

SEVerA: Verified Synthesis of Self-Evolving Agents

DEBANGSHU BANERJEE*, CHANGMING XU*, and GAGANDEEP SINGH, University of Illinois Urbana-Champaign, USA

Recent advances have demonstrated the effectiveness of self-evolving LLM agents on tasks such as program repair and scientific discovery. In this paradigm, a planner LLM synthesizes an agent program that invokes parametric models, including probabilistic generative models such as LLMs, smaller neural networks, and external tools such as SMT solvers. These components are then tuned per task to improve performance. However, unlike traditional constraint-guided program synthesis, existing self-evolving agent frameworks provide no formal guarantees of safety or correctness. Because such synthesized programs are often executed autonomously on unseen inputs, the lack of formal guarantees raises serious reliability and security concerns. To address this gap, we formulate agentic code generation as a constrained learning problem that combines hard formal specifications with soft objectives capturing task utility. We introduce Formally Guarded Generative Models (FGGM), which allow the planner LLM to specify a formal output contract for each generative-model call using first-order logic. Each FGGM call automatically wraps the underlying parametric generative model in a rejection sampler with a verified fallback, treating model outputs as samples from a proposal distribution. As a result, every returned output satisfies the specified contract for any input and any parameter setting of the underlying model. Building on FGGM, we present SEVerA (Self-Evolving Verified Agents), a three-stage framework for solving constrained learning problems arising from agent synthesis. In *Search*, the planner LLM synthesizes candidate parametric programs that may contain multiple FGGM calls. In *Verification*, we prove correctness with respect to the hard constraints for all parameter values, reducing the problem to unconstrained learning. In *Learning*, we apply scalable gradient-based optimization, including GRPO-style fine-tuning for LLMs, to improve the soft objective while preserving formal correctness. We evaluate SEVerA on constrained symbolic regression, invariant generation for Dafny programs, symbolic mathematical expression synthesis, and policy-compliant agentic tool use (τ^2 -bench). Across all tasks, SEVerA achieves zero constraint violations while simultaneously improving task performance over unconstrained and state-of-the-art baselines. Our results demonstrate that formal behavioral constraints not only guarantee correctness but also prune the space of candidate programs, steering synthesis toward higher-quality agents.

CCS Concepts: • **LLM Agents** → **Automated Verification**.

Additional Key Words and Phrases: LLM Agents, Automated Verification, Deductive Program Synthesis.

ACM Reference Format:

Debangshu Banerjee, Changming Xu, and Gagandeep Singh. 2026. SEVerA: Verified Synthesis of Self-Evolving Agents. 1, 1 (March 2026), 42 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Recent works have leveraged the code-generation capabilities of large language models (LLMs) to automatically synthesize LLM-based agents, expressed as programs, from natural language task

*These authors contributed equally to this work.

Authors' address: Debangshu Banerjee, db21@illinois.edu; Changming Xu, cx23@illinois.edu; Gagandeep Singh, ggnds@illinois.edu, University of Illinois Urbana-Champaign, USA.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2026 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM XXXX-XXXX/2026/3-ART

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

descriptions [13, 19, 47, 49, 53]. Given a user-defined task, a planner LLM generates a program that invokes parametric models, including LLMs, smaller neural networks, and external tools such as SMT solvers, and executes this program on user inputs to compute outputs. While this paradigm is powerful, it introduces serious safety and reliability concerns: the synthesized programs are executed autonomously on unseen inputs. Modern self-evolving frameworks [11, 43, 48, 51] further complicate formal safety verification as they fine-tune the parameters of models embedded in the synthesized programs to improve task-specific performance. Consequently, any correctness guarantee provided on an agentic program must continue to hold after parameter updates.

The absence of formal safety verification leads to failures across diverse domains. In program verification, agents cheat by subtly modifying the input program (e.g., altering variable initializations) so that the modified program with proposed annotations passes verification, inflating task accuracy [3]. In code repair, agents delete failing tests rather than fixing the underlying bugs [50]. In agentic tool use, unconstrained agents violate domain-specific policies, such as refund eligibility and booking-modification rules, on 65-76% of interactions [20]. In deployment, rogue agents have been observed bypassing security protections in the wild [14, 15, 25]. These are not isolated incidents; they are a consequence of evaluating synthesized agents solely on soft performance metrics without formal behavioral specifications. Natural-language task descriptions and testing on fixed inputs is not sufficient to completely prevent such failures. Agentic programs need *formal hard constraints* to ensure safety. In this paper, we aim to achieve both the formal guarantees of constraint-guided synthesis [1] and the flexibility and performance benefits of self-evolving agentic frameworks.

Key Challenges: We summarize the challenges below.

- **Balancing Safety and Performance:** A safe agent synthesis algorithm must ensure formal correctness guarantees while simultaneously improving task-specific performance. Existing work on deductive program synthesis [1, 22, 38] provides formal guarantees but typically does not optimize task-specific performance objectives. Meanwhile, popular gradient-based optimization methods for LLMs, such as GRPO [36], empirically improve output quality but provide no assurance that post-training outputs satisfy specified constraints.
- **Enforcing Constraints per Model Call:** Generative models (GMs), such as LLMs, may be used in an agentic program at different places with different constraints. For example, one call may require generating Python programs, while another may require producing symbolic mathematical expressions conforming to a specific formal grammar. Thus, we require a mechanism to specify formal output constraints for *each* GM call in an agentic program and a way to use these specifications to verify that an agentic program satisfies the user-provided behavioral specification. Recent work on constrained decoding [4, 35, 45] can enforce context-free grammar-based syntactic constraints and limited semantic constraints such as variable naming or simple type checks. However, these approaches are restricted in the types of constraints they support and require modifying the model’s decoding procedure, which limits their applicability to open-source models. Moreover, constrained decoding strategies are known to distort the output distribution, producing outputs that satisfy constraints but degrade task performance [4, 35].

Our Contributions: We summarize our contributions below.

- We introduce the novel concept of Formally Guarded Generative Models (FGGM), which allow us to define and enforce local contracts on model calls through formal input-output specifications expressed in first-order logic. FGGM operates solely on LLM output strings and can be applied to both open and closed source models. Using FGGM, we retain the formal guarantees of deductive synthesis while utilizing the performance gained through gradient-based parameter optimization.
- We introduce SEVerA: the first self-evolving agent synthesis algorithm with verifiable guarantees. The approach consists of three stages: (a) *Search*, where a planner LLM samples program strings in

a verifier-aware language such as Dafny, and use the proposed FGGMs to establish local contracts on the embedded model calls; (b) *Verify*, where the language’s built-in verifier is used to prove that the sampled programs satisfy all specified contracts over all parameters of the embedded parametric model calls; and (c) *Learn*, where verified parametric programs reduce the constrained learning problem to an unconstrained optimization problem over model parameters, enabling scalable gradient-based training without sacrificing correctness.

- We prove that SEVerA is *sound*: any agent returned satisfies the behavioral specification for all inputs and all parameter values (Theorem 5.4). We also establish a *sufficient condition* under which a verified agent exists that satisfies the hard constraints while incurring no greater task loss than any unconstrained generative model with initial parameters, with strict improvement whenever the unconstrained model violates the specification (Theorem 5.5).
- We evaluate SEVerA on tasks spanning LLM-assisted program verification, symbolic math synthesis, agentic tool use, and constrained symbolic regression. Across all benchmarks, SEVerA achieves provably zero constraint violations while outperforming state-of-the-art agents in task performance: 97.0% verification rate on HumanEvalDafny (vs. 86.9% for the best baseline), 66.0% accuracy on GSM-Symbolic (vs. 44.7% for the best constrained-decoding method), and 52.6% pass rate τ^2 -bench’s airline domain using Qwen3-8B. Notably, on τ^2 -bench airline SEVerA even beats Agent-C [20] with Claude Sonnet 4.5, a state-of-the-art constrained agent using a frontier model. These results demonstrate that behavioral constraints do not merely enforce safety but actively prune the search space of candidate programs and steer synthesis toward higher-quality agents. The code is available on github.

2 BACKGROUND

Notations and Terminology: We use Σ to denote the alphabet, a finite set of characters. All LLMs output finite strings over Σ . We use lowercase letters (x) to denote constants, bold letters (\mathbf{x}) to denote strings, capital letters (X) to denote sets (including sets of functions), and $|\mathbf{x}|$ to denote string length. In the rest of the paper, we treat neural networks as a restricted form of generative model that produces a single output for a fixed input.

LLMs: LLMs $\mathcal{L} : \Sigma^* \rightarrow \Sigma^*$ are probabilistic string generators that, given a prompt $\mathbf{p} \in \Sigma^*$ over the alphabet Σ , sample an output string $\mathbf{y} \in \Sigma^*$. Typically, the output length $|\mathbf{y}|$ is bounded by a fixed user-specified n , i.e., $|\mathbf{y}| \leq n$ or equivalently $\mathbf{y} \in \Sigma^{\leq n}$. Fixed-length generation $\mathcal{L}_n : \Sigma^* \rightarrow \Sigma^{\leq n}$ is implemented by iteratively composing two high-level steps, repeated at most n times. Each iteration samples a single character (sometimes referred to as a token) from Σ or the special end-of-sequence symbol $\langle eos \rangle \notin \Sigma$ that marks termination of generation. The first step is the **distribution prediction step**, defined as $\mathcal{L}^p : \Sigma^* \rightarrow \mathcal{P}(\Sigma \cup \langle eos \rangle)$, which maps a prefix string $\mathbf{x} \in \Sigma^*$ to a probability distribution over $\Sigma \cup \langle eos \rangle$. This distribution captures how likely each character is to appear next. The second step is the **decoding step**, defined as $\mathcal{D} : \Sigma^* \times \mathcal{P}(\Sigma \cup \langle eos \rangle) \rightarrow \Sigma \cup \langle eos \rangle$, which draws the next character using the current prefix \mathbf{x} and the predicted distribution. If the decoding step produces $\langle eos \rangle$, the generation terminates. Otherwise, the sampled character c is appended to the prefix, and generation continues. Henceforth, we drop the subscript n for simplicity.

Definition 2.1 (LLMs). For a fixed $n \in \mathbb{N}$, an LLM is a bounded string generator $\mathcal{L}_n : \Sigma^* \rightarrow \Sigma^{\leq n}$ that samples an output string $\mathbf{y} \in \Sigma^{\leq n}$ for input \mathbf{x} . \mathbf{y} on \mathbf{x} is computed recursively as $\mathbf{y}_0 = \mathbf{x}$, and $\mathbf{y}_i = \mathbf{y}_{i-1} \cdot c_i$ if $(\forall j \leq i. c_j \neq \langle eos \rangle)$ else \mathbf{y}_{i-1} where $i \in [n]$, $c_i = \mathcal{D}(\mathbf{y}_{i-1}, \mathcal{L}^p(\mathbf{y}_{i-1}))$ and $\mathbf{x} \cdot \mathbf{y} = \mathbf{y}_n$.

Rejection Sampling: Rejection sampling is a popular Monte Carlo method for generating samples from a *target distribution* π_t when direct sampling from π_t is difficult. Instead, it relies on a *proposal distribution* π_p from which sampling is easy. The support set $S(\pi_p)$ (see Definition 2.2) must cover $S(\pi_t)$, i.e., $S(\pi_t) \subseteq S(\pi_p)$. Intuitively, π_p generates candidate samples, and the rejection sampler

decides whether to accept each candidate, rejecting those that are unlikely under π_t (details in Appendix B). The rejection sampler guarantees that no accepted sample lies outside $S(\pi_t)$.

Definition 2.2. Let π be a probability distribution over Ω with density function $D_\pi : \Omega \rightarrow \mathbb{R}^+$. The support set $S(\pi)$ of π captures all points with positive probability $S(\pi) = \{x \in \Omega; |D_\pi(x) > 0\}$.

Self-evolving LLM agents: In this work, we focus on LLM agents modeled as programs $f : T_i \rightarrow T_o$ that, given any user input $x \in T_i$, compute the output $y = f(x)$. Typically, f is written in a popular imperative language such as Python. The program invokes one or more parametric models (formally defined in § 5.1), such as LLMs, as well as other tools such as an SMT solver, to compute the output y . All parametric models \mathcal{F}_p and tool calls \mathcal{F}_c are provided as a collection of pre-implemented library functions. We denote this library by $\mathcal{F} = \mathcal{F}_p \cup \mathcal{F}_c$ (formal definition in § 5.1). In the self-evolving paradigm, the program f is not handwritten. Instead, it is generated by a planner LLM based on task-specific instructions provided as a prompt. The goal is to synthesize a program f that, given a training dataset of input and optional ground-truth output pairs $D \subseteq T_i \times T_o$ and a loss function $L : T_i \times T_o \times T_o \rightarrow \mathbb{R}^+$, minimizes the aggregated loss $\frac{1}{|D|} \sum_{(x_i, y_i) \in D} L(x_i, y_i, f(x_i))$. For certain tasks, the ground truth y_i may not be available; in such cases, the loss is denoted as $L(x_i, _, f(x_i))$. Assume that a context-free grammar (CFG) G defines the set of syntactically valid programs f as $L(G) \subseteq \Sigma^*$. Given the library functions \mathcal{F} , let $S(G, \mathcal{F})$ denote the search space for $f : T_i \rightarrow T_o$ with type signature $(T_i \rightarrow T_o)$, including all parametric values of the models in \mathcal{F}_p . Self-evolving agentic frameworks aim to find the optimal solution f^* to the following optimization problem.

$$f^* = \arg \min_{f \in S(G, \mathcal{F})} \frac{1}{|D|} \times \sum_{(x_i, _) \in D} L(x_i, _, f(x_i)) \quad (1)$$

3 PROBLEM FORMULATION

High-level formulation of agent synthesis with formal constraints: Eq. 1 searches only for the optimal program f^* that minimizes the loss on the training set D . However, it provides no guarantee about how f^* performs on unseen user inputs outside D . In contrast, deductive program synthesis [1, 38] allows program generation to be controlled by formal behavioral input $\Phi : T_i \rightarrow \{T, F\}$ and output specifications $\Psi : T_i \times T_o \rightarrow \{T, F\}$. The use of formal behavioral specifications ❶ enables proving the correctness of the synthesized program for all inputs that satisfy the input specification, and ❷ can guide the search for f^* by eliminating unverified program candidates. With (Φ, Ψ) , the unconstrained learning problem in Eq. 1 can be reformulated as the constrained learning problem shown below.

$$f^* = \arg \min_{f \in S(G, \mathcal{F})} \overbrace{\frac{1}{|D|} \times \sum_{(x_i, _) \in D} L(x_i, _, f(x_i))}^{\text{soft learning objective}} \quad \text{s.t.} \quad \overbrace{\forall x \in T_i. \Phi(x) \implies \Psi(x, f(x))}^{\text{hard formal constraints}} \quad (2)$$

The central question is whether Eq. 2 can be solved while retaining the performance benefits and scalability of gradient-based methods designed for Eq. 1. Next, we illustrate how practical constraints from several application domains can be encoded as behavioral specifications (Φ, Ψ) .

Examples of formal constraints (Φ, Ψ) : The output specification Ψ characterizes the agent's expected behavior on any valid input, not just on examples in D . Hard constraints therefore detect and mitigate errors during synthesis that would not be revealed by performance evaluation on a fixed dataset D . We show encodings of (Φ, Ψ) for four different tasks from distinct domains. ❶ **Scientific Discovery:** We consider constrained symbolic regression, where the task is to recover an unknown mathematical formula from observed data. The formal constraints (Φ, Ψ) encode prior knowledge about the target formula, such as known symbolic bounds [7, 17, 28] or asymptotic behavior [27]. Any generated formula that violates these hard constraints is invalid, regardless of its empirical fit to the data. Moreover, in the presence of noise, the constraints (Φ, Ψ) help prevent overfitting, which cannot be addressed solely by optimizing a soft objective [24]. ❷ **Program Verification:** Given a

program in a verification-aware language such as Dafny with a predefined specification, the agent is expected to synthesize annotations (e.g., loop invariants, assertions, ranking functions) [29, 31, 39]. Task success depends on whether the annotated program verifies. Recent work [3] shows that agents using frontier models (e.g., Claude) may cheat by subtly modifying the input program, such as altering variable initializations, instead of producing correct annotations. [3] reports that such modifications were not detected by prior string-based checkers, leading to inflated task accuracy. In contrast, formal AST-based diff checkers that syntactically enforce equivalence between the input program and the annotated output serve as hard constraints, ensuring the agent cannot cheat even on unseen inputs.

③ **Constrained LLM Generation:** Constrained generation captures a class of problems in which LLMs are required to generate strings that follow a formal structure. Typically, these formal structures are defined using formal grammars (regular [40] or context-free grammars [35, 45]), or static analyzers such as type checkers [32, 33]. In our setting, these formal structures define the output specification Ψ . We consider the task of generating symbolic math expressions from natural language questions, where hard constraints ensure that the synthesized agent always produces at least a structurally valid expression.

④ **Agentic Tool Use:** Conversational LLM agents deployed in customer-service settings must select and invoke API tools to resolve user requests while respecting domain-specific policies, such as refund eligibility, booking-modification rules, and authentication-before-access requirements [5]. These temporal and logical policy constraints can be expressed as formal specifications that govern the ordering and content of agent actions. Agent-C [20] provides a domain-specific language for specifying such temporal properties and translates them into linear temporal logic (LTL), enabling SMT-based checking of the generated tool calls. Crucially, the hard constraint is enforced at each individual tool call: given the current tool-call trace, the Agent-C checker verifies whether the next proposed tool call complies with the defined formal policy before it is executed. Task performance (utility), in contrast, is measured over the entire multi-turn interaction trace as whether the agent successfully resolves the user request. In our formulation, Ψ uses the Agent-C checker to define a hard constraint applied to each tool call, ensuring that synthesized agents never violate domain policies at any step, regardless of user input. Overall, hard behavioral constraints capture the minimal requirements that synthesized agents must satisfy on all inputs, thereby enabling the pruning of untrusted candidates during synthesis.

4 OVERVIEW

Fig. 1 gives a high-level overview of SEVerA for solving the constrained learning problem in Eq. 2 over the search space $S(G, \mathcal{F})$. In this setting, we need to simultaneously optimize task-specific performance while guaranteeing that the resulting program provably satisfies the given behavioral constraints. This is difficult because the search space contains both discrete program structure and continuous model parameters, and because although GMs such as LLMs are powerful, they are inherently unreliable, which makes it hard to ensure conformance with behavioral specifications. To address these challenges, SEVerA retains the best of both worlds: constraint-guided program search from deductive synthesis and unconstrained gradient-based parameter optimization. Specifically, SEVerA employs a CEGIS-style loop in which a planner LLM proposes candidate parametric programs and a verifier checks their correctness with respect to the behavioral specification. Once a parametric program is verified, SEVerA applies unconstrained gradient-based optimization to tune the underlying GM parameters and reduce the task-specific loss L . The proposed Formally Guarded Generative Model (FGGM) serves as the critical link between these two strategies. The FGGM setup allows one to bind each GM call to explicit local input–output contracts and a verified non-parametric fallback program, ensuring that these contracts hold regardless of the underlying GM parameters. The verifier then leverages these local contracts to prove the correctness of the parametric program with respect to the behavioral specification. Beyond providing correctness

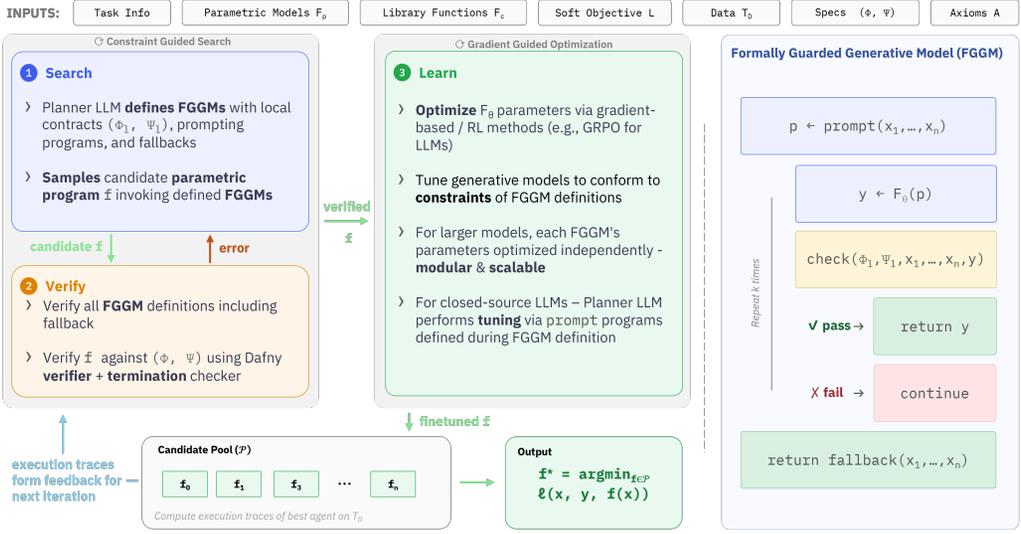


Fig. 1. Overview of SEVerA, which operates in three key steps. (1) **Search**: Given the task information, library functions \mathcal{F}_c , a list of parametric generative models (GMs) \mathcal{F}_p , and specifications (Φ, Ψ) , the planner LLM outputs a parametric agentic program in which all formal output specifications of GM calls are defined using the FGGM setup. (2) **Verify**: Given the parametric agentic code, this step first verifies all FGGM definitions proposed by the planner. Once verified, it uses their local contracts to check the program against (Φ, Ψ) . If all checks pass, the program is accepted and the unconstrained learning step is invoked on the verified program; otherwise, an error message is returned to the planner. (3) **Learn**: This step optimizes the parameters of the underlying GMs within each FGGM to improve conformance with local contracts defined in step (1), while also reducing the task-loss L over the dataset D . After tuning, SEVerA maintains a pool of verified fine-tuned agents and uses their execution traces on D to generate new candidates via step (1). All agents in the pool are valid candidate solutions, and the one with the lowest L on D is returned. Central to SEVerA is the Formally Guarded Generative Model (FGGM), which binds each GM call to local input–output contracts (Φ_l, Ψ_l) and a verified non-parametric fallback, ensuring that the contracts hold irrespective of the underlying GM parameters. This enables unconstrained gradient-based parameter optimization without compromising formal correctness.

guarantees, these contracts also extract local learning objectives by characterizing the expected outputs of each GM call. This enables gradient-based optimization to improve the GM’s conformance to its local contracts, thereby connecting constrained program search with unconstrained learning. Next, we discuss the components of SEVerA through two representative and easy-to-understand problem instances from different domains: scientific discovery and program verification.

