## PROVISIONAL PATENT APPLICATION

### UNITED STATES PATENT AND TRADEMARK OFFICE

---

### SYSTEM AND METHOD FOR ITERATIVE FORMAL VERIFICATION OF AI-GENERATED CODE USING A HALLUCINATION-VERIFICATION LOOP

---

### CROSS-REFERENCE TO RELATED APPLICATIONS

This application claims the benefit of the filing date of this provisional application under 35 U.S.C. 119(e).

---

### FIELD OF THE INVENTION

The present invention relates generally to the field of artificial intelligence and software engineering, and more particularly to systems and methods for improving the correctness of code generated by large language models (LLMs) through iterative formal verification and targeted remediation.

---

### BACKGROUND OF THE INVENTION

**The Problem of AI-Generated Code Correctness**  Large language models (LLMs) such as GPT-4, Claude, Gemini, and others have demonstrated remarkable capability in generating source code from natural language descriptions.  AI coding platforms including Cursor, Replit, Bolt.new, Lovable, Devin, GitHub Copilot, and others are integrating these models into developer workflows, enabling rapid code generation from prompts.

However, a fundamental problem remains: AI-generated code is unreliable. LLMs hallucinate — they generate plausible-looking code that contains errors, logical flaws, security vulnerabilities, and incorrect implementations.  Three independent mathematical proofs have established that hallucination is a structural inevitability of transformer-based language models:

1. Xu, Jain, and Kankanhalli (2024) proved via computational learning theory that any computable LLM must hallucinate on certain inputs, as the approximation of world models from finite training data is inherently lossy.

2. Banerjee, Sarkar, and Schwaller (2024) invoked connections to Godel's incompleteness theorems to demonstrate that any sufficiently expressive formal system (which LLMs approximate) cannot be both complete and consistent with respect to all factual claims.

3. Karpowicz (2025) established a quadrilemma proving that no LLM can simultaneously achieve truthful knowledge representation, semantic information conservation, complete revelation, and knowledge-constrained optimality. The proof demonstrates that the Log-Sum-Exp operation at the core of softmax attention makes hallucination structurally inevitable.

These results prove that hallucination cannot be eliminated from LLMs through training improvements, reinforcement learning from human feedback (RLHF), Constitutional AI, guardrails, or any other technique that modifies the generative model itself.

**Failure of Existing Approaches**   Current approaches to improving AI code correctness fall into several categories, all of which have been demonstrated to be inadequate:

**Self-Refinement.** The model is asked to review and improve its own output. Empirical evidence demonstrates that self-refinement provides essentially no improvement over single-pass generation. On the HumanEval benchmark (164 code generation tasks), self-refinement achieves only +1.2 percentage points improvement over 5 iterations (87.8% vs. 86.6% baseline). This is because the self-critic shares the generator's knowledge gaps and failure modes — the model cannot identify errors in code that it was already incapable of writing correctly.

**LLM-as-Judge.** A separate LLM instance evaluates the generated code for correctness. While initially effective (improving accuracy from 86.6% to 98.2% at k=1), LLM-as-Judge verification *regresses* at higher iteration counts. On HumanEval, accuracy drops from 99.4% at k=3 iterations to 97.0% at k=5 iterations. This regression occurs because the LLM judge produces false-positive "failure" verdicts on correct code, triggering unnecessary remediation that introduces bugs. The noise in the learned verifier accumulates over iterations, causing the system to actively degrade correct solutions.

**Random Verification.** Assigning random pass/fail verdicts to generated code causes accuracy to *decrease* from 97.6% at k=1 to 95.1% at k=3 as random "failure" signals cause the model to corrupt correct solutions.

**Guardrails and Filters.** Post-hoc detection and suppression of hallucinated content. These approaches reduce but cannot eliminate hallucination, and they constrain the generative process, reducing the model's creative and reasoning capabilities.

**Reinforcement Learning from Human Feedback (RLHF).** Training the model with human preference signals. Subject to reward hacking (model optimizes proxy objectives rather than true correctness), inter-annotator disagreement (20-30% disagreement rates), and the fundamental impossibility of eliminating hallucination through training alone.

There is therefore a need for a system and method that can reliably improve the correctness of AI-generated code without requiring modification of the underlying generative model, and that converges monotonically rather than regressing at higher iteration counts.

---

## SUMMARY OF THE INVENTION

The present invention provides a system and method for iteratively improving AI-generated code through a formal verification loop referred to herein as LUCID (Leveraging Unverified Claims Into Deliverables). The

invention is based on the insight that, since hallucination cannot be eliminated from the generative process, the architectural solution is to *separate generation from verification* and use formal (non-learned) verification as the error signal driving iterative improvement.

