan Zhai, Xuanzhe Liu, and Xin Jin. 2023. Automated Verification of an In-Production DNS Authoritative Engine. In *Proceedings of the 29th Symposium on Operating Systems Principles, SOSP 2023, Koblenz, Germany, October 23-26, 2023*, Jason Flinn, Margo I. Seltzer, Peter Druschel, Antoine Kaufmann, and Jonathan Mace (Eds.). ACM, 80–95. doi:10.1145/3600006.3613153