4.1 Example Instantiation of SEVerA

4.1.1 Constrained Symbolic Regression: In symbolic regression, each problem instance consists of possibly noisy observations of input–output pairs generated by an unknown ground-truth function f_{gt} . The goal is to recover this ground-truth function from the training data D , a popular task in scientific discovery and neuro-symbolic program synthesis [10, 37]. In constrained symbolic regression [7, 17], users encode known properties of f_{gt} beyond what is captured in D as formal constraints, ensuring that the recovered function f satisfies these constraints. This setup directly follows the constrained learning formulation in Eq. 2, where the synthesized agent f^* corresponds to the recovered function. For illustration, we consider the problem instance with ground-truth

function $f_{gt}(x) = \sqrt{1.23 \times \max(x, 0.0)}$, taken from [28]. The dataset D contains input–output pairs $(x, y) \in \mathbb{R}^2$, where $y = f_{gt}(x) + \epsilon$ is a noisy observation of $f_{gt}(x)$ and $\epsilon \sim \mathcal{N}(0, \sigma)$ denotes additive Gaussian noise. The pointwise loss function L is the commonly used Normalized Mean Squared Error (NMSE) (see Eq. 3). For this instance, the agent’s type signature τ is $(\mathbb{R} \rightarrow \mathbb{R})$, and the behavioral specifications (Φ, Ψ) define known symbolic bounds on f_{gt} (see Eq. 4) where $\text{sqr}t : \mathbb{R} \rightarrow \mathbb{R}$ is a provided library function.

$$L(x_i, y_i, f(x_i)) = \frac{(f(x_i) - y_i)^2}{C} \quad \text{where } C = \sum_{(x_i, y_i) \in D} y_i^2 \quad (3)$$

$$\Phi(x) : (x \geq 0), \Psi(x, f(x)) : ((x \leq 1) \implies (f(x) \geq \text{pow}(x, 0.8)) \wedge ((x \geq 1) \implies (f(x) \geq \text{sqr}t(x))) \quad (4)$$

Search Space: SEVerA defines the search space using a restricted subset of Dafny supporting four basic types $\{\text{int}, \text{bool}, \text{real}, \text{str}\}$, conditional blocks, while loops, and function calls. We formalize the search space in § 5.1.

Library Functions: For this task, the library function set \mathcal{F} includes common transcendental functions (e.g., sin , cos , pow , $\text{sqr}t$) as well as small parametric neural networks (NNs), enabling the search space to capture a wide range of neuro-symbolic programs. Users provide type signatures and formal contracts for all non-parametric library functions; the planner LLM defines formal contracts for each parametric model call through the FGGM setup (§ 4.2). We present representative type signatures and specifications in Eq. 5, with the complete set in Appendix C.

$$\text{sqr}t(x : \text{real}) \rightarrow \text{real}, \Phi_{\text{sqr}t}(x) : (x \geq 0), \Psi_{\text{sqr}t}(x, \text{sqr}t(x)) : (\text{sqr}t(x) \geq 0); \text{NN}_{\theta}^{(1)}(x_1 : \text{real}) \rightarrow \text{real} \quad (5)$$

$$\forall x, d_1, d_2 \in \mathbb{R}. (x \geq 1) \wedge (d_1 \geq d_2) \implies \text{pow}(x, d_1) \geq \text{pow}(x, d_2) \quad (6)$$

$$\forall x, d_1, d_2 \in \mathbb{R}. (x \leq 1) \wedge (d_1 \geq d_2) \implies \text{pow}(x, d_1) \leq \text{pow}(x, d_2) \quad (7)$$

$$\varphi_{\mathcal{A}} \implies \forall x \in T_i. (\Phi(x) \implies \Psi(x, f(x))) \quad \text{where } \varphi_{\mathcal{A}} = (\bigwedge_{\varphi \in \mathcal{A}} \varphi) \quad (8)$$

Axioms: The axioms \mathcal{A} encode known properties of library functions in \mathcal{F} , including their formal contracts as universally quantified first-order sentences, where library functions are treated as uninterpreted. For example, Eq. 7 captures a property of pow . These axioms assist automated verification as shown in Eq. 8. We provide the formal syntax for defining \mathcal{A} in § 5.1 and the complete axiom set for this task in Appendix D.

Like other synthesis frameworks with library functions [12], SEVerA assumes that all library functions have correct input–output specifications, and that they are pure and terminating. Verifying the correctness of the library functions and axioms is outside the scope of SEVerA.

4.1.2 LLM-Assisted Automated Verification. This task aims to synthesize annotations, such as inductive loop invariants, ranking functions, and assertions, that enable automatic verification and termination checking of input Dafny programs with pre-encoded input–output specifications [29, 39]. The agent has type signature $(p : \text{str}) \rightarrow \text{str}$, mapping an input Dafny program with pre-encoded specifications, represented as a string, to an annotated Dafny program. The point-wise loss function L evaluates whether the annotated program verifies within a pre-fixed time budget, assigning 1 to unverified programs and 0 to verified ones (Eq. 10). Hard constraints (Φ, Ψ) ensure that, for any parsable input program, the output remains parsable and does not change the input program beyond adding annotations (Eq. 11). We use the Dafny built-in verifier to define L , and the Dafny parser together with the AST-based formal diff checker from [3] to define (Φ, Ψ) . noDiff returns True if the second input program is syntactically equivalent to the first (ignoring annotations), and False otherwise. The library function set \mathcal{F} includes all these components. Although parsing and diff checks could be incorporated into L by penalizing invalid outputs on D , these constraints must hold for all possible inputs, not only those in D . Therefore, any candidate agent that fails to satisfy them must be filtered out.

Library Functions: We use the same search space as constrained symbolic regression. \mathcal{F} includes the verifier with a pre-fixed time budget *verify*, the Dafny parser *parse*, the diff-checker *noDiff*. We show the type signature and formal contract of representative functions in Eq. 9 while providing the rest in Appendix E.

$$\text{noDiff}(\mathbf{x} : \text{str}, \mathbf{y} : \text{str}) \rightarrow \text{bool}, \Phi(\mathbf{x}, \mathbf{y}) : \text{parse}(\mathbf{x}), \Psi(\mathbf{x}, \mathbf{y}, \text{noDiff}(\mathbf{x}, \mathbf{y})) : T \quad (9)$$

$$\mathbb{L}(\mathbf{x}, _ , \mathbf{y}) = 1 - \mathbb{I}(\text{parse}(\mathbf{y}) \wedge \text{noDiff}(\mathbf{x}, \mathbf{y}) \wedge \text{verify}(\mathbf{y})) \quad (10)$$

$$\Phi(\mathbf{x}) : \text{parse}(\mathbf{x}); \Psi(\mathbf{x}, f(\mathbf{x})) : \text{parse}(f(\mathbf{x})) \wedge \text{noDiff}(\mathbf{x}, f(\mathbf{x})) \quad (11)$$

Axioms: Example of a first-order sentence encoding known properties of *noDiff* in Eq. 12. *noDiff* is both reflexive and transitive w.r.t all parsable Dafny programs. It is easy to see that same program will always have no difference with itself. We list all the axioms in Appendix F.

$$\text{Reflexivity: } \forall \mathbf{x} \in \Sigma^*. \text{parse}(\mathbf{x}) \implies \text{noDiff}(\mathbf{x}, \mathbf{x}) \quad (12)$$

4.2 Formally Guarded Generative Model

Formally Guarded Generative Model (FGGM) setup allows the planner LLM to bind each generative model (GM) call to local input–output contracts (Φ_l, Ψ_l) that formally characterize the output. Once the FGGM is verified to be well-formed, the defined contracts can be used to prove the correctness of the agentic program against the behavioral specifications (Φ, Ψ) . Our FGGM design ensures the following essential properties. **1 Flexibility:** Unlike functions in \mathcal{F}_c , output contracts cannot be predefined by users because the expected output of a GM varies across calls depending on the prompt. In our design, the planner LLM synthesizes custom first-order formulas over terms involving non-parametric functions (e.g., *sqrt*) and predicates (e.g., *noDiff*) from \mathcal{F}_c to specify local contracts. **2 Parameter-Independent Correctness:** Provided the defined FGGM is well-formed, the specified contracts hold irrespective of the parameters of the underlying GM. This property is crucial to ensure that gradient-based parameter optimization never breaks the correctness guarantees. **3 Local Learning Objective:** Although correctness holds irrespective of parameter values, the local contracts (Φ_l, Ψ_l) extract the local learning objective for a particular GM call by characterizing its expected output. This objective can guide gradient-based optimization so that, after tuning, the GM’s conformance with (Φ_l, Ψ_l) improves (§ 4.2.3).

4.2.1 Intuition. For each FGGM with local contracts (Φ_l, Ψ_l) and underlying parametric GM \mathcal{L}_Θ , the planner LLM synthesizes two non-parametric programs using only functions from \mathcal{F}_c . The first is a fallback program f_d that satisfies (Φ_l, Ψ_l) and guards against edge cases in which the GM hallucinates or produces outputs that violate the specification. Without a verified fallback, the rejection sampler could exhaust all samples; the fallback is what makes the contract hold unconditionally for all parameter values. Intuitively, although f_d may lack the GM’s task-solving capability (low task performance), it prevents the candidate agent from entering unsafe program states during execution (high constraint satisfaction). The second is a prompting program f_p , which constructs the input to the GM (which may be just a neural network). For LLMs, f_p encodes task-specific natural-language instructions, analogous to prompt-tuning techniques [6, 21] known to improve GM performance and instruction adherence. Once the planner LLM proposes Φ_l, Ψ_l, f_p , and f_d , SEVerA automatically constructs a rejection sampler that treats outputs of underlying GM as samples from a proposal distribution π_p for an input \mathbf{p} . The local specifications (Φ_l, Ψ_l) define the support set (Definition 2.2) of the target distribution π_t on \mathbf{p} , assigning zero probability to outputs that violate (Φ_l, Ψ_l) . The rejection sampler draws a fixed number $K \geq 1$ of samples from \mathcal{L}_Θ , as specified by the user, and falls back to f_d if all samples are rejected (Fig 2). The purpose of f_p and optional conformance tuning is to improve the acceptance rate of the rejection sampler

for any input, thereby reducing reliance on the fallback f_d . We formally define the FGGM syntax, well-formedness conditions, and the computability of the output checker in § 5.2.1.

4.2.2 Examples. We show a couple of example FGGM definitions that the planner LLM suggests: one from symbolic regression and one from Dafny program verification. The first example from symbolic regression (Eq. 13) with id *boundedParam*, type signature $(l : \text{real}, u : \text{real}) \rightarrow \text{real}$, underlying parametric GM which is a neural network $NN_{\Theta}^{(2)} : \mathbb{R}^2 \rightarrow \mathbb{R}$ for this case, $\Phi_l : (l \leq u) \Psi_l$ (Eq. 13) restricts *boundedParam*(l, u) to lie within $[l, u]$. The fallback clamps any violating sample y using library functions *min* and *max*, ensuring contract preservation. The second example from Dafny program verification (Eq. 14) with id *dafnyAnnotate*, type signature $(p : \text{str}) \rightarrow \text{str}$, underlying parametric GM which is a LLM $\mathcal{L}_{\Theta} : \Sigma^* \rightarrow \Sigma^*$, $\Phi_l : \text{parse}(p)$, Ψ_l (Eq. 14) ensures *dafnyAnnotate*(p) always outputs a syntactically valid program and only adds annotations to the input program p . If all samples are rejected, the fallback safely returns the original program p , guaranteeing Ψ_l under Φ_l utilizing the **reflexivity axiom** (Eq 12). These examples show how a planner LLM can define FGGMs with fallbacks to avoid incorrect program states for diverse tasks.

```

function  $id(x_1 : T_1, \dots, x_n : T_n) : (T_o)$ 
  requires  $\Phi_l(x_1, \dots, x_n) \{$ 
  var  $p : T_{in} := f_p(x_1, \dots, x_n), y : T_o;$ 
  for  $i = 1 \dots K \{$ 
     $y := \mathcal{L}_{\Theta}(p);$ 
    if  $check_{\mathcal{A}, \Phi_l, \Psi_l}(x_1, \dots, x_n, y)$  1
      return  $y;$ 
    }
  return  $f_d(x_1, \dots, x_n, y);$  2
  }

```

Fig. 2. Auto-synthesized guarded GM with rejection sampler (1) and fallback (2).

$$\Psi_l : (l \leq f(l, u)) \wedge (f(l, u) \leq u), \quad f_d : \text{function } f_d(l:\text{real}, u:\text{real}, y:\text{real}) \{\text{return } \min(\max(l, y), u)\} \quad (13)$$

$$\Psi_l : \text{parse}(p) \wedge \text{noDiff}(p, f(p)), \quad f_d : \text{function } f_d(p:\text{str}, y:\text{str}) : (\text{str}) \{\text{return } p\} \quad (14)$$

Design Considerations. Alternative approaches for enforcing formal constraints on generative model calls exist, but they fall short of FGGM. Constrained decoding [4, 33, 35, 45] modifies the model’s internal decoding procedure, which limits its applicability to open-source models and restricts constraints to syntactic forms such as context-free grammars. Post-hoc output filtering without a fallback provides no guarantees: if all sampled outputs are rejected, the program produces no valid output. Fully verifying an agentic program is computationally intractable, as it depends on large parametric components such as LLMs. FGGM addresses all three limitations: it operates solely on model outputs, making it compatible with closed-source models; it incorporates a verified fallback that guarantees a valid output on every execution; and it decomposes verification into parameter-independent local contracts. Moreover, these local contracts can be utilized during learning to guide parameter tuning, improving conformance, and reducing reliance on the fallback.

4.2.3 Conformance Tuning. A FGGM that always uses the fallback is formally correct but practically no better than a non-parametric program f_d , conformance tuning helps avoid this by teaching the GM to satisfy the defined contract directly. The local contracts (Φ_l, Ψ_l) define the local learning objective to avoid the fallback and assign higher probability to outputs that pass the check $check_{\mathcal{A}, \Phi_l, \Psi_l}$. The objective in Eq. 20 minimizes the expected contract violation of GM samples over an input set \mathbb{P} . For any $\theta \in \Theta$, the inner term $L_{\Phi_l, \Psi_l}^{\theta}(p)$ measures the probability that a sampled output from $\mathcal{L}_{\theta}(p)$ violates $check_{\mathcal{A}, \Phi_l, \Psi_l}(\Phi_l, \Psi_l, x_1, \dots, x_n, y)$, while $L_{\Phi_l, \Psi_l}(\theta)$ aggregates this violation across inputs in \mathbb{P} . By optimizing $\theta \in \Theta$ to reduce this loss, the GM post-tuning is encouraged to concentrate the probability mass on the support defined by (Φ_l, Ψ_l) (details in § 5.2.5). Overall conformance tuning improves the acceptance rate of the rejection sampler, reducing reliance on f_d . We formalize the conformance loss in § 5.2.5 (Eq. 20).

4.3 SEVerA

SEVerA employs a three-step strategy to solve the constrained learning problem in Eq. 2: **(1) Search:** The planner LLM defines FGGMs with their prompting programs f_p , then samples candidate parametric programs that invoke these FGGMs. Direct calls to parametric GMs are disallowed; each must be wrapped within a verified FGGM.

(2) Verify: SEVerA checks all FGGM definitions for well-formedness and then verifies the candidate program against (Φ, Ψ) using the Dafny built-in verifier. If verification fails, error feedback is returned to the planner LLM to refine the candidate (CEGIS-style loop). Due to the FGGM setup, any verified program satisfies the behavioral specifications for all parameter values. **(3) Learn:** The learning step applies unconstrained gradient-based optimization to tune GM parameters, guided by the global task loss L and local conformance losses from each FGGM. For closed-source LLMs without accessible parameters, the learning step is skipped and improvements rely on prompt tuning through f_p . We formalize the search-verify-learn loop in § 5.2 (Algorithms 1 and 2).

4.3.1 Candidate Pool and Execution Feedback. SEVerA maintains a pool of verified candidate agentic programs with fine-tuned parameters (Fig. 1). SEVerA selects the best-performing candidate agent from the pool, i.e., the one with the minimal task loss L , and uses its execution traces on data D to iteratively search for new candidates via the planner LLM. The number of parametric candidate sampling steps is capped by the user. Once this limit is reached, SEVerA returns the agent in the pool with the lowest task loss L . Fig. 3 shows two candidate parametric programs sampled by the planner LLM. Each program invokes the FGGM defined in Eq. 13 with id `boundedParam` from two different locations. Fig. 3a shows the first parametric program, which fails to verify against the behavioral specification and is therefore pruned by the verifier. This candidate represents the output as a parametric affine function of the input x , which does not match the ground truth, illustrating how strict behavioral constraints eliminate poor candidates. The second candidate verifies for all parameters of the neural networks $NN_{\Theta}^{(2)}$ wrapped inside the `boundedParam` calls. After parameter tuning, it recovers the ground truth with negligible error, as discussed below. Currently, SEVerA does not share parameters across FGGM calls at different locations and instead assigns separate parameter sets to each `boundedParam` call. This design allows the system to potentially learn different parameters at different call sites to improve task-specific performance. Since both `boundedParam` calls receive constant inputs, after fine-tuning they return constant values, i.e., $a = 1.11$ and $d = 0.503$. This recovers the ground-truth function $\sqrt{1.23 \times \max(x, 0)}$ with negligible error. Note that SEVerA currently neither shares parameters nor caps the number of FGGM calls.

```
function agent(x : real) : (real) {
  var a : real := boundedParam(0.0, 1.0);
  var b : real := boundedParam(0.0, 1.0);
  var linear_x : real := a * x;
  var y : real := linear_x + b;
  return y;
}
```

(a) Unverified parametric program pruned out.

```
function agent(x : real) : (real) {
  if (x ≤ 0.0) {return 0.0;}
  var a : real := boundedParam(1.0, 1.5);
  var d : real := boundedParam(0.5, 0.8);
  assert x ≥ 0.0;
  var pow_x : real := pow(x, d);
  assert (d ≤ 0.8);  assert (d ≥ 0.5);
  assert (x ≤ 1) ⇒ (pow_x ≥ pow(x, 0.8));
  assert (x ≥ 1) ⇒ (pow_x ≥ pow(x, 0.5));
  assert (x ≥ 0) ⇒ (pow(x, 0.5) = sqrt(x));
  var y : real := a * pow_x;
  return y;
}
```

(b) Verified parametric program.

Fig. 3. Two example candidate parametric programs generated for symbolic regression.

However, SEVerA can be extended to support both these checks to reduce the computational budget required for parameter learning, as discussed in § 6.5. The parametric candidate programs and all FGGM definitions with synthesized prompting programs f_p for the program verification task are included in Appendix K.

In the next section, we formalize the search space, provide complete algorithmic details for each step, and establish theoretical guarantees for the framework.

5 TECHNICAL DETAILS

We formalize the search space in § 5.1, provide the details of the steps search, verify, and learn along with FGGM in § 5.2. We provide formal proof of soundness and a sufficient condition for agent synthesis success with non-trivial utility in § 5.3.