In one aspect, the invention provides a computer-implemented method comprising:

(a) receiving, from a generative language model, source code generated in response to a task description;

(b) extracting, by a claim extraction module, one or more testable claims about the expected behavior of the generated source code;

(c) verifying, by a formal verification engine that executes the generated source code against a test specification derived from the extracted claims, each of the one or more testable claims, wherein the formal verification engine produces a verification result for each claim comprising one of: PASS, PARTIAL, FAIL, or NOT APPLICABLE;

(d) when one or more claims receive a FAIL or PARTIAL verification result, generating, by a remediation module, a structured remediation plan that identifies specific verification failures, their root causes, and proposed code modifications;

(e) regenerating, by the generative language model, updated source code incorporating the remediation plan; and

(f) repeating steps (b) through (e) for a predetermined number of iterations or until all testable claims receive a PASS verification result.

In another aspect, the invention provides a system comprising a processor and memory storing instructions that, when executed by the processor, cause the system to expose an application programming interface (API) that receives code verification requests and returns verified code, implementing the method described above.

A key technical advantage of the present invention is its **monotonic convergence property**: unlike self-refinement, LLM-as-judge, or random verification approaches — all of which plateau or regress at higher iteration counts — the formal verification loop of the present invention improves monotonically with each iteration on isolated code generation tasks. This property arises from the zero-noise characteristic of formal verification: a formal verifier (test suite execution) produces exact, noise-free error signals within its decidable domain, whereas learned verifiers (LLM-based judgment, reward models) produce noisy signals with positive correlation to the generator's own errors.

---

**DETAILED DESCRIPTION OF THE INVENTION**

**1. System Architecture Overview**    The LUCID system comprises four principal components that operate in an iterative loop:

**1.1 Generator Module.**    Any large language model (LLM) capable of generating source code from task descriptions. The system is model-agnostic — it functions as a meta-architecture composable with any generative model including but not limited to GPT-4, Claude, Gemini, Llama, Mistral, DeepSeek, and any future models. The generator module receives a task description (e.g., a function signature with docstring, a GitHub issue description, a natural language specification) and produces candidate source code. The generator operates at standard or elevated temperature settings to maximize the diversity and creativity of generated solutions.

**1.2 Claim Extraction Module.** A processing stage that decomposes the generator's output into individual, testable claims about the code's expected behavior. Each claim is a decidable predicate over the space of codebases — a statement that can be verified as true or false by executing the code. Claims are categorized by type (functionality, security, data handling, performance, compliance) and assigned a severity level (critical, high, medium, low). The extraction module may use an LLM or rule-based parsing to decompose the generated output.

In certain embodiments, the claim extraction module extracts claims from: - Function signatures and docstrings (expected input/output behavior) - Inline comments and assertions - Error handling patterns (expected exception behavior) - Data validation patterns (expected type and format constraints) - Hallucinated specifications (in the specification-level variant described in Section 3)

**1.3 Formal Verification Engine.** The core innovation of the system. Unlike learned verifiers (LLM-as-judge, reward models, discriminators), the formal verification engine uses *execution-based ground truth* to verify claims. For each testable claim, the engine:

(a) Generates or retrieves a test specification corresponding to the claim. In the code generation variant, the test specification is the ground-truth test suite associated with the task. In the specification-level variant, the verification engine generates executable tests from the extracted claims.

(b) Executes the generated source code against the test specification in a sandboxed environment. The sandbox provides isolation from the host system, preventing generated code from accessing the network, filesystem, or other system resources.

(c) Captures the execution result, including:

- Pass/fail status for each test case
- Error messages and stack traces for failing tests
- Standard output and standard error streams
- Execution time and resource utilization

(d) Returns a structured verification result comprising: a verdict (PASS, PARTIAL, FAIL, or NOT APPLICABLE) for each claim; the specific error messages and stack traces for failing claims; and a quantitative specification gap metric equal to one minus the fraction of passing claims.

The formal verification engine has a critical property: **zero verification noise** within its decidable domain. A formal verifier (a decision procedure for a decidable property class) returns the correct verdict for all claims in its domain, by definition of decidability. This means the prediction error signal driving the iterative loop is exact — there are no false positives (marking correct code as incorrect) or false negatives (marking incorrect code as correct) within the domain of decidable properties.