5.1 Search Space

Program Search Space: We use a restricted subset of the popular verification-aware language Dafny [26] to define the search space of candidate programs. In this restricted subset, we allow conditional blocks, while loops, assignment statements, assertion statements, function calls, and annotations including invariant clauses and ranking functions (decreases clauses), along with first-order input–output specifications. We present the BNF grammar G_D capturing this subset in Appendix A. We allow four basic types $T = \{\text{bool}, \text{int}, \text{real}, \text{string}\}$, basic built-in arithmetic functions $\mathcal{B}_a = \{+, -, \times, /\}$, boolean functions $\mathcal{B}_b = \{\&\&, ||, !, \implies\}$, and relations $\mathcal{B}_r = \{==, \leq, \geq, <, >\}$ with their usual semantics. We use a planner LLM $\mathcal{L}_p : \Sigma^* \rightarrow \Sigma^*$ to search the program space by sampling candidate programs as strings. We assume that \mathcal{L}_p 's output space can represent all valid Dafny programs $L(G_D) \subseteq \Sigma^*$ and that the type *string* can represent all strings Σ^* over Σ . SEVerA provides built-in predicates $\text{lex}_T(x : \text{string}) \rightarrow \text{bool}$ to check whether a string can be parsed as a specific type $T \in \{\text{bool}, \text{int}, \text{real}\}$. Using Dafny as the language allows us to leverage its built-in verifier to check candidate programs against a first-order input–output specification, as well as its built-in termination checker. This choice is motivated by several factors: (1) Dafny's built-in verification infrastructure directly supports the formal guarantees central to SEVerA, (2) its type system aligns with the basic types T used in our specification language, and (3) it supports annotations such as loop invariants and ranking functions that our planner LLM can synthesize to aid the verifier. While the framework is not language-specific, using a verification-aware host language avoids re-implementing a custom verifier and termination checker.

Library Functions: We categorize the set of library functions \mathcal{F} into two disjoint sets: \mathcal{F}_c , the set of non-parametric functions, and \mathcal{F}_p , the set of parametric functions. \mathcal{F}_c includes all built-in functions. Each user specified non-parametric function of arity n is specified by a tuple $(\text{id}(f_n), \tau_{f_n}, f_n^\Phi, f_n^\Psi)$, where $\text{id}(f_n)$ is the function name, $\tau_{f_n} = (T_1 \times \dots \times T_n \rightarrow T_0)$ is the type signature with $T_i \in T$, and (f_n^Φ, f_n^Ψ) are first-order input–output specifications defining a formal over-approximate semantics of f_n . The function $f_n : T_1 \times \dots \times T_n \rightarrow T_0$ halts for any inputs $x_1 \in T_1, \dots, x_n \in T_n$ and computes $r = f_n(x_1, \dots, x_n)$. The computed output r , together with (x_1, \dots, x_n) , always satisfies $\forall x_1 \in T_1, \dots, \forall x_n \in T_n. f_n^\Phi(x_1, \dots, x_n) \implies f_n^\Psi(x_1, \dots, x_n, r)$. \mathcal{A} contains all these formal contracts.

Parametric Functions: Each f_n^Θ in \mathcal{F}_p is specified by the tuple $(\text{id}(f_n^\Theta), \tau_{f_n^\Theta}, \theta_0, \Theta)$, where $\text{id}(f_n^\Theta)$ denotes the function name, $\tau_{f_n^\Theta}$ is the type signature, Θ is a continuous set of parameters, and $\theta_0 \in \Theta$ denotes the default parameters. Here, $f_n^\Theta = \{f_n^\theta \mid \theta \in \Theta\}$ represents an infinite set of functions with the same type signature $\tau_{f_n^\Theta}$, where each individual function $f_n^\theta \in f_n^\Theta$ is instantiated by a concrete parameter value $\theta \in \Theta$. Although f_n^Θ represents a set of functions, we sometimes abuse this notation to denote a single function corresponding to $f_n^{\theta_0}$ with the default parameters.

The input $\mathcal{I} = (\mathcal{F}, \mathcal{A}, L, D, \Phi, \Psi, I)$ to SEVerA includes the library functions \mathcal{F} , axioms \mathcal{A} , point-wise differentiable loss L , dataset D , behavioral specifications (Φ, Ψ) , and textual task information I used by the planner LLM \mathcal{L}_p . We use $\mathcal{P}_C : \Sigma^* \rightarrow \{T, F\}$ to denote the parser for checking the syntactic validity of input programs, $\mathcal{V}_C[\Phi, \Psi] : \Sigma^* \rightarrow \{T, F\}$ to denote the verifier that checks a program against the specifications (Φ, Ψ) , and $\mathcal{T}_C : \Sigma^* \rightarrow \{T, F\}$ to denote the termination checker. Each of \mathcal{P}_C , $\mathcal{V}_C[\Phi, \Psi]$, and \mathcal{T}_C is instantiated with user-provided context $C = (G_D, \mathcal{F}, \mathcal{A})$. Both $\mathcal{V}_C[\Phi, \Psi]$ and \mathcal{T}_C are restricted to a fixed timeout and return *false* if they cannot verify (Φ, Ψ) or prove termination within the allotted time budget.

5.2 Key Steps

5.2.1 *FGGM Formal Definition.* Eq. 15 specifies the formal syntax for defining an FGGM.

$$\langle FGGM \rangle ::= \langle id \rangle \text{ ";" } \langle GMid \rangle \text{ ";" } \overbrace{\langle typeSig \rangle \text{ ";" } \langle \Phi_I \rangle \text{ ";" } \langle \Psi_I \rangle \text{ ";" } \langle f_p \rangle \text{ ";" } \langle f_d \rangle \text{ ";" } \langle info \rangle}^{\text{signature + formal contract}} \quad (15)$$

prompt fallback

The planner LLM defines a guarded GM following Eq. 15, where: (a) *id*: is the unique identifier of the created FGGM; (b) *GMid*: identifies the underlying GM used as the proposal distribution; (c) *typeSig*: defines the type signature of the FGGM, including input variables with their types and the output type, e.g., $(x_1 : T_1, \dots, x_n : T_n) \rightarrow T_o$; (d) (Φ_I, Ψ_I) : are first-order formulas whose terms can include non-parametric library functions from \mathcal{F}_c ; (e) f_p : is the program that constructs the GM input from the typed input variables and the task description; (f) f_d : is a non-parametric fallback program that satisfies the local specifications (Φ_I, Ψ_I) ; and (g) *info*: is a textual description of the defined FGGM, utilized by the planner LLM during the generation of the parametric program. SEVerA parses these inputs, validates the definition, and automatically synthesizes a rejection sampler with fallback, as shown in Fig. 2 where \mathcal{L}_Θ is the parametric GM specified by *GMid*. The defined FGGM is then invoked in the candidate program via its identifier *id*. For valid definitions, SEVerA synthesizes $check_{\mathcal{A}, \Phi_I, \Psi_I}$, which checks whether a concrete input-output pair (x_1, \dots, x_n, y) satisfies the proposed contract (Φ_I, Ψ_I) , as elaborated in § 5.2.3.

We provide the BNF grammar for parsing the FGGM definition in Appendix L. Next, we discuss correctness checks for a proposed FGGM definition. Provided the SEVerA validates definition including the fallback f_d , $\forall \theta \in \Theta. \forall x_1 \in T_1 \dots \forall x_n \in T_n. \Phi_I(x_1, \dots, x_n) \implies \Psi_I(x_1, \dots, x_n, id(x_1, \dots, x_n))$ holds ensuring formal contracts are preserved for all parametric values $\theta \in \Theta$.

5.2.2 *FGGM Well-formedness.* Given an FGGM definition $G = (id_G, f_G^\Theta, \tau_G, \Phi_I, \Psi_I, f_p, f_d, info)$ following Eq. 15, this step checks whether the definition is well-formed. Let the parametric GM identified by *GMid* be $\mathcal{L}_\Theta : T_{in} \rightarrow T_o$. SEVerA first checks the well-formedness of the type signature $\tau_G = (T_1 \times \dots \times T_n) \rightarrow T_o$ and the well-formedness of the local contracts (Φ_I, Ψ_I) . We then verify that all terms in Φ_I and subsequently in Ψ_I that involve non-parametric functions from \mathcal{F}_c are valid (details in Appendix M). The G definition is considered valid provided both f_p and f_d are type-checked and terminating, and the fallback f_d satisfies the contract $\mathcal{V}_C[\Phi_I, \Psi_I](f_d)$. After the validity check, SEVerA constructs $check_{\mathcal{A}, \Phi_I, \Psi_I} : (T_1 \times \dots \times T_n \times T_o) \rightarrow \{T, F\}$ that checks whether a concrete input-output tuple (x_1, \dots, x_n, y) satisfies the local output contract Ψ_I (see § 5.2.3). We provide the pseudocode of $check_{\mathcal{A}, \Phi_I, \Psi_I}$ along with a soundness proof showing $\varphi_{\mathcal{A}} \implies (\forall x_1 \in T_1, \dots, \forall x_n \in T_n, \forall y \in T_o. \Phi_I(x_1, \dots, x_n) \implies (check_{\mathcal{A}, \Phi_I, \Psi_I}(x_1, \dots, x_n, y) \implies \Psi_I(x_1, \dots, x_n, y)))$ in Appendix R. Furthermore, for quantifier-free Ψ_I , the checker is complete (Lemma 5.1).

LEMMA 5.1 (COMPLETENESS OF CHECKER). *For any valid FGGM G with quantifier-free Ψ_I , $\forall x_1 \in T_1, \dots, \forall x_n \in T_n, \forall y \in T_o. \Phi_I(x_1, \dots, x_n) \implies (check_{\mathcal{A}, \Phi_I, \Psi_I}(x_1, \dots, x_n, y) \iff \Psi_I(x_1, \dots, x_n, y))$.*

PROOF SKETCH. We prove this by structural induction on the quantifier-free formula Ψ_l . *Base case:* If Ψ_l is an atomic predicate over terms built from computable library functions in \mathcal{F}_c , then under any concrete substitution (x_1, \dots, x_n, y) , $\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}$ evaluates the same computable terms and returns the same boolean value as Ψ_l . *Inductive step:* Assume the claim holds for formulas ϕ and ψ . For compound formulas $\neg\phi$, $\phi \wedge \psi$, and $\phi \vee \psi$, $\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}$ recursively evaluates the subformulas and applies the corresponding boolean operator, and returns the result. Thus, by induction on the structure of Ψ_l , $\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}(x_1, \dots, x_n, y)$ returns true iff $\Psi_l(x_1, \dots, x_n, y)$ holds whenever $\Phi_l(x_1, \dots, x_n)$ holds. Hence, the checker is complete. Formal details in Appendix R. \square

For any valid FGGM G , the rejection sampler $\text{id}_G : T_1 \times \dots \times T_n \rightarrow T_o$ with fallback synthesized by SEVerA (Fig. 2) satisfies the defined contract (Φ_l, Ψ_l) for all well-typed inputs (x_1, \dots, x_n) . Crucially the contract preservation holds over all parameters $\theta \in \Theta$ of the underlying parametric model f_G^Θ . The function id_G terminates on all inputs since the number of samples K is finite and f_p , f_d , and $\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}$ are all terminating. We provide a soundness proof sketch in Theorem 5.2.

THEOREM 5.2 (VALID LOCAL CONTRACT). *For any valid FGGM G with (Φ_l, Ψ_l) and parametric GM f_G^Θ , $\varphi_{\mathcal{A}} \implies (\forall \theta \in \Theta, \forall x_1 \in T_1, \dots, \forall x_n \in T_n. \Phi_l(x_1, \dots, x_n) \implies \Psi_l(x_1, \dots, x_n, \text{id}_G(x_1, \dots, x_n)))$.*

PROOF SKETCH. Let $r = \text{id}_G(x_1, \dots, x_n)$. Based on the rejection sampler with the checker $\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}$, the following condition always holds for the output r (Eq. 16). Eq. 17 follows from the validity of the fallback f_d , and Eq. 18 follows from the soundness of the checker $\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}$. To simplify the notation, we used $(\forall x_i \in T_i)$ to denote $(\forall x_1 \in T_1, \dots, \forall x_n \in T_n)$ and y_f is the final rejected sample if f_d is used.

$$\varphi_{\mathcal{A}} \implies (\forall \theta \in \Theta. (\forall x_i \in T_i). \Phi_l(x_1, \dots, x_n) \implies \text{check}_{\mathcal{A}, \Phi_l, \Psi_l}(x_1, \dots, x_n, r) \vee (r = f_d(x_1, \dots, x_n, y_f))) \quad (16)$$

$$\varphi_{\mathcal{A}} \implies \forall \theta \in \Theta. (\forall x_i \in T_i). \Phi_l(x_1, \dots, x_n) \wedge (r = f_d(x_1, \dots, x_n, y_f)) \implies \Psi_l(x_1, \dots, x_n, r) \quad (17)$$

$$\varphi_{\mathcal{A}} \implies (\forall \theta \in \Theta. (\forall x_i \in T_i). (\Phi_l(x_1, \dots, x_n) \wedge \text{check}_{\mathcal{A}, \Phi_l, \Psi_l}(x_1, \dots, x_n, r)) \implies \Psi_l(x_1, \dots, x_n, r)) \quad (18)$$

$$\varphi_{\mathcal{A}} \implies (\forall \theta \in \Theta. (\forall x_i \in T_i). \Phi_l(x_1, \dots, x_n) \implies \Psi_l(x_1, \dots, x_n, r)) \quad \text{Using Eq. (16, 17, 18)}$$

We provide the detailed formal proof in Appendix R. \square

5.2.3 Computability of $\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}$ (Fig. 2) on concrete values. The auto-synthesized rejection sampler (Fig. 2) requires checking whether a sample y from \mathcal{L}_Θ , given a concrete input (x_1, \dots, x_n) , satisfies the planner LLM-defined Ψ_l . We show that for any *quantifier-free* Ψ_l , this check can be computed efficiently for all concrete values (x_1, \dots, x_n, y) . Since all library functions in \mathcal{F}_c are computable, $\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}$ can evaluate every term on the concrete inputs (x_1, \dots, x_n, y) . Satisfiability can then be decided in linear time in the number of terms in Ψ_l . We provide formal completeness for the quantifier-free fragment in Lemma 5.1. For Ψ_l with quantifiers, we instantiate $\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}$ with the axioms \mathcal{A} . Then, $\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}$ substitutes the free variables in Ψ_l with their corresponding concrete values, encodes the axioms \mathcal{A} and Φ_l similarly to Eq. 8, and submits the resulting formula to an SMT solver. First-order formulas over the basic types \mathbb{T} and uninterpreted library functions may lie outside decidable fragments. Therefore, we bound the SMT solver with a user-specified timeout. $\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}$ returns true only if the solver succeeds within the timeout; otherwise, it returns false. This guarantees soundness; however, in the presence of quantifiers, $\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}$ may be incomplete and may incorrectly reject valid samples y .

5.2.4 Search and Verify. The search and verification steps jointly implement a CEGIS-style loop that explores candidate parametric programs while enforcing the behavioral specification. In the **search** phase, the planner LLM \mathcal{L}_p first synthesizes a set of FGGM definitions $\mathcal{G} = \{G_1, \dots, G_m\}$ following Eq. 15. Each FGGM wraps a parametric model with local contracts and a verified fallback. Using these FGGMs as callable functions, the planner then samples a candidate parametric program

$P_{\{\Theta\}}$ where $\{\Theta\} = \{\Theta_1, \dots, \Theta_k\}$ represents the set of all optimizable parameters in $P_{\{\Theta\}}$. As mentioned in § 4.3.1, FGGM calls at different program locations do not share parameters and add their own parameter in $\{\Theta\}$. In the verification phase, SEVerA first checks the well-formedness of every $G_i \in \mathcal{G}$. If any of the definitions fail to verify it returns an error message. If all definitions are valid, the verification context $C = (G_D, \mathcal{F}, \mathcal{A})$ is extended with the synthesized FGGMs and their local contracts. The updated context $C' = (G_D, \mathcal{F} \cup \mathcal{G}, \mathcal{A} \cup \mathcal{A}_{\mathcal{G}})$ adds all FGGMs \mathcal{G} along with $\mathcal{A}_{\mathcal{G}}$ containing local contracts $(G_i^{\Phi_l}, G_i^{\Psi_l})$ of each FGGM $G_i \in \mathcal{G}$. The candidate program $P_{\{\Theta\}}$ is then checked for syntactic validity, termination, and compliance with the behavioral specification (Φ, Ψ) . If verification succeeds, SEVerA accepts the parametric program $P_{\{\Theta\}}$ and ensures that $P_{\{\Theta\}}$ satisfies (Φ, Ψ) over all parameters (Theorem 5.3). Algorithm 1 summarizes the combined search–verification loop.

Given the input specification I and planner \mathcal{L}_p , the algorithm iteratively proposes candidate solutions within a fixed budget δ (Lines 1–4). In each iteration, the planner first synthesizes a set of FGGM definitions \mathcal{G} (Line 6), which are then checked for well-formedness using *validateFGGM* (Line 7). Within \mathcal{G} , the planner LLM \mathcal{L}_p also synthesizes the prompting functions f_p , enabling each FGGM-specific customization. If the definitions are valid, the planner samples a candidate parametric program $P_{\{\Theta\}}$ that may invoke these FGGMs (Line 11), and the verification context is extended with their contracts (Line 12). The candidate program is then validated for syntactic correctness, termination, and compliance with the behavioral specification using *validateProgram* (Line 13). If verification succeeds, the program is returned (Lines 14–15); otherwise, the loop continues with the error feedback until the search budget is exhausted (Lines 16–18).

Algorithm 1 searchVerify

```

1: Input:  $I = (\mathcal{F}, \mathcal{A}, \Phi, \Psi, I)$ , planner  $\mathcal{L}_p$ , budget  $\delta$ 
2: Input: Previous attempts feed-back  $I'$ 
3: Output: Verified Parametric Program  $P_{\{\Theta\}}$  else  $\perp$ 
4:  $i \leftarrow 0$ ,  $\text{err} \leftarrow \{\}$ 
5: while  $i < \delta$  do
6:    $\mathcal{G} \leftarrow \text{defineFGGM}(\mathcal{L}_p, I, \mathcal{F}, \mathcal{A}, I', \text{err})$ 
7:    $\text{err} \leftarrow \text{validateFGGM}(\mathcal{G}, \mathcal{F}_c, \mathcal{A})$ 
8:   if  $\text{err} \neq \{\}$  then
9:      $i \leftarrow i + 1$ ; continue
10:  end if
11:   $P_{\{\Theta\}} \leftarrow \text{sampleProgram}(\mathcal{L}_p, I, \mathcal{F}, \mathcal{A}, I', \mathcal{G}, \text{err})$ 
12:   $C \leftarrow (G_D, \mathcal{F} \cup \mathcal{G}, \mathcal{A} \cup \mathcal{A}_{\mathcal{G}})$ 
13:   $\text{err} \leftarrow \text{validateProgram}(P_{\{\Theta\}}, \mathcal{P}_C, \mathcal{T}_C, \mathcal{V}_C[\Phi, \Psi])$ 
14:  if  $\text{err} = \{\}$  then
15:    return  $P_{\{\Theta\}}$ 
16:  end if
17:   $i \leftarrow i + 1$ 
18: end while
19: return  $\perp$ 

```

THEOREM 5.3. *If $P_{\{\Theta\}} \neq \perp$ with FGGM set \mathcal{G} , then $\forall \theta_1 \in \Theta_1, \dots, \forall \theta_k \in \Theta_k. \varphi_{\mathcal{A}} \implies (\forall x_1 \in T_1, \dots, \forall x_n \in T_n. \Phi(x_1, \dots, x_n) \implies \Psi(x_1, \dots, x_n, P_{\{(\theta_1, \dots, \theta_k)/\Theta\}}(x_1, \dots, x_n)))$.*

PROOF SKETCH. Follows from Theorem 5.2 with details in Appendix R. \square

5.2.5 Learn. Once a candidate parametric program $P_{\{\Theta\}}$ passes the search–verification loop (Algorithm 1), the verification step guarantees that $P_{\{\Theta\}}$ satisfies the behavioral specification (Φ, Ψ) for all parameter values (Theorem 5.3). Hence, the learning step can freely select parameters without violating (Φ, Ψ) . Finding the optimal solution to the general formulation in Eq. 19 over $\{\Theta\}$ for a general loss L is practically intractable. Instead, SEVerA employs a scalable gradient-descent–based optimization to efficiently search over $\{\Theta\}$.

$$(\theta_1^*, \dots, \theta_k^*) = \arg \min_{(\theta_1, \dots, \theta_k) \in \{\Theta\}} \sum_{(x, y) \in D} L(x, y, P_{\{(\theta_1, \dots, \theta_k)/\Theta\}}(x)) \quad (19)$$

Conformance loss L_{Φ_l, Ψ_l} for each local contract (Φ_l, Ψ_l) encourages the parameters to produce outputs that satisfy their local contracts. By improving conformance with these contracts, SEVerA increases the success rate of the corresponding rejection sampler and thereby avoids triggering the static non-parametric fallback program. Since the fallback cannot be tuned, reducing its usage allows the gradient-based search to more effectively guide the parameters toward solutions that

reduce the final loss L . Our experiments in § 6.3 substantiate this observation, showing that reducing reliance on the fallback improves task performance. For a single FGGM with local contracts (Φ_l, Ψ_l) and parameter θ , we define the conformance loss over an input set \mathbb{P} in Eq. 20. The inner term $L_{\Phi_l, \Psi_l}^{\theta}(p)$ measures the probability that a sampled output from $\mathcal{L}_{\theta}(p)$ violates $\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}$, while $L_{\Phi_l, \Psi_l}(\theta)$ aggregates this violation across inputs.

$$\underbrace{L_{\Phi_l, \Psi_l}(\theta) = \frac{1}{|\mathbb{P}|} \times \sum_{p \in \mathbb{P}} L_{\Phi_l, \Psi_l}^{\theta}(p)}_{\text{Aggregated constraint violation over } \mathbb{P}}, \quad \text{where } L_{\Phi_l, \Psi_l}^{\theta}(p) = \underbrace{\mathbb{E}_{y \sim \mathcal{L}_{\theta}(p)} 1 - \mathbb{I}(\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}(x_1, \dots, x_n, y))}_{\text{Expected constraint violation on single GM input } p} \quad (20)$$

The augmented objective, which includes a local conformance loss for each FGGM call, is defined in Eq. 21, where $\lambda \in \mathbb{R}^+$ is a user-provided constant.

$$\sum_{(x, y) \in D} L(x, y, P_{\{(\theta_1, \dots, \theta_k)/\Theta\}}(x)) + \lambda \times \sum_{i \in k, (x, y) \in D} L_{\Phi_l, \Psi_l}^i(\theta_i) \quad (21)$$

Here, the first term is the task-specific loss over the dataset D , and the second term aggregates the local conformance losses (Eq. 20) of the FGGMs appearing in the program. Optimizing all parameters jointly according to Eq. 21, while feasible for smaller neural networks, becomes practically expensive for agents with multiple interdependent LLM calls. To address this challenge, we leverage a key structural property: the conformance loss $L_{\Phi_l, \Psi_l}^i(\theta_i)$ for each FGGM depends only on the local inputs to that FGGM call and the outputs of its underlying GM, not on the parameters of other FGGMs. This independence arises because each FGGM’s rejection sampler and $\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}$ checker operate solely on the local contract (Φ_l^i, Ψ_l^i) and the corresponding GM’s outputs. As a result, the conformance term decomposes naturally across FGGM calls. Following this observation, for larger models we approximate the augmented loss by optimizing each FGGM id_{θ_i} with parameter θ_i independently, as shown in Eq. 22. Here, the first term includes the task-specific loss L only when the output y_l of id_{θ_i} matches the final output y of the agent on input x , where $(x, _) \in D$. The modularized loss (Eq. 22) enables efficient parallelized training across different θ_i .

$$L_a(\theta_i) = \frac{1}{|\mathbb{P}|} \times \sum_{p \in \mathbb{P}} L_a^{\theta_i}(p); \quad L_a^{\theta_i}(p) = \mathbb{E}_{y_l \sim \text{id}_{\theta_i}(p)} L(x, _ y) \times \mathbb{I}(y = y_l) + \lambda \times L_{\Phi_l, \Psi_l}^{\theta_i}(p) \quad (22)$$

$$\mathcal{R}(p, y_l) = 1 - \text{Sigmoid}(L(x, _ y) \times \mathbb{I}(y = y_l) + \lambda \times (1 - \mathbb{I}(\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}(p, y_l)))) \quad (23)$$

For generative models with accessible parameters, the objective in Eq. 22 can be optimized using reinforcement learning–based fine-tuning methods such as GRPO [36] (details in Appendix N). Eq. 23 defines the reward for an input-output pair (p, y_l) , with the goal of learning an output distribution over y_l for inputs $p \in \mathbb{P}$ that maximizes total reward $\mathcal{R}(p, y_l)$. The reward function combines two signals: (1) the task-specific loss L , which drives the GM toward outputs that reduce overall task error, and (2) the conformance indicator $\mathbb{I}(\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}(p, y_l))$, which rewards outputs satisfying the local contract. The sigmoid transformation maps the combined penalty to the $[0, 1]$ range, providing a smooth reward landscape. Outputs that violate the local contract receive near-zero reward due to the $(1 - \mathbb{I}(\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}(p, y_l)))$ penalty term, strongly discouraging the GM from producing outputs that would be rejected by the rejection sampler. The reward function need not be differentiable w.r.t. θ_i , as long as the outputs y_l remain differentiable w.r.t. θ_i [36]. For closed-source models with inaccessible parameters, the learning step is skipped, and performance improvements rely solely on prompt tuning through the synthesized f_p programs within each FGGM.

5.2.6 SEVerA. Algorithm 2 summarizes the overall workflow of SEVerA. The algorithm maintains a pool \mathcal{P} of verified candidate agents and iteratively improves it through a search–verify–learn loop. In each iteration, SEVerA first invokes the `searchVerify` procedure (Algorithm 1), which performs CEGIS-style discrete search to synthesize FGGM definitions and candidate parametric program $P_{\{\Theta\}}$ that satisfy the behavioral specification (Φ, Ψ) . If verification succeeds, the program is passed to the `LEARN` procedure, which optimizes the parameters of the underlying generative models using the dataset D and task loss L . For closed-source models with inaccessible parameters, this step is

skipped, and performance improvements rely solely on prompt tuning through the synthesized f_p programs in each FGGM.

The resulting tuned agent is added to the candidate pool, and its execution traces on D are used to construct feedback I' (see Appendix P) that guides the planner in subsequent search iterations. After the search budget Δ is exhausted, SEVerA returns the agent in the pool that achieves the lowest task loss on the dataset. If \mathcal{P} is empty, SEVerA returns \perp .

5.3 Theoretical Results

The soundness of SEVerA guarantees that any agent program $f^* \neq \perp$ returned satisfies the behavioral specifications (Φ, Ψ) . The `searchVerify` step generates candidate programs $P_{\{\Theta\}}$ that satisfy (Φ, Ψ) for all $(\theta_1, \dots, \theta_k) \in \{\Theta\}$. Consequently, the parameters $(\theta_1^*, \dots, \theta_k^*)$ predicted by the learning step preserve constraint satisfaction. The pool \mathcal{P} is either empty or contains only programs that satisfy the behavioral specifications.

THEOREM 5.4 (SOUNDNESS). *If $f^* \neq \perp$ then, $\varphi_{\mathcal{A}} \implies \forall x_1 \in T_1, \dots, \forall x_n \in T_n. \Phi(x_1, \dots, x_n) \implies \Psi(x_1, \dots, x_n, f^*(x_1, \dots, x_n))$.*

PROOF SKETCH. $(f^* \neq \perp) \implies (f^* \in \mathcal{P})$ and all program $f \in \mathcal{P}$ satisfies $\varphi_{\mathcal{A}} \implies \forall x_1 \in T_1, \dots, \forall x_n \in T_n. \Phi(x_1, \dots, x_n) \implies \Psi(x_1, \dots, x_n, f(x_1, \dots, x_n))$ following Theorem 5.3. We provide the formal details in Appendix R. \square

Next, we characterize a sufficient condition that ensures the existence of a candidate solution $f \in S(G, \mathcal{F})$ that satisfies (Φ, Ψ) and achieves no worse loss L than any generative model $f_n^\theta \in \mathcal{F}_p$ with any prompting function $f_p \in S(G, \mathcal{F}_c)$ and initial parameters θ_0 . The high-level idea is as follows. If there exists a non-parametric program $f_d \in S(G, \mathcal{F}_c)$ that satisfies (Φ, Ψ) , we can use it as a fallback together with any type-correct generative model f_n^θ (i.e., with output type T_o) to construct an FGGM. A program that invokes this FGGM on the input and returns its output satisfies (Φ, Ψ) . Moreover, if `check $_{\mathcal{A}, \Phi, \Psi}$` is complete and does not incorrectly reject valid samples from f_n^θ , then the FGGM always returns valid samples from f_n^θ . Lemma 5.1 shows that `check $_{\mathcal{A}, \Phi, \Psi}$` is complete when Ψ is quantifier-free. Assume: (i) L penalizes constraint-violating outputs more than constraint-satisfying ones, i.e., $\forall x \in T_i, \forall y, y' \in T_o. (\Phi(x) \wedge \neg\Psi(x, y) \wedge \Psi(x, y')) \implies L(x, _ , y') < L(x, _ , y)$, (ii) Ψ is quantifier-free, and (iii) there exists a valid non-parametric program f_d satisfying (Φ, Ψ) . Then there exists a program $f \in S(G, \mathcal{F})$ that satisfies (Φ, Ψ) and incurs no greater loss than any type-correct generative model $f_n^\theta : T_i \rightarrow T_o$ with initial parameters. Moreover, if $f_n^\theta(x)$ violates $\Psi(x, f_n^\theta(x))$ for some $(x, _) \in D$, the improvement in loss is strict.

THEOREM 5.5 (SUFFICIENT SUCCESS CONDITION). *If (i) L penalizes constraint violations, i.e., $\forall x \in T_i, \forall y, y' \in T_o. (\Phi(x) \wedge \neg\Psi(x, y) \wedge \Psi(x, y')) \implies L(x, _ , y') < L(x, _ , y)$, (ii) Ψ is quantifier-free, (iii) there exists a non-parametric program $f_d \in S(G, \mathcal{F}_c)$ such that f_d satisfies (Φ, Ψ) . Then, for any type-correct generative model $f_n^\theta : T_i \rightarrow T_o$ with initial parameters and any prompting program f_p , there*

Algorithm 2 SEVerA

```

1: Input:  $I = (\mathcal{F}, \mathcal{A}, L, D, \Phi, \Psi, I)$ , planner  $\mathcal{L}_p$ 
2: Input: total budget  $\Delta$ , budget per candidate  $\delta$ 
3: Output: Best agent  $f^*$  or  $\perp$  if search fails
4:  $\mathcal{P} \leftarrow \{\}$  ▷ Pool of verified candidates
5:  $I' \leftarrow \emptyset$ 
6: while  $\delta \leq \Delta$  do
7:    $P_{\{\Theta\}} \leftarrow \text{searchVerify}(\mathcal{F}, \mathcal{A}, \Phi, \Psi, I, \mathcal{L}_p, \delta, I')$ 
8:    $\Delta \leftarrow \Delta - \delta$ 
9:   if  $P_{\{\Theta\}} = \perp$  then
10:     continue
11:   end if
12:   if  $\{\Theta\} \neq \{\}$  then
13:      $(\theta_1^*, \dots, \theta_k^*) \leftarrow \text{learn}(P_{\{\Theta\}}, L, D)$ 
14:      $\mathcal{P} \leftarrow \mathcal{P} \cup \{P_{\{(\theta_1^*, \dots, \theta_k^*)/\Theta\}}\}$ 
15:   else
16:      $\mathcal{P} \leftarrow \mathcal{P} \cup \{P_{\{\Theta\}}\}$ 
17:   end if
18:    $I' \leftarrow \text{collectFeedback}(P_{\{\Theta\}}, D)$ 
19: end while
20:  $f^* \leftarrow \arg \min_{f \in \mathcal{P}} \sum_{(x, y) \in D} L(x, y, f(x))$ 
21: return  $f^*$ 

```

exists a program $f \in S(G, \mathcal{F})$ such that: (1) f satisfies (Φ, Ψ) , and (2) $L(f) \leq L(f_n^\theta)$. Moreover, if there exists $(x, _) \in D$ such that $\neg\Psi(x, f_n^\theta(x))$, then $L(f) < L(f_n^\theta)$ where $L(f) = \sum_{(x, _) \in D} L(x, _ f(x))$.

PROOF SKETCH. Define $G = (id, f_n^\theta, \tau, \Phi_l, \Psi_l, f_p, f_d, info)$ where $(\Phi_l, \Psi_l) := (\Phi, \Psi)$. Since f_d satisfies (Φ, Ψ) , the FGGM is valid. Define the program:

function $f(x_1 : T_1, \dots, x_n : T_n) : T_o$ **{return** $id(x_1, \dots, x_n);$ **}**

For any input satisfying Φ , by completeness of $check_{\mathcal{A}, \Phi_l, \Psi_l}$ (Lemma 5.1), id returns $f_n^\theta(x)$ whenever $\Psi(x, f_n^\theta(x))$ holds, and otherwise returns $f_d(x)$, which satisfies Ψ . For each $(x, _) \in D$, $f(x)$ either equals $f_n^\theta(x)$ (if valid) or $f_d(x)$. By assumption (i), such replacement cannot increase loss, so $L(f) \leq L(f_n^\theta)$. If f_n^θ violates Ψ on some $(x, _) \in D$, the improvement is strict. Formal details in Appendix R. \square

Theorem 5.5 provides a constructive argument: the FGGM mechanism can always improve upon a bare GM call by filtering out constraint-violating outputs and substituting the fallback. The three assumptions are mild in practice. Assumption (i) requires only that constraint-satisfying outputs are preferred over violating ones, which holds for any reasonable loss function in constrained settings. Assumption (ii), requiring quantifier-free Ψ , ensures completeness of the checker $check_{\mathcal{A}, \Phi_l, \Psi_l}$ (Lemma 5.1) so that no valid GM outputs are incorrectly rejected. In our experiments, the symbolic regression task uses quantifier-free specifications, while the program verification task uses library-function predicates (e.g., *parse*, *noDiff*) that are computable on concrete inputs, effectively behaving as quantifier-free checks at runtime. Assumption (iii), the existence of a valid fallback, is a natural requirement: if no non-parametric program can satisfy the specification, then the specification may be too strong for the given library. The theorem does not guarantee that SEVerA will *find* the optimal program, as the planner LLM may not propose the right candidate within the search budget. Rather, it establishes that the search space $S(G, \mathcal{F})$ is rich enough to contain solutions that dominate unconstrained GM calls.

6 EVALUATION

We evaluate SEVerA on four tasks spanning scientific discovery, program verification, mathematical reasoning, and agentic tool use. The evaluation seeks to answer three main questions: whether formal constraints improve safety, whether learning remains effective under those constraints, and how conformance tuning with local constraints improves performance.

6.1 Experimental Setup

6.1.1 Tasks and Datasets. We instantiate SEVerA on four tasks, each representative of a different kind of constrained agent synthesis problem.

LLM-Assisted Program Verification (DafnyBench). The agent receives a Dafny program with pre-encoded input-output specifications and must synthesize annotations that enable verification against the pre-encoded specification inside the input program [29]. The behavioral specification of the agent requires the output to remain parsable and equivalent to the input modulo added annotations. The task performance is measured by verification success within a fixed time budget. The underlying LLM is Claude Sonnet 4.5 [2], a closed-source model whose weights are not accessible; consequently, the learning step (i.e., parameter tuning) is not applicable to this task.

Symbolic Math Synthesis (GSM-Symbolic). GSM-Symbolic is a benchmark of grade-school math word problems generated from symbolic templates [30]. We use it to evaluate the LLM’s ability to synthesize correct symbolic mathematical expressions. Behavioral specification of the agent enforces that each generated expression is syntactically valid with respect to a formal grammar, and task performance measures whether the answer is equivalent to the ground-truth expression.

Agentic Tool Use (τ^2 -bench). τ^2 -bench [5] evaluates conversational LLM agents in realistic customer-service scenarios. We use the retail and airlines domains. The agent must select and invoke API tools to resolve user requests while respecting domain-specific policies. Unlike the other tasks, the hard constraint and the utility metric operate at different granularities. The hard constraint is enforced at each individual tool call: before every tool call is executed, the Agent-C checker [20] verifies that the proposed call complies with temporal policy-compliance rules (e.g., refund eligibility, booking-modification policies) given the trace of prior calls. Task performance (pass rate) is measured over the entire multi-turn interaction trace as the percentage of interactions where the agent successfully resolves the user request.

Constrained Symbolic Regression (SymReg). In symbolic regression, each synthesis instance consists of noisy input–output observations generated by an unknown ground-truth function f_{gt} , along with a behavioral specification encoding prior knowledge about f_{gt} . The goal is to synthesize a neuro-symbolic program (i.e., the agent) that fits the observed input–output pairs while provably satisfying the behavioral specifications. We consider a total of 35 synthesis tasks, each with 2 different training sets corresponding to varying noise levels (5%, 10%), and behavioral specifications that encode symbolic bounds on the ground truth. Each training set includes 600 noisy samples and 400 test samples. All synthesis instances are drawn from popular benchmarks [7, 17].

6.1.2 Implementation Details. All experiments are conducted on NVIDIA A100 GPUs (40 GB). Across all tasks we use the same high-level pipeline: planner-guided program search with a budget of $\Delta = 10$ iterations, FGGM-based rejection sampling with a fixed sample budget $K = 5$, deductive verification with bounded verifier and SMT timeouts, and hyperparameter tuning when model parameters are accessible. The planner LLM is provided with 3 ground-truth examples as in-context examples (2 for τ^2 -bench, due to the long traces); these examples are excluded from all train and test splits. The full planner prompt template is given in Appendix Q, and details on how execution feedback is constructed and provided to the planner between iterations are given in Appendix P.

6.1.3 Learning Implementation. Once a program verifies, SEVerA uses Dafny’s built-in transpiler [9] to translate the verified code into Python. Parameter tuning and experiments are then conducted on the transpiled code. SEVerA assumes that Dafny’s transpiler is sound for this restricted subset of the language.

6.1.4 Metrics. For each task, we report: (a) the *constraint violation rate* on a held-out test set of unseen inputs, measuring the fraction of outputs that violate the behavioral specification (Φ, Ψ); (b) *task performance*, measured by mean squared error for SymReg, verification success rate for Dafny, answer accuracy for GSM-Symbolic, and pass rate for τ^2 -Bench; and (c) *wall-clock time*.

6.1.5 Baseline. For all four tasks outside GSM-symbolic, we consider state-of-the-art (SOTA) hand-crafted agents. For GSM-Symbolic, we use the SOTA constrained decoder [4]. Some hand-crafted agents (e.g., Agent-C on τ^2 -Bench) manually enforce hard constraints; however, being static, they achieve lower task performance. We also consider self-evolving frameworks without constraints, which do not verify the generated programs against behavioral specifications and accept all planner programs provided they are syntactically valid and type-checked. Consequently, such frameworks provide no guarantees about the correctness or safety of the generated code. To ensure fairness, unconstrained self-evolving frameworks are only allowed to use the same set of library functions.

6.1.6 Verification Success. For program verification, tool calling, and GSM-Symbolic, SEVerA consistently finds a valid program satisfying (Φ, Ψ) within 10 attempts. For symbolic regression, it produces a verified program in 33 out of 35 cases, whereas the baselines generate at least 11 cases where the generated program violates the behavioral specification on test inputs.

Table 2. Dafny program verification results on HumanEvalDafny and DafnyBench. *Verif. & NoDiff* counts a program as correct only if it both verifies and passes the AST-based diff check against the original; *Verif.* counts verification alone.

Dataset	Method	Ver. & NoDiff (%) \uparrow	Ver. (%) \uparrow	Vio. (%) \downarrow	Time (s)
HumanEvalDafny	LLM (Claude Sonnet 4.5)	73.7	76.8	8.1	9.8
	DafnyBench baseline	86.9	87.9	4.0	16.1
	SEVerA (w/o constraints)	84.8	88.9	5.1	15.7
	SEVerA	97.0 (+10.1)	97.0 (+9.1)	0.0 (-4.0)	18.2
DafnyBench	LLM (Claude Sonnet 4.5)	68.7	71.1	10.3	10.3
	DafnyBench baseline	81.6	84.0	8.2	20.1
	SEVerA (w/o constraints)	79.2	84.8	7.9	18.4
	SEVerA	89.1 (+7.5)	89.1 (+5.1)	0.0 (-8.2)	25.6

6.1.7 *Dataset splits.* Table 1 summarizes the dataset sizes for each synthesis task. For each task, a small training set is used for the learning step and the remaining examples form the held-out test set on which all reported metrics are computed. GSM-Symbolic uses a larger training fraction (50%) because its learning step performs GRPO-based fine-tuning, which requires more training data than the other tasks.

Dataset	Total	Train
HumanEvalDafny	135	30
DafnyBench	760	50
τ^2 -bench Retail	114	15
τ^2 -bench Airline	50	10
GSM-Symbolic	100	50
Symbolic Regression	600	400

Table 1. Dataset sizes and train/test splits.