This property distinguishes the formal verification engine from all learned verification approaches: - Learned discriminators have non-zero approximation error - Reward models have irreducible noise from inter-annotator disagreement (20-30%) - Self-critique shares systematic biases with the generator

**1.4 Remediation Module.** When verification failures are detected, the remediation module generates a structured remediation plan. The remediation module receives: - The current source code - The verification results (which claims failed) - The specific error messages and stack traces for each failure - The extracted claims for context

The remediation module (which may use an LLM) produces: - A diagnosis of each failure's root cause - A prioritized list of code modifications to address the failures - Preservation directives to maintain code that already passes verification

The remediation module feeds its output back to the generator module, which produces updated source code incorporating the remediation plan. This updated code then enters the next iteration of the verification loop.

**2. The Iterative Verification Loop**    The core method of the invention operates as follows:

```
PROCEDURE VerifyAndImprove(task_description, max_iterations):

    Step 1: GENERATE
        code_0 = Generator.generate(task_description)

    Step 2: EXTRACT
        claims = ClaimExtractor.extract(code_0, task_description)

    Step 3: VERIFY
        results = FormalVerifier.verify(code_0, claims)

    Step 4: CHECK CONVERGENCE
        IF results.all_passed OR iteration >= max_iterations:
            RETURN code_current

    Step 5: REMEDIATE
        plan = Remediator.generate_plan(code_current, results, claims)

    Step 6: REGENERATE
        code_next = Generator.regenerate(code_current, plan)

    Step 7: ITERATE
        SET code_current = code_next
        GO TO Step 2
```

**Early Exit Optimization.** The loop includes an early exit condition: if all testable claims pass verification at any iteration, the loop terminates immediately. This means that tasks solved at k=1 do not consume additional API tokens at higher iteration counts, making the system more efficient than approaches that always run for a fixed number of iterations.

**Context Preservation.** At each iteration, the regeneration step receives the full context of prior iterations: the original task description, the current code, the verification results, and the remediation plan. This enables the generator to make targeted fixes rather than regenerating from scratch, improving convergence speed.

**3. Specification-Level Variant (Specification Gap Minimization)**    In an alternative embodiment, the system operates at the specification level rather than the individual function level. In this variant:

**Step 1: Hallucinate.** The LLM generates an unconstrained specification for a software system. The prompt format (e.g., Terms of Service for a not-yet-existing application) is deliberately chosen to maximize the breadth of generated claims across multiple dimensions: functionality, security, data handling, performance, and legal compliance.

**Step 2: Extract.** The hallucinated output is decomposed into individual, testable claims. Each claim is categorized and assigned a severity level. The claim extraction module produces a structured specification comprising typically 50-200 individual claims.

**Step 3: Verify.** Each extracted claim is tested against the actual codebase through automated analysis (static analysis, test execution, bounded model checking, type checking, or other formal methods appropriate to the claim type). Claims receive verdicts: PASS, PARTIAL, FAIL, or NOT APPLICABLE.

**Step 4: Remediate.** Verification failures generate a structured remediation plan. The gap between hallucinated claims and verified reality is quantified as the *specification gap* — a metric defined as:

```
specification_gap = 1 - (count_of_passing_claims / count_of_applicable_claims)
```

The specification gap is a direct analog of prediction error in the predictive processing framework from computational neuroscience (Friston, 2010). The system identifies which failures to address, prioritized by severity and feasibility.

**Step 5: Regenerate.** After remediation (which involves modifying the codebase to address verification failures), the system generates updated claims. Verified truths (PASS claims) are retained. Failed claims are revised based on the new state of the codebase.

The loop repeats, and the specification gap converges monotonically: empirically, specification-reality alignment increases with each iteration (57.3% -> 69.8% -> 83.2% -> 90.8% across iterations 3-6 on a production codebase of 30,000 lines of TypeScript, 200+ files, 91 extracted claims).

**4. The Monotonic Convergence Property**   A key technical property of the present invention is that accuracy improves monotonically with iteration count when operating on isolated code generation tasks. This is established both theoretically and empirically.

**Theoretical Basis.** The monotonic convergence property arises from three conditions:

(C1) **Remediation is non-expansive:** Remediation does not increase the specification gap. That is, fixing one failure does not break another passing test (on isolated tasks where cross-dependencies are absent).

(C2) **Remediation is strictly contractive on failures:** There exists a contraction rate gamma in (0,1) such that each failing claim gets strictly closer to passing after remediation.

(C3) **Regeneration is benign:** The regenerated specification does not introduce more new failing claims than are fixed by remediation.

Under these conditions, the specification gap converges geometrically:

```
specification_gap(t) <= (gamma + beta)^t * specification_gap(0)
```

where gamma is the contraction rate and beta is the regeneration novelty bound, with gamma + beta < 1.