6.1.8 *Task-specific details.* For Dafny, the underlying LLM is Claude Sonnet 4.5 [2], a closed-source model; because model weights are inaccessible, Dafny experiments use only the search and verification stages of SEVerA and do not perform parameter tuning. For τ^2 -bench, parameter tuning is also omitted because the long, multi-turn execution traces make gradient-based fine-tuning prohibitively expensive. For GSM-Symbolic, the learning step uses GRPO [36] with LoRA adapters [18] on Qwen3-8B [42]. Additional hyperparameters (LoRA rank, solver timeouts, etc.) are reported in Appendix O.

6.2 Main Results

We begin with the central question motivating SEVerA:

RQ1 (Safety). *What goes wrong without formal constraints, and does SEVerA eliminate these failures?*

We compare unconstrained agent synthesis against SEVerA across all four tasks to show that hard behavioral specifications are necessary: without them, agents silently produce invalid outputs on unseen inputs that are undetectable by soft performance metrics.

6.2.1 *LLM-Assisted Program Verification (Dafny).* Table 2 presents results on two Dafny program-verification benchmarks and offers the clearest illustration of why formal constraints are indispensable. The gap between the *Verif.* and *Verif. & NoDiff* columns reveals the core safety issue (**RQ1**): without hard constraints, agents inflate their verification rate by modifying the original program. The LLM baseline reports 76.8% verification on HumanEvalDafny, yet only 73.7% of those outputs actually preserve the input program. Overall, 8.1% of all outputs are violations that would go undetected by a metric that checks verification alone. The DafnyBench baseline and SEVerA without constraints reduce but do not eliminate this problem (4.0% and 5.1% violations, respectively). Only the full SEVerA pipeline, which enforces the *NoDiff* behavioral specification, achieves a 0% violation rate on both benchmarks while simultaneously attaining the highest *Verif.*

Table 3. τ^2 -bench results with Qwen3-8B on the retail and airline domains. Higher pass rate is better; lower violation rate is better.

Domain	Method	Pass Rate (%) \uparrow	Violation (%) \downarrow	Time (s)
Retail	LLM (Qwen3-8B)	11.3	76.3	146.6
	Agent-C (Qwen3-8B)	42.2	0.0	234.7
	SEVerA (w/o constraints)	49.4	10.3	238.3
	SEVerA	53.6 (+11.4)	0.0 (0.0)	212.4
Airline	LLM (Qwen3-8B)	13.2	68.4	184.8
	Agent-C (Qwen3-8B)	39.4	0.0	272.6
	SEVerA (w/o constraints)	44.7	25.5	268.4
	SEVerA	52.6 (+13.2)	0.0 (0.0)	241.1

& NoDiff rate (97.0% on HumanEvalDafny, 89.1% on DafnyBench). This confirms that FGGM-based rejection sampling with a verified fallback does not merely filter outputs but actively steers the planner toward higher-quality candidates. This result is particularly notable because it is achieved without any parameter tuning (the underlying model, Claude Sonnet 4.5, is closed-source), the gains come entirely from making constraints visible to the planner and enforcing them at synthesis time. The overhead introduced by SEVerA is modest. The full pipeline takes 18.2 s per instance on HumanEvalDafny compared to 9.8 s for LLM generation, roughly a $1.9\times$ factor that includes Dafny verification and up to k LLM calls. On DafnyBench the factor is $2.5\times$ (25.6 s vs. 10.3 s). Given that the LLM baseline produces untrustworthy outputs on 8 to 10% of inputs, we believe this overhead is justified for the formal guarantees it provides.

6.2.2 Agentic Tool Use (τ^2 -bench). Table 3 reports performance on τ^2 -bench. The LLM baseline (Qwen3-8B without constraints) violates domain policies on 76.3% of retail interactions and 68.4% of airline interactions, yielding pass rates of only 11.3% and 13.2%, respectively. Both Agent-C and SEVerA reduce the violation rate to 0%, confirming that formal policy enforcement is necessary. SEVerA without constraints illustrates an intermediate point: by running the search loop without informing the planner of the policy specification, SEVerA (w/o constraints) improves pass rates to 49.4% (retail) and 44.7% (airline), but still incurs 10.3% and 25.5% violation rates, respectively. Enforcing constraints eliminates all violations while further improving pass rates to 53.6% and 52.6%. This gap reflects SEVerA’s search-verify-learn design. By synthesizing agent programs that are verified to satisfy the per-call policy specification, SEVerA can explore a wider space of compliant strategies rather than relying solely on runtime constraint checking. The per-call hard constraint ensures that every individual tool call is policy-compliant, while the search loop optimizes the agent’s overall strategy to maximize the end-of-trace pass rate. Impressively, SEVerA with the small open-weight Qwen3-8B outperforms Agent-C [20] with Claude Sonnet 4.5 (52.6 vs 47.3). SEVerA is also faster than Agent-C in both domains (212.4s vs. 234.7s in retail and 241.1s vs. 272.6s in airline). The LLM baseline is fastest because it performs no policy checking, but its outputs are unusable in practice given the violation rates.

6.2.3 Symbolic Math Synthesis (GSM-Symbolic). GSM-Symbolic uses Qwen3-8B as the underlying LLM. Qwen3-8B is open-weight, making parameter tuning feasible. This allows us to ask:

RQ2 (Training Impact). *Does parameter tuning improve task performance under formal constraints?*

Table 4 compares SEVerA against the LLM baseline and CRANE [4], a state-of-the-art constrained-decoding method, on GSM-Symbolic using Qwen3-8B. The LLM baseline (Qwen3-8B) achieves

Table 4. GSM-Symbolic results with Qwen3-8B. Higher accuracy is better; lower violation rate is better.

Method	Accuracy (%) \uparrow	Violation (%) \downarrow	Time (s)
LLM (Qwen3-8B)	38.3	10.6	10.9
CRANE [4]	44.7	2.1	12.4
SEVerA (no parameter tuning)	53.2 (+8.5)	0.0 (-2.1)	18.8
SEVerA (with parameter tuning)	66.0 (+21.3)	0.0 (-2.1)	16.7

Table 5. % of synthesised instances that violate behavioral specifications on the test set.

Method	ϵ (5%) \downarrow	ϵ (10%) \downarrow
PySR	62.86%	62.86%
LLM-SR	31.43%	34.29%
SEVerA	0.00%	0.00%

Table 6. NMSE values averaged over all instances that do not violate constraints on test data.

Comparison	$\epsilon = 0.05$		$\epsilon = 0.1$	
	SEVerA	Baseline	SEVerA	Baseline
vs PySR	0.0593	0.3113	0.0879	0.1217
vs LLM-SR	0.0200	0.1580	0.0205	0.1583

only 38.3% accuracy with a 10.6% constraint violation rate (**RQ1**), confirming that a vanilla LLM frequently produces outputs that are syntactically or semantically invalid. CRANE improves accuracy to 44.7% and reduces violations to 2.1% by augmenting the output grammar with reasoning tokens.¹SEVerA without parameter tuning already achieves 53.2% accuracy and eliminates all violations, an 8.5% accuracy improvement over CRANE with strictly stronger correctness guarantees. This improvement stems from the FGGM mechanism: the rejection sampler with a verified fallback ensures that every output satisfies the formal grammar and the semantic specification, while the synthesized prompting program f_p guides the underlying LLM toward correct answers. We also compare the performance of SEVerA with and without parameter tuning (**RQ2**). We fine-tune Qwen3-8B via GRPO with LoRA adapters under the same formal constraints. This provides a further 12.8% accuracy gain (from 53.2% to 66.0%) while maintaining the zero-violation guarantee. The tuned model produces outputs that more often pass the FGGM checker $check_{\mathcal{A}, \Phi, \Psi}$ on the first sample, reducing reliance on the fallback and thereby improving answer quality. Notably, the tuned variant is also faster than the untuned one (16.7 s vs. 18.8 s), because higher conformance reduces the number of rejection-sampling iterations needed per call.

6.2.4 Constrained Symbolic Regression. In constrained symbolic regression, we consider the SOTA genetic programming-based method PySR [8], as well as the existing self-evolving method LLM-SR [37], which samples parametric programs with optimizable PyTorch parameters and refines them to minimize loss. We denote the noise level by ϵ . As shown in Table 5, both baselines frequently produce programs that violate the behavioral specifications on the test data, with violations occurring in up to 62.86% of synthesis instances for PySR and up to 34.29% for LLM-SR. For both noise levels, SEVerA finds a verified solution in 33 out of 35 instances. Table 6 presents a pairwise comparison between SEVerA and the baselines in terms of normalized mean squared error (NMSE) on the test dataset. We compute the average NMSE over only those instances where the baseline-generated programs satisfy the behavioral constraints, pruning instances that correspond to invalid solutions. SEVerA achieves significantly lower NMSE values compared to both baselines.

6.3 Constraint Decomposition

The previous section showed that formal constraints improve both safety and performance. We now ask a finer-grained question:

¹CRANE only ensures that at every decoding step, the generated partial output is a valid prefix of the target grammar.

Table 7. Constraint decomposition ablation on GSM-Symbolic (Qwen3-8B). *Local* = FGGM contracts only; *Global* = behavioral specification only; *Full* = both.

Configuration	Accuracy (%) \uparrow	Violation (%) \downarrow
LLM (Qwen3-8B)	38.3	10.6
SEVerA (no parameter tuning)	53.2 (+14.9)	0.0 (-10.6)
SEVerA (local only parameter tuning)	55.3 (+17.0)	0.0 (-10.6)
SEVerA (global only parameter tuning)	61.7 (+23.4)	0.0 (-10.6)
SEVerA (full)	66.0 (+27.7)	0.0 (-10.6)

RQ3 (Constraint Decomposition). *What are the individual contributions of parameter tuning with local FGGM contracts versus global loss L?*

We ablate each component independently on GSM-Symbolic. Table 7 decomposes the 12.8% gain that parameter tuning provides (from 53.2% to 66.0%, as shown in Table 4) into its two sources. Starting from SEVerA without parameter tuning (53.2%), adding *local-only* tuning, which optimizes the FGGM conformance loss so that the model’s outputs more frequently satisfy per-call contracts, yields a modest improvement to 55.3%. This 2.1% gain indicates that parameter tuning alone helps the model produce syntactically valid outputs more reliably, reducing reliance on rejection sampling. *Global-only* tuning, which optimizes the task loss L without the FGGM conformance component, provides a substantially larger gain to 61.7%. This 8.5% improvement reflects the direct benefit of training the model to produce semantically correct answers. The full SEVerA pipeline combines both losses and achieves 66.0%, a 4.3% gain. This demonstrates that the two training signals are complementary: global tuning improves answer correctness while local conformance tuning ensures outputs satisfy per-call contracts, jointly yielding the best overall accuracy.

6.4 Example Synthesized Agent

We provide a complete agent program synthesised by SEVerA for the DafnyBench task in Appendix K. The agent defines three FGGMs: *initialFGGM*, *diffErrorFGGM*, and *verifierErrorFGGM*. Each sharing the local output contract $\Psi_l := \text{noDiff}(\text{base_program}, \cdot)$ and differing only in their prompting function f_p , which specialises the LLM prompt for initial annotation, diff-checker repair, and verifier-error repair, respectively. The resulting agent program is verified to provably comply with the diff-checker specification for all inputs and all parameter values, ensuring that no synthesised output modifies the original program beyond adding annotations.

6.5 Discussion

Across all four tasks, SEVerA achieves zero constraint violations on held-out test inputs while simultaneously improving task performance over every baseline. These results confirm that formal behavioral specifications do not merely enforce safety but actively prune the space of candidate programs, steering synthesis toward higher-quality agents. SEVerA’s runtime overhead is competitive: on τ^2 -bench it is *faster* than Agent-C, and on Dafny and GSM-Symbolic the full pipeline introduces a modest 1.9–2.5 \times slowdown relative to the LLM baseline, attributable to the verification and rejection-sampling steps. One current limitation is that the formulation is *resource-unaware*: behavioral specifications constrain functional correctness but do not account for computational resources such as LLM calls, token usage, or wall-clock budget. Extending SEVerA with resource-bound specifications [23], for example by encoding token or call budgets as additional hard constraints within the FGGM framework, is a promising direction for future work.

7 RELATED WORK

Self-Evolving Agents and Their Risks. Recent work has explored automatic agent synthesis and self-improvement, from code-based action unification [49] and open-ended skill libraries [47] to architecture search [13, 19, 53], gradient-free symbolic updates [51], co-evolution with world models [11], and reinforcement-learning-driven skill refinement [43, 48]. However, none of these provide formal behavioral guarantees, a gap made urgent by demonstrated risks of coding-agent exploits [14, 15, 25], test-case exploitation [50], and verification cheating [3]. SEVerA addresses this by imposing and verifying hard formal specifications before deployment.

Runtime Monitoring and Shielding. Rather than verifying at synthesis time, runtime approaches enforce safety on observed executions. Shield synthesis from neural policies [52] and temporal-logic runtime checking of tool-call sequences [20] are representative examples, but both are limited to observed traces and cannot guarantee safety on unseen inputs. SEVerA instead verifies the synthesized program for *all* inputs and parameter values at synthesis time. The monitoring checks can be used to define the specification as done with [20].

Neuro-Symbolic Program Synthesis. Neuro-symbolic methods, including neural program generation from examples [34], type-directed differentiable programming [46], wake-sleep library learning [10], vision-language program composition [41], and LLM-guided symbolic regression [37], synthesize programs that invoke neural components but lack formal correctness guarantees over the composed system. SEVerA extends this paradigm with FGGMs that bind each model call to a verified local contract, enabling end-to-end verification.

Constrained Decoding. Constrained decoding restricts LLM outputs to a formal language through grammar-based [35, 44, 45], diffusion-based [40], programming-abstraction [6, 21], reasoning-preserving [4], type-based [32], and semantic [33] techniques. These approaches require access to the model’s decoding internals (precluding closed-source models) and enforce per-token or per-output constraints, not program-level behavioral specifications. SEVerA’s FGGM mechanism operates on model outputs via rejection sampling, supports rich first-order logic specifications verified at the program level, and is model-agnostic.

Deductive Program Synthesis. Classical deductive synthesis, including syntax-guided [1], semantics-guided [22], counterexample-guided [38], abstraction-guided [16], and component-based [12] approaches, provides strong correctness guarantees but targets deterministic, non-parametric programs. SEVerA bridges this gap with a CEGIS-style search–verify loop for program structure while using FGGMs to encapsulate parametric generative models within verified contracts.

8 CONCLUSION

We presented SEVerA, a framework for synthesizing self-evolving LLM agents with formal behavioral guarantees. By introducing Formally Guarded Generative Models (FGGMs), SEVerA binds each generative model call to a verified local contract backed by a rejection sampler with a provably correct fallback, ensuring that specifications hold for all inputs and all parameter values. This design decomposes the constrained learning problem into a CEGIS-style discrete search over program structure followed by unconstrained gradient-based parameter optimization, retaining the scalability of modern fine-tuning methods such as GRPO while providing end-to-end correctness guarantees. Our evaluation across constrained symbolic regression, LLM-assisted Dafny verification, symbolic math synthesis, and policy-compliant agentic tool use demonstrates that SEVerA achieves zero constraint violations on all tasks while simultaneously improving task performance over unconstrained and state-of-the-art baselines. These results show that formal behavioral constraints are not merely a safety mechanism but an active guide that prunes the search space and steers synthesis toward higher-quality agents.

REFERENCES

- [1] Rajeev Alur, Rastislav Bodik, Garvit Juniwal, Milo M. K. Martin, Mukund Raghothaman, Sanjit A. Seshia, Rishabh Singh, Armando Solar-Lezama, Emina Torlak, and Abhishek Udupa. 2013. Syntax-guided synthesis. In *2013 Formal Methods in Computer-Aided Design*. 1–8. <https://doi.org/10.1109/FMCAD.2013.6679385>
- [2] Anthropic. 2025. System Card: Claude Sonnet 4.5. <https://www-cdn.anthropic.com/963373e433e489a87a10c823c52a0a013e9172dd.pdf>. Accessed: 2026-03-18.
- [3] Debangshu Banerjee, Olivier Bouissou, and Stefan Zetsche. 2026. DafnyPro: LLM-Assisted Automated Verification for Dafny Programs. arXiv:2601.05385 [cs.SE] <https://arxiv.org/abs/2601.05385>
- [4] Debangshu Banerjee, Tarun Suresh, Shubham Ugare, Sasa Misailovic, and Gagandeep Singh. 2025. CRANE: Reasoning with constrained LLM generation. In *Forty-second International Conference on Machine Learning*. <https://openreview.net/forum?id=wKs9ffHYxCV>
- [5] Victor Barres, Honghua Dong, Soham Ray, Xujie Si, and Karthik Narasimhan. 2025. τ^2 -Bench: Evaluating Conversational Agents in a Dual-Control Environment. arXiv:2506.07982 [cs.AI] <https://arxiv.org/abs/2506.07982>
- [6] Luca Beurer-Kellner, Marc Fischer, and Martin Vechev. 2023. Prompting Is Programming: A Query Language for Large Language Models. *Proc. ACM Program. Lang.* 7, PLDI, Article 186 (June 2023), 24 pages. <https://doi.org/10.1145/3591300>
- [7] Iwo Bładek and Krzysztof Krawiec. 2019. Solving symbolic regression problems with formal constraints. In *Proceedings of the Genetic and Evolutionary Computation Conference (Prague, Czech Republic) (GECCO '19)*. Association for Computing Machinery, New York, NY, USA, 977–984. <https://doi.org/10.1145/3321707.3321743>
- [8] Miles Cranmer. 2023. Interpretable Machine Learning for Science with PySR and SymbolicRegression.jl. arXiv:2305.01582 [astro-ph.IM] <https://arxiv.org/abs/2305.01582>
- [9] Dafny Language Community. 2023. Integrating Dafny and Python Code. <https://dafny.org/v3.10.0/DafnyRef/integration-py/IntegrationPython>. Accessed: 2026-03-17.
- [10] Kevin Ellis, Catherine Wong, Maxwell Nye, Mathias Sablé-Meyer, Lucas Morales, Luke Hewitt, Luc Cary, Armando Solar-Lezama, and Joshua B. Tenenbaum. 2021. DreamCoder: bootstrapping inductive program synthesis with wake-sleep library learning. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation (Virtual, Canada) (PLDI 2021)*. Association for Computing Machinery, New York, NY, USA, 835–850. <https://doi.org/10.1145/3453483.3454080>
- [11] Tianqing Fang, Hongming Zhang, Zhisong Zhang, Kaixin Ma, Wenhao Yu, Haitao Mi, and Dong Yu. 2025. WebEvolver: Enhancing Web Agent Self-Improvement with Coevolving World Model. *arXiv preprint arXiv:2504.21024* (2025).
- [12] Yu Feng, Ruben Martins, Jacob Van Geffen, Isil Dillig, and Swarat Chaudhuri. 2017. Component-based synthesis of table consolidation and transformation tasks from examples. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation (Barcelona, Spain) (PLDI 2017)*. Association for Computing Machinery, New York, NY, USA, 422–436. <https://doi.org/10.1145/3062341.3062351>
- [13] Hongcheng Gao, Yue Liu, Yufei He, Longxu Dou, Chao Du, Zhijie Deng, Bryan Hooi, Min Lin, and Tianyu Pang. 2025. FlowReasoner: Reinforcing Query-Level Meta-Agents. arXiv:2504.15257 [cs.AI] <https://arxiv.org/abs/2504.15257>
- [14] The Guardian. 2026. Rogue AI agents published passwords and bypassed security protections. News investigation.
- [15] Chengquan Guo, Xun Liu, Chulin Xie, Andy Zhou, Yi Zeng, Zinan Lin, Dawn Song, and Bo Li. 2024. RedCode: Risky Code Execution and Generation Benchmark for Code Agents. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*. <https://openreview.net/forum?id=mAG68wdggA>
- [16] Sankha Narayan Guria, Jeffrey S. Foster, and David Van Horn. 2023. Absynthe: Abstract Interpretation-Guided Synthesis. *Proc. ACM Program. Lang.* 7, PLDI, Article 171 (June 2023), 24 pages. <https://doi.org/10.1145/3591285>
- [17] Christian Haider and Gabriel Kronberger. 2022. Shape-Constrained Symbolic Regression with NSGA-III. In *Computer Aided Systems Theory – EUROCAST 2022: 18th International Conference, Las Palmas de Gran Canaria, Spain, February 20–25, 2022, Revised Selected Papers* (Las Palmas de Gran Canaria, Spain). Springer-Verlag, Berlin, Heidelberg, 164–172. https://doi.org/10.1007/978-3-031-25312-6_19
- [18] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Liang Wang, Weizhu Chen, et al. 2022. Lora: Low-rank adaptation of large language models. *Iclr* 1, 2 (2022), 3.
- [19] Shengran Hu, Cong Lu, and Jeff Clune. 2025. Automated Design of Agentic Systems. In *The Thirteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=t9U3LW7JVX>
- [20] Adharsh Kamath, Sishen Zhang, Calvin Xu, Shubham Ugare, Gagandeep Singh, and Sasa Misailovic. 2025. Enforcing Temporal Constraints for LLM Agents. arXiv:2512.23738 [cs.PL] <https://arxiv.org/abs/2512.23738>
- [21] Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan A, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller, Matei Zaharia, and Christopher Potts. 2024. DSPy: Compiling Declarative Language Model Calls into State-of-the-Art Pipelines. In *The Twelfth International Conference on Learning Representations*. <https://openreview.net/forum?id=sY5N0zY5Od>
- [22] Jinwoo Kim, Qinheping Hu, Loris D’Antoni, and Thomas Reps. 2021. Semantics-guided synthesis. *Proc. ACM Program. Lang.* 5, POPL, Article 30 (Jan. 2021), 32 pages. <https://doi.org/10.1145/3434311>