**Empirical Verification.** On the HumanEval benchmark (164 Python code generation tasks):

| Condition | k=1 | k=3 | k=5 | Convergence |
|---|---|---|---|---|
| Baseline (single-pass) | 86.6% | — | — | N/A |
| Self-Refine | 87.2% | 87.2% | 87.8% | Plateaus |
| LLM-as-Judge | 98.2% | 99.4% | 97.0% | **Regresses** |
| LUCID (present invention) | 98.8% | 100% | 100% | **Monotonic** |

LUCID is the only condition whose accuracy increases monotonically with iteration count, reaching 100% at k=3 and maintaining 100% at k=5. Zero regressions were observed across all iteration counts.

On the SWE-bench Lite benchmark (300 real-world software engineering tasks from GitHub):

| Condition | Resolved | Rate | vs. Baseline |
|---|---|---|---|
| Baseline k=1 | 55/300 | 18.3% | — |
| LUCID k=1 | 75/300 | 25.0% | +36.4% relative |
| LUCID best (k=1 or k=3) | 91/300 | 30.3% | +65.5% relative |

Head-to-head comparison: LUCID improved 23 tasks that baseline failed, with only 3 regressions (7.7:1 improvement-to-regression ratio).

On SWE-bench, the iterative loop is non-monotonic (14 tasks solved at k=1 regressed at k=3) because SWE-bench patches modify large, interconnected codebases where condition (C1) may be violated — fixing one test may break another. This distinction between monotonic convergence on isolated tasks and non-monotonic behavior on interconnected systems is predicted by the theoretical framework and does not diminish the core advantage: LUCID produces a net improvement at every iteration count.

**5. Precision-Weighted Verification** In a preferred embodiment, the system implements precision-weighted verification, where verification failures are weighted by their severity:

| Severity | Precision Weight |
|---|---|
| Critical | 4 |
| High | 3 |
| Medium | 2 |
| Low | 1 |

The precision-weighted specification gap is:

```
weighted_gap = sum(weight_i * (1 - verdict_i)) / sum(weight_i)
```

where the sums are over all applicable claims. This ensures that critical failures (e.g., security vulnerabilities) contribute more to the gap metric and are prioritized in remediation, while low-severity cosmetic issues have less influence.

This precision weighting is analogous to precision weighting in the predictive processing framework from computational neuroscience, where the brain assigns different confidence levels to prediction errors based on the reliability of the sensory channel.

**6. API-Based Delivery** In a preferred embodiment, the system is exposed as an application programming interface (API) that receives code verification requests and returns verified code. The API architecture comprises:

**Request Format:** - Task description (natural language or structured format) - Source code to verify (or empty for generation-from-scratch) - Maximum iteration count (default: 3) - Verification mode (code-level or specification-level) - Language and framework context

**Response Format:** - Verified source code - Verification results (per-claim verdicts) - Specification gap metric (before and after) - Number of iterations performed - Cost of verification

**Integration Patterns:** The API is designed for integration into AI coding platforms as a verification layer:

(a) **Inline verification:** The platform sends each AI-generated code output through the LUCID API before presenting it to the user. The user receives code that has been iteratively verified and improved.

(b) **Background verification:** The platform displays the initial AI output immediately and runs LUCID verification in the background. If improvements are found, the user is notified.

(c) **CI/CD integration:** The LUCID API is called as part of the continuous integration pipeline, verifying AI-generated code before it is merged into the codebase.

(d) **IDE plugin:** A development environment plugin that automatically runs LUCID verification on AI-generated code suggestions before presenting them to the developer.

**7. Sandboxed Execution Environment**    The formal verification engine executes generated code in a sandboxed environment to prevent malicious or erroneous code from affecting the host system. The sandbox provides:

- **Process isolation:** Generated code runs in a separate process with restricted permissions
- **Network isolation:** No network access from within the sandbox
- **Filesystem isolation:** Generated code can only access a temporary directory
- **Resource limits:** CPU time, memory, and disk quotas prevent resource exhaustion
- **Timeout enforcement:** Each test execution has a configurable timeout (default: 30 seconds for unit tests, 300 seconds for integration tests)

For complex verification tasks (e.g., SWE-bench), the sandbox uses containerization (e.g., Docker) to provide full operating system-level isolation, enabling verification against the exact dependency environment of the target repository.

**8. Formal Test Generation from LLM Output**    In an embodiment where ground-truth tests do not exist (i.e., the system is generating code for a novel task rather than solving a benchmark problem), the formal verification engine generates executable tests from the extracted claims. This process comprises:

(a) For each extracted claim, the claim extraction module produces a natural language description of the expected behavior.