- [23] Tristan Knoth, Di Wang, Nadia Polikarpova, and Jan Hoffmann. 2019. Resource-guided program synthesis. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Phoenix, AZ, USA) (PLDI 2019). Association for Computing Machinery, New York, NY, USA, 253–268. <https://doi.org/10.1145/3314221.3314602>
- [24] G. Kronberger, F. O. de Franca, B. Burlacu, C. Haider, and M. Kommenda. 2022. Shape-Constrained Symbolic Regression—Improving Extrapolation with Prior Knowledge. *Evolutionary Computation* 30, 1 (03 2022), 75–98. https://doi.org/10.1162/evco_a_00294 arXiv:https://direct.mit.edu/evco/article-pdf/30/1/75/1995582/evco_a_00294.pdf
- [25] Eunkyoo Lee, Donghyeon Kim, Wonyoung Kim, and Insu Yun. 2025. Takedown: How It’s Done in Modern Coding Agent Exploits. arXiv:2509.24240 [cs.CR] <https://arxiv.org/abs/2509.24240>
- [26] K. Rustan M. Leino. 2010. Dafny: an automatic program verifier for functional correctness. In *Proceedings of the 16th International Conference on Logic for Programming, Artificial Intelligence, and Reasoning* (Dakar, Senegal) (LPAR’10). Springer-Verlag, Berlin, Heidelberg, 348–370.
- [27] Li Li, Minjie Fan, Rishabh Singh, and Patrick Riley. 2019. Neural-Guided Symbolic Regression with Asymptotic Constraints. arXiv:1901.07714 [cs.LG] <https://arxiv.org/abs/1901.07714>
- [28] Wenqiang Li, Weijun Li, Linjun Sun, Min Wu, Lina Yu, Jingyi Liu, Yanjie Li, and Songsong Tian. 2023. Transformer-based model for symbolic regression via joint supervised learning. In *The Eleventh International Conference on Learning Representations*. <https://openreview.net/forum?id=ULzyv9M1j5>
- [29] Chloe R 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. *Transactions on Machine Learning Research* (2025). <https://openreview.net/forum?id=yBgTVWccIX>
- [30] Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. 2024. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models. *arXiv preprint arXiv:2410.05229* (2024).
- [31] Eric Mugnier, Emmanuel Anaya Gonzalez, Nadia Polikarpova, Ranjit Jhala, and Zhou Yuanyuan. 2025. Laurel: Unblocking Automated Verification with Large Language Models. *Proc. ACM Program. Lang.* 9, OOPSLA1, Article 134 (April 2025), 27 pages. <https://doi.org/10.1145/3720499>
- [32] Niels Müндler, Jingxuan He, Hao Wang, Koushik Sen, Dawn Song, and Martin Vechev. 2025. Type-Constrained Code Generation with Language Models. *Proc. ACM Program. Lang.* 9, PLDI, Article 171 (June 2025), 26 pages. <https://doi.org/10.1145/3729274>
- [33] Shaan Nagy, Timothy Zhou, Nadia Polikarpova, and Loris D’Antoni. 2026. ChopChop: A Programmable Framework for Semantically Constraining the Output of Language Models. *Proc. ACM Program. Lang.* 10, POPL, Article 66 (Jan. 2026), 28 pages. <https://doi.org/10.1145/3776708>
- [34] Emilio Parisotto, Abdel rahman Mohamed, Rishabh Singh, Lihong Li, Dengyong Zhou, and Pushmeet Kohli. 2016. Neuro-Symbolic Program Synthesis. arXiv:1611.01855 [cs.AI] <https://arxiv.org/abs/1611.01855>
- [35] Kanghee Park, Jiayu Wang, Taylor Berg-Kirkpatrick, Nadia Polikarpova, and Loris D’Antoni. 2024. Grammar-Aligned Decoding. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*. <https://openreview.net/forum?id=5G7ve8E1Lu>
- [36] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. arXiv:2402.03300 [cs.CL] <https://arxiv.org/abs/2402.03300>
- [37] Parshin Shojaee, Kazem Meidani, Shashank Gupta, Amir Barati Farimani, and Chandan K. Reddy. 2025. LLM-SR: Scientific Equation Discovery via Programming with Large Language Models. In *The Thirteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=m2nmp8P5in>
- [38] Armando Solar-Lezama. 2013. Program Sketching. *International Journal on Software Tools for Technology Transfer* 15, 5 (2013), 475–495. <https://doi.org/10.1007/s10009-012-0249-7>
- [39] Chuyue Sun, Ying Sheng, Oded Padon, and Clark Barrett. 2024. Clover: Closed-Loop Verifiable Code Generation. In *AI Verification: First International Symposium, SAIV 2024, Montreal, QC, Canada, July 22–23, 2024, Proceedings* (Montreal, QC, Canada). Springer-Verlag, Berlin, Heidelberg, 134–155. https://doi.org/10.1007/978-3-031-65112-0_7
- [40] Tarun Suresh, Debangshu Banerjee, Shubham Ugare, Sasa Misailovic, and Gagandeep Singh. 2025. DINGO: Constrained Inference for Diffusion LLMs. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*. <https://openreview.net/forum?id=KaYMGsnZ4R>
- [41] Didac Suris, Sachit Menon, and Carl Vondrick. 2023. ViperGPT: Visual Inference via Python Execution for Reasoning. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*. 11854–11864. <https://doi.org/10.1109/ICCV51070.2023.01092>
- [42] Qwen Team. 2025. Qwen3 Technical Report. arXiv:2505.09388 [cs.CL] <https://arxiv.org/abs/2505.09388>
- [43] Wanxin Tian, Shijie Zhang, Kevin Zhang, Xiaowei Chi, Chunkai Fan, Junyu Lu, Yulin Luo, Qiang Zhou, Yiming Zhao, Ning Liu, Siyu Lin, Zhiyuan Qin, Xiaozhu Ju, Shanghang Zhang, and Jian Tang. 2025. SEEA-R1: Tree-Structured Reinforcement Fine-Tuning for Self-Evolving Embodied Agents. arXiv:2506.21669 [cs.AI] <https://arxiv.org/abs/2506.21669>

21669

- [44] Shubham Ugare, Rohan Gumaste, Tarun Suresh, Gagandeep Singh, and Sasa Misailovic. 2025. IterGen: Iterative Semantic-aware Structured LLM Generation with Backtracking. In *The Thirteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=ac93gRzxxV>
- [45] Shubham Ugare, Tarun Suresh, Hangoo Kang, Sasa Misailovic, and Gagandeep Singh. 2025. SynCode: LLM Generation with Grammar Augmentation. *Transactions on Machine Learning Research* (2025). <https://openreview.net/forum?id=HiUZtgAPoH>
- [46] Lazar Valkov, Dipak Chaudhari, Akash Srivastava, Charles Sutton, and Swarat Chaudhuri. 2018. HOUDINI: lifelong learning as program synthesis. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems (Montréal, Canada) (NIPS'18)*. Curran Associates Inc., Red Hook, NY, USA, 8701–8712.
- [47] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023. Voyager: An Open-Ended Embodied Agent with Large Language Models. arXiv:2305.16291 [cs.AI] <https://arxiv.org/abs/2305.16291>
- [48] Jiong Xiao Wang, Qiaoqing Yan, Yawei Wang, Yijun Tian, Soumya Smruti Mishra, Zhichao Xu, Megha Gandhi, Panpan Xu, and Lin Lee Cheong. 2026. Reinforcement Learning for Self-Improving Agent with Skill Library. arXiv:2512.17102 [cs.AI] <https://arxiv.org/abs/2512.17102>
- [49] Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. 2024. Executable code actions elicit better LLM agents. In *Proceedings of the 41st International Conference on Machine Learning (Vienna, Austria) (ICML '24)*. JMLR.org, Article 2054, 25 pages.
- [50] Ziqian Zhong, Aditi Raghunathan, and Nicholas Carlini. 2025. ImpossibleBench: Measuring LLMs' Propensity of Exploiting Test Cases. *arXiv preprint arXiv:2510.20270* (2025).
- [51] Wangchunshu Zhou, Yixin Ou, Shengwei Ding, Long Li, Jialong Wu, Tiannan Wang, Jiamin Chen, Shuai Wang, Xiaohua Xu, Ningyu Zhang, Huajun Chen, and Yuchen Eleanor Jiang. 2024. Symbolic Learning Enables Self-Evolving Agents. arXiv:2406.18532 [cs.CL] <https://arxiv.org/abs/2406.18532>
- [52] He Zhu, Zikang Xiong, Stephen Magill, and Suresh Jagannathan. 2019. An inductive synthesis framework for verifiable reinforcement learning. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation (Phoenix, AZ, USA) (PLDI 2019)*. Association for Computing Machinery, New York, NY, USA, 686–701. <https://doi.org/10.1145/3314221.3314638>
- [53] Mingchen Zhuge, Wenyi Wang, Louis Kirsch, Francesco Faccio, Dmitrii Khizbullin, and Jürgen Schmidhuber. 2024. GPTSwarm: Language Agents as Optimizable Graphs. In *Forty-first International Conference on Machine Learning*. <https://openreview.net/forum?id=uTC9AFXlhg>

A RESTRICTED DAFNY GRAMMAR

$\langle \text{program} \rangle ::= \langle \text{method_decl} \rangle \mid \langle \text{function_decl} \rangle$
 $\langle \text{method_decl} \rangle ::= \text{method } \langle \text{ident} \rangle (\langle \text{params} \rangle) \langle \text{specs} \rangle \{ \langle \text{stmt_list} \rangle \}$
 $\langle \text{function_decl} \rangle ::= \text{function } \langle \text{ident} \rangle (\langle \text{params} \rangle) : \langle \text{type} \rangle \langle \text{specs} \rangle \langle \text{expr} \rangle$
 $\langle \text{params} \rangle ::= \epsilon \mid \langle \text{param} \rangle \mid \langle \text{param} \rangle, \langle \text{params} \rangle$
 $\langle \text{param} \rangle ::= \langle \text{ident} \rangle : \langle \text{type} \rangle$
 $\langle \text{type} \rangle ::= \text{int} \mid \text{bool} \mid \text{string} \mid \text{real}$
 $\langle \text{specs} \rangle ::= \epsilon \mid \langle \text{spec} \rangle \langle \text{specs} \rangle$
 $\langle \text{spec} \rangle ::= \text{requires } \langle \text{formula} \rangle \mid \text{ensures } \langle \text{formula} \rangle$
 $\langle \text{stmt_list} \rangle ::= \epsilon \mid \langle \text{stmt} \rangle \langle \text{stmt_list} \rangle$
 $\langle \text{stmt} \rangle ::= \langle \text{assign_stmt} \rangle \mid \langle \text{if_stmt} \rangle \mid \langle \text{while_stmt} \rangle \mid \langle \text{assert_stmt} \rangle \mid \langle \text{call_stmt} \rangle$
 $\langle \text{assign_stmt} \rangle ::= \langle \text{ident} \rangle := \langle \text{expr} \rangle;$
 $\langle \text{call_stmt} \rangle ::= \langle \text{ident} \rangle (\langle \text{args} \rangle);$
 $\langle \text{args} \rangle ::= \epsilon \mid \langle \text{expr} \rangle \mid \langle \text{expr} \rangle, \langle \text{args} \rangle$
 $\langle \text{if_stmt} \rangle ::= \text{if } (\langle \text{formula} \rangle) \{ \langle \text{stmt_list} \rangle \}$
 $\quad \mid \text{if } (\langle \text{formula} \rangle) \{ \langle \text{stmt_list} \rangle \} \text{ else } \{ \langle \text{stmt_list} \rangle \}$
 $\langle \text{while_stmt} \rangle ::= \text{while } (\langle \text{formula} \rangle) \langle \text{loop_annots} \rangle \{ \langle \text{stmt_list} \rangle \}$
 $\langle \text{loop_annots} \rangle ::= \epsilon \mid \langle \text{loop_annot} \rangle \langle \text{loop_annots} \rangle$
 $\langle \text{loop_annot} \rangle ::= \text{invariant } \langle \text{formula} \rangle \mid \text{decreases } \langle \text{expr} \rangle$
 $\langle \text{assert_stmt} \rangle ::= \text{assert } \langle \text{formula} \rangle;$
 $\langle \text{expr} \rangle ::= \langle \text{aexpr} \rangle \mid \langle \text{sexpr} \rangle$
 $\langle \text{aexpr} \rangle ::= \langle \text{int_lit} \rangle \mid \langle \text{ident} \rangle \mid \langle \text{aexpr} \rangle \langle \text{arith_op} \rangle \langle \text{aexpr} \rangle \mid (\langle \text{aexpr} \rangle) \mid \langle \text{ident} \rangle (\langle \text{args} \rangle)$
 $\langle \text{sexpr} \rangle ::= \langle \text{string_lit} \rangle \mid \langle \text{ident} \rangle \mid (\langle \text{sexpr} \rangle) \mid \langle \text{ident} \rangle (\langle \text{args} \rangle)$
 $\langle \text{formula} \rangle ::= \langle \text{expr} \rangle \langle \text{rel_op} \rangle \langle \text{expr} \rangle$
 $\quad \mid ! \langle \text{formula} \rangle$
 $\quad \mid (\langle \text{formula} \rangle)$
 $\quad \mid \langle \text{ident} \rangle (\langle \text{args} \rangle)$
 $\quad \mid \text{forall } \langle \text{bound_vars} \rangle :: \langle \text{formula} \rangle$
 $\quad \mid \text{exists } \langle \text{bound_vars} \rangle :: \langle \text{formula} \rangle$
 $\langle \text{bound_vars} \rangle ::= \langle \text{ident} \rangle : \langle \text{type} \rangle \mid \langle \text{ident} \rangle : \langle \text{type} \rangle, \langle \text{bound_vars} \rangle$
 $\langle \text{arith_op} \rangle ::= + \mid - \mid * \mid /$
 $\langle \text{rel_op} \rangle ::= = \mid \neq \mid < \mid > \mid \leq \mid \geq$
 $\langle \text{literal} \rangle ::= \langle \text{int_lit} \rangle \mid \langle \text{string_lit} \rangle \mid \text{true} \mid \text{false}$
 $\langle \text{int_lit} \rangle ::= 0 \mid 1 \mid \dots \mid 9 \text{ (and sequences thereof)}$
 $\langle \text{string_lit} \rangle ::= \text{''} \langle \text{char} \rangle^* \text{''}$
 $\langle \text{ident} \rangle ::= \langle \text{letter} \rangle \mid \langle \text{letter} \rangle \langle \text{ident} \rangle$
 $\langle \text{letter} \rangle ::= a \mid b \mid \dots \mid z \mid A \mid B \mid \dots \mid Z$
 $\langle \text{digit} \rangle ::= 0 \mid 1 \mid \dots \mid 9$

B REJECTION SAMPLING DETAILS

Rejection Sampling Details. Given a target distribution π_t and a proposal distribution π_p with $S(\pi_t) \subseteq S(\pi_p)$, assume there exists $M \geq 1$ such that $\forall z \in \Omega, D_{\pi_t}(z) \leq M \cdot D_{\pi_p}(z)$. The procedure samples $z \sim \pi_p$ and $u \sim \text{Uniform}(0, 1)$, and accepts z iff $u \leq \frac{D_{\pi_t}(z)}{M \cdot D_{\pi_p}(z)}$. Otherwise, the sample is rejected and the process repeats. Accepted samples follow π_t . Moreover, any $z \notin S(\pi_t)$ is always rejected since $D_{\pi_t}(z) = 0$, ensuring all accepted samples lie in $S(\pi_t)$. The acceptance probability is $1/M$, so efficiency depends on how tightly π_p upper-bounds π_t . For this paper, we assume $M = 1$.

C LIBRARY FUNCTIONS FOR SYMBOLIC REGRESSION

Library Function Set \mathcal{F} .

- $\text{abs}(x : \text{real}) \rightarrow r : \text{real}$ ensures $r \geq 0$ ensures $r = x \vee r = -x$
- $\text{max}(x : \text{real}, y : \text{real}) \rightarrow r : \text{real}$ ensures $(x \leq r) \wedge (y \leq r)$ ensures $(r = x) \vee (r = y)$
- $\text{min}(x : \text{real}, y : \text{real}) \rightarrow r : \text{real}$ ensures $(x \geq r) \wedge (y \geq r)$ ensures $(r = x) \vee (r = y)$
- $\text{sin}(x : \text{real}) \rightarrow r : \text{real}$ ensures $-1 \leq r \leq 1$ ensures $(x = 0) \Rightarrow (r = 0)$
- $\text{cos}(x : \text{real}) \rightarrow r : \text{real}$ ensures $-1 \leq r \leq 1$ ensures $(x = 0) \Rightarrow (r = 1)$
- $\text{pi}(x : \text{real}) \rightarrow r : \text{real}$ requires $x \geq 0$ ensures $3.141592653589790 \cdot x \leq r$ ensures $r \leq 3.141592653589793 \cdot x$
- $\text{pow}(x : \text{real}, d : \text{real}) \rightarrow r : \text{real}$ requires $x \geq 0$ requires $(x \neq 0) \vee (d \geq 0)$ ensures $r \geq 0$ ensures $(x > 0) \Rightarrow (r > 0)$ ensures $(d = 0) \Rightarrow (r = 1)$ ensures $(x \geq 1 \wedge d \geq 0) \Rightarrow (r \geq 1)$ ensures $(x \geq 1 \wedge d \leq 0) \Rightarrow (r \leq 1)$ ensures $(x \leq 1 \wedge d \geq 0) \Rightarrow (r \leq 1)$ ensures $(x \leq 1 \wedge d \leq 0) \Rightarrow (r \geq 1)$
- $\text{sqrt}(x : \text{real}) \rightarrow r : \text{real}$ requires $x \geq 0$ ensures $r \geq 0$ ensures $r \cdot r = x$
- $\text{exp}(x : \text{real}) \rightarrow r : \text{real}$ ensures $r \geq 0$ ensures $(x \leq 0) \Rightarrow (r \leq 1)$ ensures $(x \geq 0) \Rightarrow (r \geq 1)$ ensures $(x = 1) \Rightarrow (2.71 \leq r \leq 2.72)$ ensures $r \geq (1 + x)$
- $\text{log}(x : \text{real}) \rightarrow r : \text{real}$ requires $x > 0$ ensures $(x \geq 1) \Rightarrow (r \geq 0)$ ensures $(x \leq 1) \Rightarrow (r \leq 0)$ ensures $(x = 1) \Rightarrow (r = 0)$
- $\text{neural1}(x : \text{real}) \rightarrow r : \text{real}$
- $\text{neural2}(x1 : \text{real}, x2 : \text{real}) \rightarrow r : \text{real}$
- $\text{neural3}(x1 : \text{real}, x2 : \text{real}, x3 : \text{real}) \rightarrow r : \text{real}$

D ADDITIONAL AXIOMS FOR SYMBOLIC REGRESSION

D.1 Axiom Syntax

$$\langle \text{axiom} \rangle ::= \langle \text{formula} \rangle$$

Axiom Set \mathcal{A} .

- **Interval Multiplication Rule:** $\forall x, y, x_{lb}, x_{ub}, y_{lb}, y_{ub} \in \mathbb{R}. (x_{lb} \leq x \leq x_{ub}) \wedge (y_{lb} \leq y \leq y_{ub}) \Rightarrow$

$$\text{min}(x_{lb} \times y_{lb}, x_{lb} \times y_{ub}, x_{ub} \times y_{lb}, x_{ub} \times y_{ub}) \leq x \times y \leq \text{max}(x_{lb} \times y_{lb}, x_{lb} \times y_{ub}, x_{ub} \times y_{lb}, x_{ub} \times y_{ub}).$$

- $\forall x, d_1, d_2 \in \mathbb{R}. (0 \leq x \leq 1) \wedge (d_1 \leq d_2) \Rightarrow \text{pow}(x, d_1) \geq \text{pow}(x, d_2).$
- $\forall x, d_1, d_2 \in \mathbb{R}. (x \geq 1) \wedge (d_1 \leq d_2) \Rightarrow \text{pow}(x, d_1) \leq \text{pow}(x, d_2).$
- $\forall x \in \mathbb{R}. \text{exp}(x) = \text{pow}(\text{exp}(1), x).$
- $\forall x \in \mathbb{R}. (x > 0) \Rightarrow (\forall r \in \mathbb{R}. r = \text{log}(x) \iff \text{exp}(r) = x).$
- $\forall x \in \mathbb{R}. (x \geq 0) \Rightarrow (\text{sqrt}(x) = \text{pow}(x, 0.5)).$

E LIBRARY FUNCTIONS FOR DAFNY PROGRAM VERIFICATION

Library Function Set \mathcal{F} .

- $claude(prompt : str) \rightarrow r : str$
- $noDiff(base_code : str, annotated_code : str) \rightarrow r : bool$
- $noDiffWithErrorMsg(base_code : str, annotated_code : str) \rightarrow r : str$
- $dafnyVerifier(dafny_code : str) \rightarrow r : bool$
- $dafnyVerifierWithErrorMsg(dafny_code : str) \rightarrow r : str$
- $extractDafnyCode(text : str) \rightarrow r : str$
- $isSubstring(sub : str, full : str) \rightarrow r : bool$

F ADDITIONAL AXIOMS FOR DAFNY PROGRAM VERIFICATION

Axiom Set \mathcal{A} .

- **Reflexivity:** $\forall a, b \in \Sigma^*. (a = b) \implies noDiff(a, b)$.
- **Transitivity:** $\forall a, b, c \in \Sigma^*. noDiff(a, b) \wedge noDiff(b, c) \implies noDiff(a, c)$.
- $\forall a, b \in \Sigma^*. noDiffWithErrorMsg(a, b) = \text{“”} \iff noDiff(a, b)$.
- $\forall c \in \Sigma^*. dafnyVerifierWithErrorMsg(c) = \text{“”} \iff dafnyVerifier(c)$.

G LIBRARY FUNCTIONS FOR τ^2 -BENCH

Library Function Set \mathcal{F} .

- $qwen(prompt : str) \rightarrow r : str$
- $agentCCheck(tool : str, trace : str, domain : str) \rightarrow r : bool$
- $isToolCall(response : str) \rightarrow r : bool$
- $parseToolCallName(response : str) \rightarrow r : str$

H ADDITIONAL AXIOMS FOR τ^2 -BENCH

Axiom Set \mathcal{A} .

- **Non-tool-call safety:** $\forall r, t, d \in \Sigma^*. \neg isToolCall(r) \implies agentCCheck(r, t, d)$.
- **Transfer-to-human satisfiability:** $\forall r, t, d \in \Sigma^*. isToolCall(r) \wedge parseToolCallName(r) = \text{“transfer_to_human_agents”} \implies agentCCheck(r, t, d)$.

I LIBRARY FUNCTIONS FOR GSM-SYMBOLIC

Library Function Set \mathcal{F} .

- $qwen(prompt : str) \rightarrow r : str$
- $gsmParser(expression : str) \rightarrow r : bool$
- $gsmParserWithErrorMsg(expression : str) \rightarrow r : str$
- $extractExpression(text : str) \rightarrow r : str$

J ADDITIONAL AXIOMS FOR GSM-SYMBOLIC

Axiom Set \mathcal{A} .

- $gsmParser(\text{“} \ll \gg \text{”})$.
- $\forall e \in \Sigma^*. gsmParserWithErrorMsg(e) = \text{“”} \iff gsmParser(e)$.

K EXAMPLE AGENTS FOR LLM ASSISTED PROGRAM VERIFICATION

The following is an example verified agent program for the Dafny program verification task. The agent defines and invokes three FGGMs—*initialFGGM* (Fig. 4), *diffErrorFGGM* (Fig. 5), and *verifierErrorFGGM* (Fig. 6)—each with local contract $\Psi_i : noDiff(base_program, \cdot)$ and a fallback that returns the original program (justified by the reflexivity axiom). All three share $GMid := claude$ and differ only in the prompting function f_p : the first generates an initial annotation with context-specific guidance, the second repairs after a diff-checker failure, and the third repairs

after a verification failure with error-specific feedback. The agent program (Fig. 7) dispatches to the appropriate FGGM on each loop iteration.

FGGM Definition:

```

id := initialFGGM
GMid := claude
typeSig := (base_program : str) → str
ΦI := true
ΨI := noDifff(base_program, f(base_program))
fp := function fp(base_program : str) : (str) {
  var g : str :=
    "You are an expert in Dafny verification. Given the following unannotated Dafny code,"
    + "add appropriate loop invariants, assertions, and decreases clauses to make it verify.\n\n"
    + "IMPORTANT RULES:\n"
    + "1. Do NOT modify the method signatures, requires clauses, or ensures clauses\n"
    + "2. Only add invariants, assertions, and decreases clauses\n"
    + "3. Do NOT change the logic or control flow\n"
    + "4. Keep invariants MINIMAL - only what's necessary for postconditions\n"
    + "5. Add decreases clauses for termination\n";
  if (isSubstring("lemma", base_program)) {
    g := g
    + "6. LEMMAS: If a lemma has a non-trivial postcondition, add explicit proof steps:\n"
    + "  - Use assertions to unfold function definitions\n"
    + "  - Add case analysis or pattern matching if needed\n"
    + "  - Call other lemmas to establish intermediate facts\n"
    + "  - Do NOT leave lemma bodies empty unless postcondition is trivial\n"; }
  if (isSubstring("<", base_program) ∧ isSubstring(">", base_program)) {
    g := g
    + "7. GENERICS: When calling generic functions/lemmas, propagate type constraints:\n"
    + "  - If a called function requires T(00), declare it in the caller too\n"
    + "  - Explicitly instantiate type parameters when needed\n"; }
  if (isSubstring("set<", base_program) ∨ isSubstring("set", base_program)) {
    g := g
    + "8. SETS: Use simple, direct assertions for set reasoning:\n"
    + "  - Assert set equality decompositions: A == B + (A - B)\n"
    + "  - Use cardinality constraints: |A| == |B| + |C|\n"
    + "  - Avoid complex witness variables or case analysis\n"
    + "  - Let Dafny's built-in set theory automation work\n"; }
  if (isSubstring("multiset", base_program)) {
    g := g
    + "9. MULTISSETS: Match the structure of postconditions directly\n"
    + "  - Use simple decomposition assertions\n"
    + "  - Avoid unnecessary intermediate steps\n"; }
  return g + "\nBase program:\n" + base_program
    + "\n\nReturn ONLY the annotated Dafny code in a ```dafny code block."; }
fd := function fd(base_program : str, y : str) : (str)
  ensures noDifff(base_program, fd(base_program, y))
{
  return base_program;
}

```

Fig. 4. FGGM definition for *initialFGGM*. The prompting function builds context-specific guidance by checking for lemmas, generics, sets, and multisets in the base program. The fallback returns the original program, satisfying Ψ_I by the reflexivity axiom.

FGGM Definition:

```

id := diffErrorFGGM
GMid := claude
typeSig := (base_program : str, diff_error : str, prev_attempt : str) → str
Φl := true
Ψl := noDiff(base_program, f(base_program, diff_error, prev_attempt))
fp := function fp(base_program : str, diff_error : str, prev_attempt : str) : (str) {
  return
    "The previous Dafny annotation modified the base program logic incorrectly.\n\n"
    + "noDiff error:\n" + diff_error + "\n\n"
    + "Base program:\n" + base_program + "\n\n"
    + "Previous attempt:\n" + prev_attempt + "\n\n"
    + "Please add annotations WITHOUT changing:\n"
    + "- Method signatures\n"
    + "- Requires/ensures clauses\n"
    + "- Program logic or control flow\n"
    + "Only add invariants, assertions, and decreases clauses.\n"
    + "Return ONLY the corrected Dafny code in a ``` dafny code block."; }
fd := function fd(base_program : str, diff_error : str, prev_attempt : str, y : str) : (str)
  ensures noDiff(base_program, fd(base_program, diff_error, prev_attempt, y))
{
  return base_program;
}

```

Fig. 5. FGGM definition for *diffErrorFGGM*. Invoked when the previous iteration’s output failed the diff checker. The prompt includes the diff error and previous attempt, instructing the model to preserve the base program logic.

```

FGGM Definition:
id := verifierErrorFGGM
GMid := claude
typeSig := (base_program : str, verify_error : str, prev_attempt : str) → str
ΦI := true
ΨI := noDiff(base_program, f(base_program, verify_error, prev_attempt))
fp := function fp(base_program : str, verify_error : str, prev_attempt : str) : (str) {
  var sg : str :=
    "Focus on:\n"
    + "- Strengthening invariants MINIMALLY\n"
    + "- Adding simple decomposition assertions\n"
    + "- Ensuring decreases clauses are correct\n";

  if (isSubstring("timeout", verify_error) ∨ isSubstring("resource", verify_error)) {
    sg := "TIMEOUT/RESOURCE ISSUE: Break down the proof into smaller steps:\n"
        + "- Add strategic intermediate assertions\n"
        + "- Call helper lemmas to establish sub-goals\n"
        + "- Simplify complex invariants\n"; }

  if (isSubstring("lemma", verify_error) ∨ isSubstring("postcondition", verify_error)) {
    sg := "LEMMA PROOF ISSUE: Add explicit proof steps in lemma body:\n"
        + "- Assert intermediate facts that lead to postcondition\n"
        + "- Unfold recursive function definitions\n"
        + "- Use case analysis if needed\n"; }

  if (isSubstring("invariant", verify_error)) {
    sg := "INVARIANT ISSUE: Adjust loop invariants:\n"
        + "- Ensure invariants are maintained after each iteration\n"
        + "- Add bounds and relationships incrementally\n"
        + "- Keep invariants minimal but sufficient\n"; }

  if (isSubstring("type", verify_error) ∨ isSubstring("constraint", verify_error)) {
    sg := "TYPE CONSTRAINT ISSUE: Check generic type parameters:\n"
        + "- Propagate nonemptiness constraints like T(00)\n"
        + "- Ensure type parameters match between caller and callee\n"; }

  return
    "The previous Dafny annotation preserved the base logic but failed verification.\n\n"
    + "Verification error:\n" + verify_error + "\n\n"
    + sg + "\n"
    + "Base program:\n" + base_program + "\n\n"
    + "Previous attempt:\n" + prev_attempt + "\n\n"
    + "Please fix the annotations. Do NOT modify method signatures, requires, or ensures clauses.\n"
    + "Return ONLY the corrected Dafny code in a `` ` ` dafny code block"; }

fd := function fd(base_program : str, verify_error : str, prev_attempt : str, y : str) : (str)
  ensures noDiff(base_program, fd(base_program, verify_error, prev_attempt, y))
{
  return base_program;
}

```

Fig. 6. FGGM definition for *verifierErrorFGGM*. Invoked when the previous iteration passed the diff checker but failed verification. The prompting function analyzes the verification error to provide targeted guidance (timeout, lemma, invariant, or type issues).

```

method agent(base_program : str) returns (r : str)
  ensures noDiff(base_program, r)
{
  var max_attempts : int := 5;
  var attempt : int := 0;
  var best : str := base_program;
  var best_verified : bool := false;

  var prev_diff_err : str := "";
  var prev_ver_err : str := "";
  var prev_code : str := "";

  while (attempt < max_attempts)
    invariant noDiff(base_program, best)
    decreases max_attempts - attempt
  {
    var y : str;
    if (attempt = 0) {
      y := initialFGGM(base_program);
    } else {
      if (prev_diff_err ≠ "") {
        y := diffErrorFGGM(base_program, prev_diff_err, prev_code);
      } else {
        y := verifierErrorFGGM(base_program, prev_ver_err, prev_code);
      }
    }

    prev_diff_err := "";
    prev_ver_err := "";
    prev_code := y;

    var diff_err : str := noDiffWithErrorMsg(base_program, y);
    if (diff_err = "") {
      var v_err : str := dafnyVerifierWithErrorMsg(y);
      if (v_err = "") {
        r := y; return;
      } else {
        if (¬best_verified) { best := y; }
        prev_ver_err := v_err;
      }
    } else {
      prev_diff_err := diff_err;
    }
    attempt := attempt + 1;
  }
  r := best; return;
}

```

Fig. 7. Verified agent program for Dafny annotation synthesis. On the first iteration, *initialFGGM* (Fig. 4) generates an initial annotation. On subsequent iterations, the agent dispatches to *diffErrorFGGM* (Fig. 5) or *verifierErrorFGGM* (Fig. 6) depending on the previous error type. Each FGGM guarantees *noDiff*(*base_program*, *y*), maintaining the loop invariant. Upon exhausting all attempts, the agent returns *best*, which satisfies the postcondition by the loop invariant.

L FGGM SYNTAX

The syntax of a first-order guarded generative model (FGGM) can be described using the following BNF grammar. We use the production rules ($\langle \text{id} \rangle$, $\langle \text{type} \rangle$, $\langle \text{spec} \rangle$, $\langle \text{program} \rangle$, $\langle \text{string_lit} \rangle$, $\langle \text{formula} \rangle$) from Appendix A.

$$\begin{aligned}
 \langle id \rangle &::= \langle \text{id} \rangle \\
 \langle GMid \rangle &::= \langle \text{id} \rangle \\
 \langle typeSig \rangle &::= "(" \langle \text{typedVars} \rangle ")" \rightarrow \langle \text{type} \rangle \\
 \langle \text{typedVars} \rangle &::= \langle \text{typedVar} \rangle \mid \langle \text{typedVar} \rangle, \langle \text{typedVars} \rangle \\
 \langle \text{typedVar} \rangle &::= \langle \text{id} \rangle : \langle \text{type} \rangle \\
 \langle \Phi_l \rangle &::= "requires" \langle \text{formula} \rangle \\
 \langle \Psi_l \rangle &::= "ensures" \langle \text{formula} \rangle \\
 \langle f_p \rangle &::= \langle \text{program} \rangle \\
 \langle f_d \rangle &::= \langle \text{program} \rangle \\
 \langle \text{info} \rangle &::= \langle \text{string_lit} \rangle
 \end{aligned}$$

M FGGM VALIDITY CHECK

Algorithm 3 checks whether an FGGM definition is valid. It verifies that the type signature and local contracts are syntactically well-formed and that all terms using library functions satisfy their input specifications. It then checks that the prompting program and fallback program are type-correct and terminating. Finally, the fallback program is verified using the deductive verifier to ensure it satisfies the local contract for all inputs. The FGGM definition is accepted only if no errors are found.

Algorithm 3 validateFGGM

```

1: Input: FGGM definition  $G = (id, GMid, typeSig, \Phi_I, \Psi_I, f_p, f_d)$ 
2: Input: Library functions  $\mathcal{F}_c$ , axioms  $\mathcal{A}$ 
3: Output:  $\{\}$  if valid else set of errors  $err$ 
4:  $err \leftarrow \{\}$ 
5: if  $\neg WellTyped(typeSig)$  then
6:    $err \leftarrow err \cup \{\text{Invalid type signature}\}$ 
7: end if
8: if  $\neg WellFormedFormula(\Phi_I)$  then
9:    $err \leftarrow err \cup \{\text{Invalid input contract}\}$ 
10: end if
11: if  $\neg WellFormedFormula(\Psi_I)$  then
12:    $err \leftarrow err \cup \{\text{Invalid output contract}\}$ 
13: end if
14: if  $\neg CheckTerms(\Phi_I, \mathcal{F}_c)$  then
15:    $err \leftarrow err \cup \{\text{Invalid library terms in } \Phi_I\}$ 
16: end if
17: if  $\neg CheckTerms(\Psi_I, \mathcal{F}_c)$  then
18:    $err \leftarrow err \cup \{\text{Invalid library terms in } \Psi_I\}$ 
19: end if
20: if  $\neg TypeCheck(f_p)$  or  $\neg Terminates(f_p)$  then
21:    $err \leftarrow err \cup \{\text{Invalid prompting program}\}$ 
22: end if
23: if  $\neg TypeCheck(f_d)$  or  $\neg Terminates(f_d)$  then
24:    $err \leftarrow err \cup \{\text{Invalid fallback program}\}$ 
25: end if
26: if  $\neg Verify(\forall x_1 \dots x_n, y. \Phi_I(x_1 \dots x_n) \Rightarrow \Psi_I(x_1 \dots x_n, f_d(x_1 \dots x_n, y)))$  then
27:    $err \leftarrow err \cup \{\text{Fallback violates contract}\}$ 
28: end if
29: return  $err$ 

```

Algorithm 4 checks whether a concrete input–output tuple satisfies the FGGM output contract. It first substitutes the concrete values into the specification Ψ_I . If the resulting formula is quantifier-free, the checker directly evaluates it by computing all terms using the library functions. If the specification contains quantifiers, the substituted formula is combined with the axioms and input specification and submitted to an SMT solver with a timeout. The procedure returns true only if the solver confirms satisfiability within the time bound.

Algorithm 4 $check_{\mathcal{A}, \Phi_I, \Psi_I}$ Contract Checker

```

1: Input: Concrete values  $(x_1, \dots, x_n, y)$ 
2: Input: Local contracts  $(\Phi_I, \Psi_I)$ , axioms  $\mathcal{A}$ 
3: Output:  $T$  if contract satisfied else  $F$ 
4:  $\phi \leftarrow \text{Substitute}(\Psi_I, \{x_1, \dots, x_n, y\})$ 
5: if  $\text{QuantifierFree}(\phi)$  then
6:   if  $\phi = \neg\phi'$  then return  $\neg check_{\mathcal{A}, \Phi, \phi'}(x_1, \dots, x_n, y)$ 
7:   else if  $\phi = \phi' \vee \psi$  then return  $check_{\mathcal{A}, \Phi, \phi'}(x_1, \dots, x_n, y) \vee check_{\mathcal{A}, \Phi, \psi}(x_1, \dots, x_n, y)$   $\triangleright$  Recursively evaluate
   on subformulas
8:   else if  $\phi = \phi' \wedge \psi$  then return  $check_{\mathcal{A}, \Phi, \phi'}(x_1, \dots, x_n, y) \wedge check_{\mathcal{A}, \Phi, \psi}(x_1, \dots, x_n, y)$ 
9:   else if  $\text{isAtomic}(\phi)$  then return  $\phi(x_1, \dots, x_n, y)$   $\triangleright$  Compute the atomic predicate  $\phi$  value on  $(x_1, \dots, x_n, y)$   $T$ 
   or  $F$ 
10:  end if
11: end if
12:  $\psi \leftarrow \varphi_{\mathcal{A}} \wedge \Phi_I(x_1, \dots, x_n) \wedge \phi$ 
13:  $res \leftarrow \text{SMTSolve}(\psi, time)$ 
14: if  $res = SAT$  then
15:   return  $T$ 
16: else
17:   return  $F$ 
18: end if

```

N BACKGROUND ON GRPO

Group Relative Policy Optimization (GRPO): GRPO is a reinforcement learning–based fine-tuning method for generative models that optimizes a policy by comparing outputs *relative to other samples from the same input*, rather than relying on an explicit value function. Let $\pi_{\theta}(y | p)$ denote the parametric policy for an FGGM id_{θ_i} , where $p \in \mathbb{P}$ is the input prompt and $y_l \sim \pi_{\theta}(\cdot | p)$ is a sampled output. For each input p , GRPO samples a *group* of G candidate outputs $\{y^{(1)}, \dots, y^{(G)}\}$ from the current policy. Each sample is assigned a scalar reward $r^{(j)} = \mathcal{R}(p, y^{(j)})$. Instead of using absolute rewards directly, GRPO computes *relative advantages* within the group:

$$A^{(j)} = \frac{r^{(j)} - \mu_p}{\sigma_p}, \quad \mu_p = \frac{1}{G} \sum_{j=1}^G r^{(j)}, \quad \sigma_p = \sqrt{\frac{1}{G} \sum_{j=1}^G (r^{(j)} - \mu_p)^2 + \epsilon} \quad (24)$$

Here μ_p and σ_p are the mean and standard deviation of rewards within the group.

$$\max_{\theta} \mathbb{E}_{p \sim \mathbb{P}} \left[\frac{1}{G} \sum_{j=1}^G \min \left(\rho^{(j)} A^{(j)}, \text{clip}(\rho^{(j)}, 1 - \epsilon, 1 + \epsilon) A^{(j)} \right) \right] \quad (25)$$

where $\rho^{(j)} = \frac{\pi_{\theta}(y^{(j)} | p)}{\pi_{\theta_{\text{old}}}(y^{(j)} | p)}$ is the importance ratio. This clipped objective, similar to PPO, ensures stable policy updates. In our setting, each FGGM id_{θ_i} defines a policy $\pi_{\theta_i}(y_l | p)$ over outputs y_l for input p . The reward function defined in Eq. 23:

$$\mathcal{R}(p, y_l) = 1 - \text{Sigmoid}(L(x_{_, y}) \times \mathbb{I}(y = y_l) + \lambda \times (1 - \mathbb{I}(check_{\mathcal{A}, \Phi_I, \Psi_I}(p, y_l)))) \quad (26)$$

O HYPERPARAMETERS

Table 8 lists all hyperparameters used across the four evaluation tasks. The top block reports the shared pipeline settings applied uniformly across tasks; the middle blocks report task-specific training configurations for GSM-Symbolic and Symbolic Regression; the bottom block reports the per-task verifier and solver timeouts.