(b) The test generation module (which may use an LLM with low temperature for determinism) translates each claim into one or more executable test cases.

(c) The generated tests are validated for syntactic correctness and executability before being used for verification.

(d) The test suite is refined across iterations: tests that consistently pass are retained, and new tests may be generated for newly extracted claims.

This test generation capability enables the system to operate in domains where pre-existing test suites do not exist, extending its applicability beyond benchmark evaluation to general code generation tasks.

**9. Ablation Evidence for Component Contributions**    The contribution of each system component has been empirically validated through ablation studies on HumanEval (164 tasks):

| Ablation | k=1 | k=3 | k=5 | Effect |
|---|---|---|---|---|
| Full system | 98.8% | 100% | 100% | Monotonic convergence |
| No extraction | 100% | 100% | 100% | Converged at k=1 (simple tasks) |
| No context | 99.4% | 99.4% | 100% | Slower convergence |
| No remediation | 99.4% | 99.4% | 99.4% | **Plateaus** (cannot reach 100%) |
| Learned verification | 98.2% | 97.6% | 98.2% | Non-monotonic |
| Random verification | 97.6% | 95.1% | 97.0% | **Diverges** |

**Key findings from ablation:**

1. The formal verification engine is the most critical component. Replacing it with random verdicts causes accuracy to *decrease* with more iterations.

2. The remediation module enables the final convergence step. Without it, accuracy plateaus at 99.4% — the system cannot solve the hardest tasks.

3. Learned (LLM-based) verification is unreliable: it produces non-monotonic behavior (dips at k=3 before recovering at k=5) due to false-positive failures.

4. Context preservation accelerates convergence but does not affect final accuracy.

**10. Model-Agnostic Architecture**   A key property of the present invention is that it is *composable with any generative model*. The system does not require a specific LLM; it functions as a meta-architecture that can sit on top of any generator:

- The hallucination stage can use GPT-4, Claude, Gemini, Llama, Mistral, DeepSeek, or any future model
- The extract and verify stages are model-agnostic
- The remediation stage produces structured output that can guide any generation process

This composability means the system is not competing with transformer architectures or other model architectures. It is a *verification layer* that enhances any generator, providing the iterative verification loop that no current model natively implements.

**11. Relationship to Predictive Processing Neuroscience**   The present invention is architecturally inspired by and maps onto the predictive processing framework from computational neuroscience (Friston, 2010; Seth, 2021; Clark, 2023). The brain operates as a prediction machine: it generates predictions (top-down), compares them to sensory input (bottom-up), computes prediction error, and updates its model. The brain does not suppress its generative (hallucinatory) process — it constrains it through sensory verification.

The LUCID system implements this same architecture:

| Brain Component | LUCID Component | Function |
|---|---|---|
| Cortical prediction | LLM generation | Generate candidates |
| Sensory comparison | Formal verification | Compare against ground truth |

| Brain Component | LUCID Component | Function |
| --- | --- | --- |
| Prediction error | Specification gap | Quantify discrepancy |
| Belief updating | Remediation | Address discrepancies |
| Updated prediction | Regeneration | Produce improved candidates |

This mapping is not merely analogical but mathematically precise: LUCID's specification gap is equivalent to Friston's prediction error, its iterative convergence corresponds to free energy minimization, and its use of formal verification as the grounding signal solves the "who verifies the verifier?" problem that undermines learned discriminators, reward models, and self-critique.

---

## ABSTRACT OF THE DISCLOSURE

A computer-implemented system and method for improving the correctness of source code generated by artificial intelligence models through an iterative formal verification loop. The system receives AI-generated source code, extracts testable claims about the code's expected behavior, formally verifies each claim by executing the code against test specifications in a sandboxed environment, generates structured remediation plans for verification failures, and regenerates improved code incorporating the remediation. The loop repeats until all claims pass or a maximum iteration count is reached. Unlike approaches using self-refinement or learned verification (LLM-as-judge), which plateau or regress at higher iteration counts, the present invention converges monotonically due to the zero-noise property of formal (execution-based) verification. Empirical evaluation demonstrates 100% accuracy on HumanEval (164 tasks) at k=3 iterations, compared to 87.8% for self-refinement and 97.0% for LLM-as-judge at k=5. On SWE-bench Lite (300 real-world tasks), the system improves resolution rate by 65.5% over baseline. The system is model-agnostic and is delivered as an API for integration into AI coding platforms.