Table 8. Hyperparameter settings for all evaluation tasks.

Category	Hyperparameter	Value
Shared pipeline	Planner search budget (Δ)	10
	FGGM rejection-sample budget (K)	5
	In-context examples	3 (2 for τ^2 -bench)
GRPO / LoRA (GSM-Symbolic)	Base model	Qwen3-8B
	LoRA rank	16
	LoRA alpha (α)	32
	GRPO training epochs	5
	GRPO batch size	64
	GRPO learning rate	1×10^{-5}
	GRPO generations per prompt (G)	8
Parameter tuning (Symbolic Regression)	Optimizer	Adam
	Learning rate	0.05
	Backpropagation steps	40
Verifier timeouts	Dafny — <code>dafny verify</code> timeout	240 s per subprocess call
	τ^2 -bench — <code>Z3 solver.check()</code> timeout	10 s per call (up to 2 retries)
	GSM-Symbolic — <code>Z3 solver.check()</code> timeout	5 s per call

P PLANNER FEEDBACK

After each search–verify–learn iteration, SEVerA constructs structured feedback I' from the execution traces of the current best candidate agent on the training data D (Algorithm 2, line 31). The planner LLM receives the task-performance score of the previous program along with per-example error diagnostics. Specifically, for each training example on which the previous candidate failed, SEVerA provides the planner with the example input, the agent’s output, and the corresponding error message (e.g., a verification error, a grammar violation, or an incorrect answer). The planner is then asked to generate a one-to-two sentence description and a suggested fix for each failed example. The resulting list of per-example descriptions and suggestions is collected and fed back to the planner as part of the context I' when it samples the next candidate program. This feedback mechanism enables the planner to identify recurring failure patterns and adjust its synthesized FGGM definitions and prompting programs f_p accordingly in subsequent iterations.

Q PLANNER LLM PROMPT

The following is the prompt template used by the planner LLM to synthesize candidate agent programs. Placeholders in braces (e.g., `{task_description}`) are populated at each iteration with the task specification, postcondition, available operator library, agent signature, and feedback from the previous iteration’s execution traces (Appendix P).

```

You are an expert programmer tasked with writing a Dafny agent function.

## TASK
{task_description}

## CRITICAL POSTCONDITION
{postcondition_description}

## AVAILABLE OPERATORS
You can use these operators in your agent:

```dafny
{library_functions}
```

## AGENT SIGNATURE AND POSTCONDITION
```dafny
{agent_signature}
```

## CONSTRAINTS
- Fully Typed: All variables must have type annotations
- No imports: All operators are already available in scope
- No recursion: Don't call agent() recursively
- Always satisfy postcondition: Every return path must satisfy the postcondition
- No f-strings: Do not use format strings (f-strings); use the addition operator for concatenation
{task_specific_constraints}

## FEW-SHOT EXAMPLES
{few_shot_examples}

## PREVIOUS AGENT CODE
{previous_agent_code}

## PREVIOUS ATTEMPT FEEDBACK
{feedback}

## SCORING INFORMATION
{scoring_info}

## OUTPUT
Write ONLY the agent function inside a ```dafny``` block.

Make sure to:
- Check the postcondition before every return
- If the task fails, return a safe fallback value
- LLM calls are stateless, calling the LLM again does not maintain previous context
- Predefine variables before using them to avoid syntax errors
{task_specific_output_tips}

```

Fig. 8. Planner LLM prompt template. The planner receives this prompt at each search–verify–learn iteration, with placeholders filled from the task specification and feedback from the previous iteration.

R PROOF DETAILS

LEMMA 5.1 (COMPLETENESS OF CHECKER). *For any valid FGGM G with quantifier-free Ψ_l , $\forall x_1 \in T_1, \dots, \forall x_n \in T_n, \forall y \in T_0$. $\Phi_l(x_1, \dots, x_n) \implies (\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}(x_1, \dots, x_n, y) \iff \Psi_l(x_1, \dots, x_n, y))$.*

PROOF OF LEMMA 5.1. We prove the claim by structural induction on the quantifier-free formula Ψ_l . Fix arbitrary (x_1, \dots, x_n, y) . We show:

$$\text{check}_{\mathcal{A}, \Phi_l, \Psi_l}(x_1, \dots, x_n, y) \iff \Psi_l(x_1, \dots, x_n, y).$$

Base case. Suppose Ψ_l is an atomic predicate $P(x_1, \dots, x_n, y)$ over terms constructed from computable library functions in \mathcal{F}_c . Under the substitution (x_1, \dots, x_n, y) , $\text{check}_{\mathcal{A}, \Phi, P}$ evaluates each term

by executing the corresponding functions in \mathcal{F}_c . Since these functions are computable $check_{\mathcal{A},\Phi,P}$ exactly computes the atomic predicate $P(x_1, \dots, x_n, y)$ on (x_1, \dots, x_n, y) which reduces to T or F .

Inductive step. Assume the claim holds for formulas ϕ and ψ i.e. $check_{\mathcal{A},\Phi,\phi}(x_1, \dots, x_n, y) \iff \phi(x_1, \dots, x_n, y)$ and $check_{\mathcal{A},\Phi,\psi}(x_1, \dots, x_n, y) \iff \psi(x_1, \dots, x_n, y)$ We show it holds for compound formulas:

- If $\Psi_l = \neg\phi$, then by definition of $check_{\mathcal{A},\Phi_l,\Psi_l}$,

$$check_{\mathcal{A},\Phi_l,\Psi_l}(x_1, \dots, x_n, y) = \neg check_{\mathcal{A},\Phi,\phi}(x_1, \dots, x_n, y).$$

By the inductive hypothesis,

$$check_{\mathcal{A},\Phi_l,\Psi_l}(x_1, \dots, x_n, y) \iff \neg(check_{\mathcal{A},\Phi,\phi}(x_1, \dots, x_n, y)) \iff \neg(\phi(x_1, \dots, x_n, y)) \iff \Psi_l(x_1, \dots, x_n, y)$$

- If $\Psi_l = \phi \wedge \psi$, then

$$check_{\mathcal{A},\Phi_l,\Psi_l}(x_1, \dots, x_n, y) = check_{\mathcal{A},\Phi,\phi}(x_1, \dots, x_n, y) \wedge check_{\mathcal{A},\Phi,\psi}(x_1, \dots, x_n, y),$$

which, by the inductive hypothesis,

$$check_{\mathcal{A},\Phi_l,\Psi_l}(x_1, \dots, x_n, y) \iff (check_{\mathcal{A},\Phi,\phi}(x_1, \dots, x_n, y) \wedge check_{\mathcal{A},\Phi,\psi}(x_1, \dots, x_n, y)) \iff (\phi(x_1, \dots, x_n, y) \wedge \psi(x_1, \dots, x_n, y)) \iff \Psi_l(x_1, \dots, x_n, y)$$

- If $\Psi_l = \phi \vee \psi$, the argument is analogous to $\Psi_l = \phi \wedge \psi$.

□

LEMMA R.1 (CHECKER SOUNDNESS). $\varphi_{\mathcal{A}} \implies (\forall x_1 \in T_1, \dots, \forall x_n \in T_n, \forall y \in T_o. \Phi_l(x_1, \dots, x_n) \implies (check_{\mathcal{A},\Phi_l,\Psi_l}(x_1, \dots, x_n, y) \implies \Psi_l(x_1, \dots, x_n, y)))$.

PROOF. For quantifier-free Ψ_l soundness follows from Lemma 5.1. For Ψ_l with quantifier $check_{\mathcal{A},\Phi_l,\Psi_l}$ returns true if $(\varphi_{\mathcal{A}} \wedge \Phi_l(x_1, \dots, x_n)) \wedge \Psi_l(x_1, \dots, x_n, y)$ is evaluated to true within timeout. Hence for any (x_1, \dots, x_n, y)

$$\begin{aligned} check_{\mathcal{A},\Phi_l,\Psi_l}(x_1, \dots, x_n, y) &\implies (\varphi_{\mathcal{A}} \wedge \Phi_l(x_1, \dots, x_n)) \wedge \Psi_l(x_1, \dots, x_n, y) \\ (\varphi_{\mathcal{A}} \wedge check_{\mathcal{A},\Phi_l,\Psi_l}(x_1, \dots, x_n, y)) &\implies \Phi_l(x_1, \dots, x_n) \wedge \Psi_l(x_1, \dots, x_n, y) \\ (\varphi_{\mathcal{A}} \wedge check_{\mathcal{A},\Phi_l,\Psi_l}(x_1, \dots, x_n, y) \wedge \Phi_l(x_1, \dots, x_n)) &\implies \Psi_l(x_1, \dots, x_n, y) \\ (\varphi_{\mathcal{A}} \wedge \Phi_l(x_1, \dots, x_n)) &\implies (check_{\mathcal{A},\Phi_l,\Psi_l}(x_1, \dots, x_n, y) \implies \Psi_l(x_1, \dots, x_n, y)) \\ \varphi_{\mathcal{A}} \implies (\Phi_l(x_1, \dots, x_n)) &\implies ((check_{\mathcal{A},\Phi_l,\Psi_l}(x_1, \dots, x_n, y) \implies \Psi_l(x_1, \dots, x_n, y))) \end{aligned} \quad (27)$$

From Eq 27 it follows, $\varphi_{\mathcal{A}} \implies (\forall x_1 \in T_1, \dots, \forall x_n \in T_n, \forall y \in T_o. \Phi_l(x_1, \dots, x_n) \implies (check_{\mathcal{A},\Phi_l,\Psi_l}(x_1, \dots, x_n, y) \implies \Psi_l(x_1, \dots, x_n, y)))$. □

THEOREM 5.2 (VALID LOCAL CONTRACT). *For any valid FGGM G with (Φ_l, Ψ_l) and parametric GM f_G^Θ , $\varphi_{\mathcal{A}} \implies (\forall \theta \in \Theta, \forall x_1 \in T_1, \dots, \forall x_n \in T_n. \Phi_l(x_1, \dots, x_n) \implies \Psi_l(x_1, \dots, x_n, id_G(x_1, \dots, x_n)))$.*

PROOF. Let $r = id_G(x_1, \dots, x_n)$. Based on the rejection sampler with the checker $check_{\mathcal{A},\Phi_l,\Psi_l}$ (Fig 2). r is returned from the if branch if $check_{\mathcal{A},\Phi_l,\Psi_l}(x_1, \dots, x_n, y) \{ \text{return } y; \}$ or from the fallback $r = f_d(x_1, \dots, x_n, y)$. Using the semantics of an if block, the following condition always holds for the output r where y_f is the final rejected sample.

$$\varphi_{\mathcal{A}} \implies (\forall \theta \in \Theta. (\forall x_i \in T_i). \Phi_l(x_1, \dots, x_n) \implies (check_{\mathcal{A},\Phi_l,\Psi_l}(x_1, \dots, x_n, r) \vee (r = f_d(x_1, \dots, x_n, y_f)))) \quad (28)$$

Eq. 29 follows from the validity of the fallback f_d , and Eq. 18 follows from the soundness of the checker $\text{check}_{\mathcal{A}, \Phi, \Psi}$ (lemma R.1). To simplify the notation, we used $(\forall x_i \in T_i)$ to denote $(\forall x_1 \in T_1, \dots, \forall x_n \in T_n)$.

$$\varphi_{\mathcal{A}} \implies \forall \theta \in \Theta. (\forall x_i \in T_i). \Phi_I(x_1, \dots, x_n) \wedge (r = f_d(x_1, \dots, x_n, y_f)) \implies \Psi_I(x_1, \dots, x_n, r) \quad (29)$$

$$\varphi_{\mathcal{A}} \implies (\forall \theta \in \Theta. (\forall x_i \in T_i). (\Phi_I(x_1, \dots, x_n) \wedge \text{check}_{\mathcal{A}, \Phi, \Psi}(x_1, \dots, x_n, r)) \implies \Psi_I(x_1, \dots, x_n, r)) \quad (30)$$

$$\varphi_{\mathcal{A}} \implies (\forall \theta \in \Theta. (\forall x_i \in T_i). \Phi_I(x_1, \dots, x_n) \implies \Psi_I(x_1, \dots, x_n, r)) \quad \text{Using Eq. (28, 29, 30)}$$

□

THEOREM 5.3. *If $P_{\{\Theta\}} \neq \perp$ with FGGM set \mathcal{G} , then $\forall \theta_1 \in \Theta_1, \dots, \forall \theta_k \in \Theta_k. \varphi_{\mathcal{A}} \implies (\forall x_1 \in T_1, \dots, \forall x_n \in T_n. \Phi(x_1, \dots, x_n) \implies \Psi(x_1, \dots, x_n, P_{\{(\theta_1, \dots, \theta_k)/\Theta\}}(x_1, \dots, x_n)))$.*

PROOF. Let $\mathcal{G} = \{G_1, \dots, G_k\}$ and $\mathcal{A}_{\mathcal{G}}$ denote the set of first-order sentences defining all the local contracts as shown below.

$\mathcal{A}_{\mathcal{G}} = \{\varphi_{G_1}, \dots, \varphi_{G_k}\}$ where φ_{G_i} defined below

$$\begin{aligned} \varphi_{G_i} &: \forall x_1 \in T_1, \dots, x_n \in T_n. G_i^{\Psi_I}(x_1, \dots, x_n) \implies G_i^{\Psi_I}(x_1, \dots, x_n, G_i(x_1, \dots, x_n)) \\ (\varphi_{\mathcal{A}} \implies (\forall \theta_1 \in \Theta_1. \varphi_{G_1})) &\wedge \dots \wedge (\varphi_{\mathcal{A}} \implies (\forall \theta_k \in \Theta_k. \varphi_{G_k})) \quad \text{from Theorem 5.2} \\ (\varphi_{\mathcal{A}} \implies \bigwedge_{i=1}^k (\forall \theta_i \in \Theta_i. \varphi_{G_i})) &\quad (31) \end{aligned}$$

If $P_{\{\Theta\}} \neq \perp$ then $\mathcal{V}_C[\Phi, \Psi]$ with $C = (G_D, \mathcal{F} \cup \mathcal{G}, \mathcal{A} \cup \mathcal{A}_{\mathcal{G}})$ enforces

$$\begin{aligned} (P_{\{\Theta\}} \neq \perp) &\implies ((\varphi_{\mathcal{A}} \wedge (\bigwedge_{i=1}^k \forall \theta_i \in \Theta_i. \varphi_{G_i})) \implies \\ \forall x_1 \in T_1 \dots \forall x_n \in T_n. \Phi(x_1, \dots, x_n) &\implies \Psi(x_1, \dots, x_n, P_{\{(\theta_1, \dots, \theta_k)/\Theta\}}(x_1, \dots, x_n)) \quad (32) \\ (P_{\{\Theta\}} \neq \perp) &\implies ((\varphi_{\mathcal{A}} \implies \\ \forall \theta_1 \in \Theta_1, \dots, \forall \theta_k \in \Theta_k, \forall x_1 \in T_1 \dots \forall x_n \in T_n. & \\ \Phi(x_1, \dots, x_n) \implies \Psi(x_1, \dots, x_n, P_{\{(\theta_1, \dots, \theta_k)/\Theta\}}(x_1, \dots, x_n)) &\text{ Using Eq. 31 and Eq. 32} \end{aligned}$$

□

THEOREM 5.4 (SOUNDNESS). *If $f^* \neq \perp$ then, $\varphi_{\mathcal{A}} \implies \forall x_1 \in T_1, \dots, \forall x_n \in T_n. \Phi(x_1, \dots, x_n) \implies \Psi(x_1, \dots, x_n, f^*(x_1, \dots, x_n))$.*

PROOF. $f^* \leftarrow \arg \min_{f \in \mathcal{P}} \sum_{(x,y) \in D} L(x, y, f(x))$. If $f^* \neq \perp$, then by definition $(\mathcal{P} \neq \{\}) \wedge (f^* \in \mathcal{P})$. By Theorem 5.3, every program $f \in \mathcal{P}$ that satisfies $\varphi_{\mathcal{A}}$ also satisfies

$$\forall x_1 \in T_1, \dots, \forall x_n \in T_n. \Phi(x_1, \dots, x_n) \implies \Psi(x_1, \dots, x_n, f(x_1, \dots, x_n)).$$

□

THEOREM 5.5 (SUFFICIENT SUCCESS CONDITION). *If (i) L penalizes constraint violations, i.e., $\forall x \in T_i, \forall y, y' \in T_o. (\Phi(x) \wedge \neg \Psi(x, y) \wedge \Psi(x, y')) \implies L(x, _ , y') < L(x, _ , y)$, (ii) Ψ is quantifier-free, (iii) there exists a non-parametric program $f_d \in S(G, \mathcal{F}_c)$ such that f_d satisfies (Φ, Ψ) . Then, for any type-correct generative model $f_n^\theta : T_i \rightarrow T_o$ with initial parameters and any prompting program f_p , there exists a program $f \in S(G, \mathcal{F})$ such that: (1) f satisfies (Φ, Ψ) , and (2) $L(f) \leq L(f_n^\theta)$. Moreover, if there exists $(x, _) \in D$ such that $\neg \Psi(x, f_n^\theta(x))$, then $L(f) < L(f_n^\theta)$ where $L(f) = \sum_{(x, _) \in D} L(x, _ , f(x))$.*

PROOF. Let $f_n^\theta : T_i \times \dots \times T_n \rightarrow T_o$ be any type-correct generative model. By assumption (iii), there exists $f_d \in S(G, \mathcal{F}_c)$ such that

$$\varphi_{\mathcal{A}} \implies \forall x_1, \dots, x_n. \Phi(x_1, \dots, x_n) \implies \Psi(x_1, \dots, x_n, f_d(x_1, \dots, x_n)).$$

Construction. Define

$$G = (id, f_n^\theta, \tau, \Phi_I, \Psi_I, f_p, f_d, info),$$

with $(\Phi_l, \Psi_l) := (\Phi, \Psi)$ and any well-typed f_p . Since f_d satisfies the contract and Ψ is quantifier-free, G is valid and $check_{\mathcal{A}, \Phi_l, \Psi_l}$ is complete (Lemma 5.1).

Define the program $f \in S(G, \mathcal{F})$ as:

function $f(x_1 : T_1, \dots, x_n : T_n) : T_o$ {**return** $id(x_1, \dots, x_n)$; }.

Correctness. For any (x_1, \dots, x_n) such that $\Phi(x_1, \dots, x_n)$ holds, by Theorem 5.2,

$$\Psi(x_1, \dots, x_n, f(x_1, \dots, x_n)).$$

Output of f . Fix (x_1, \dots, x_n) with $\Phi(x_1, \dots, x_n)$. By completeness,

$$check_{\mathcal{A}, \Phi_l, \Psi_l}(x_1, \dots, x_n, y) \iff \Psi(x_1, \dots, x_n, y).$$

Hence,

$$f(x_1, \dots, x_n) = \begin{cases} f_n^\theta(x_1, \dots, x_n) & \text{if } \Psi(x_1, \dots, x_n, f_n^\theta(x_1, \dots, x_n)), \\ f_d(x_1, \dots, x_n) & \text{otherwise.} \end{cases}$$

Loss comparison. For any $(x, _) \in D$:

Case 1: $\Psi(x, f_n^\theta(x))$. Then $f(x) = f_n^\theta(x)$ and

$$L(x, _, f(x)) = L(x, _, f_n^\theta(x)).$$

Case 2: $\neg\Psi(x, f_n^\theta(x))$. Then $f(x) = f_d(x)$ and $\Psi(x, f(x))$ holds. By assumption (i),

$$L(x, _, f(x)) < L(x, _, f_n^\theta(x)).$$

Thus,

$$\forall (x, _) \in D. \quad L(x, _, f(x)) \leq L(x, _, f_n^\theta(x)),$$

which implies

$$L(f) \leq L(f_n^\theta).$$

Strict improvement. If $\exists (x, _) \in D$ such that $\neg\Psi(x, f_n^\theta(x))$, then Case 2 holds for at least one point, yielding

$$L(f) < L(f_n^\theta).$$

Conclusion. The constructed program f satisfies (Φ, Ψ) and achieves no greater loss than f_n^θ , with strict improvement when violations occur. Note the program verification setup satisfies this sufficient condition. \square