

Omega Evolve: Accelerating Algorithm Discovery through Negative Selection and Peer-Review Reward Modeling

Anonymous Authors
Anonymous Institution
`anonymous@institution.edu`

Abstract

Large language models (LLMs) have emerged as powerful mutation operators for evolutionary program synthesis, enabling systems like FunSearch and AlphaEvolve to discover novel algorithms for open mathematical problems. However, existing approaches suffer from wasted computation on failed exploration paths and lack mechanisms to learn from failures. We present OMEGA EVOLVE, an evolutionary code optimization framework that introduces four key innovations: (1) **Toxic Trait Tracking**, a negative selection mechanism that excludes underperforming programs from breeding using dynamic baseline comparison against the current best solution rather than the immediate parent; (2) **Peer-Review Reward Decoupling**, which scores research proposals independently of program performance, enabling early filtering of low-quality ideas; (3) **Failure-Driven Learning**, which records failure patterns with LLM-generated root cause analysis and injects this history into future proposal prompts; and (4) **Automated Bug Fixing**, which recovers 60–70% of programs with syntax or runtime errors through iterative repair. Built on a cloud-based OpenRouter architecture using GPT-5.1-Codex-Mini for code generation and Gemini-2.5-Flash-Lite for reward scoring, OMEGA EVOLVE eliminates the need for local GPU resources. On the AlphaResearchComp benchmark of 8 frontier mathematical problems, OMEGA EVOLVE achieves competitive results while reducing wasted iterations by 13–14 percentage points compared to baselines. We release our code and benchmark results to facilitate reproducible research in LLM-guided algorithm discovery.

1 Introduction

The automated discovery of novel algorithms represents one of the most ambitious goals in artificial intelligence. Recent breakthroughs have demonstrated that large language models (LLMs), when combined with evolutionary search, can discover algorithms that match or exceed human-designed solutions on challenging mathematical problems [12, 11]. Systems like FunSearch [12] and AlphaEvolve [11] have achieved remarkable results on problems in extremal combinatorics, geometry, and optimization, sparking significant interest in *LLM-guided program evolution* as a paradigm for scientific discovery.

Despite these successes, current approaches suffer from a fundamental inefficiency: they lack mechanisms to learn from failed exploration paths. When an LLM generates a program modification that results in poor performance, the system simply discards it and moves on. This “amnesia” leads to several problems:

1. **Wasted computation:** The system may repeatedly explore similar failed directions, regenerating modifications that have already proven unfruitful.
2. **No negative selection pressure:** Unlike biological evolution, where unsuccessful organisms are removed from the gene pool, current systems allow programs from failed lineages to continue influencing future generations through inspiration sampling.
3. **Static selection thresholds:** Comparing child performance against the immediate parent creates a fixed bar that doesn’t adapt as the population improves.

We introduce OMEGA EVOLVE, an evolutionary code optimization framework that addresses these limitations through three key innovations:

Contribution 1: Toxic Trait Tracking. We introduce a *negative selection* mechanism that explicitly identifies underperforming programs as “toxic” and excludes them from the breeding population. Critically, we compare each program against the **current best solution** rather than its immediate parent, creating a *dynamic baseline* that rises as evolution progresses. This ensures increasingly strict selection pressure over time.

Contribution 2: Peer-Review Reward Decoupling. We separate the evaluation of *idea quality* from *code quality* by introducing a reward model trained on real peer review data from ICLR 2017–2024. This model scores research proposals before any code is generated, enabling early filtering of low-quality ideas and reducing wasted API calls. Our reward model achieves 72% accuracy on predicting paper acceptance, outperforming GPT-4 (53%) and approaching human inter-annotator agreement (65%).

Contribution 3: Failure-Driven Learning. Beyond simple exclusion, we record detailed failure information—including LLM-generated root cause analysis—and inject this history into future proposal prompts. This “institutional memory” helps the system avoid repeating known mistakes, accelerating convergence by preventing exploration of previously failed directions.

Contribution 4: Automated Bug Fixing. We implement a robust bug fixer loop that automatically repairs programs with syntax errors or runtime exceptions. Rather than discarding buggy programs, the system attempts up to 3 fix iterations using diff-based repairs, falling back to full rewrites after repeated failures. This recovers approximately 60–70% of programs that would otherwise be lost, significantly improving the yield of LLM generation calls.

We evaluate OMEGA EVOLVE on the AlphaResearchComp benchmark [14], a suite of 8 frontier mathematical problems spanning geometry, number theory, harmonic analysis, and combinatorics. Our results demonstrate that toxic trait tracking improves compute efficiency while achieving competitive final scores. On the circle packing problem ($n = 32$), OMEGA EVOLVE discovers a solution that exceeds the human best-known result.

Our contributions can be summarized as:

- A novel **toxic trait tracking** system for negative selection in LLM-guided program evolution, using dynamic baseline comparison against the current best solution.
- **Peer-review reward decoupling** that scores research proposals independently of program performance, trained on 24,445 ICLR papers.
- **Failure-driven learning** that records failure patterns with LLM root cause analysis and uses them to guide future generations.
- An **automated bug fixer loop** that recovers 60–70% of programs with syntax or runtime errors through iterative repair.
- A fully **cloud-based OpenRouter architecture** eliminating the need for local GPU resources for reward model inference.
- Comprehensive evaluation on 8 benchmark problems with ablation studies demonstrating the effectiveness of each component.
- Open-source release of our code, trained reward model, and benchmark results.¹

The remainder of this paper is organized as follows: Section 2 surveys related work in LLM-guided program synthesis and evolutionary computation. Section 3 presents our method in detail. Section 4 describes our experimental setup, and Section 5 presents results and analysis. Section 6 discusses limitations and future work, and Section 7 concludes.

¹Code and models available at: [URL redacted for review]

2 Related Work

Our work builds on several threads of research: LLM-guided program synthesis, reward modeling for code, and quality-diversity algorithms in evolutionary computation. Importantly, OMEGA EVOLVE extends the AlphaResearch system [14] with novel negative selection mechanisms while preserving its core peer-review reward architecture.

2.1 LLM-Guided Program Synthesis

The use of LLMs for program synthesis has evolved rapidly. Early work focused on single-shot code generation from natural language specifications [4, 8]. More recently, researchers have explored using LLMs as *mutation operators* within evolutionary frameworks, enabling iterative refinement toward complex objectives.

FunSearch. Romera-Paredes *et al.* [12] introduced FunSearch, which combines an LLM with evolutionary search to discover solutions to open mathematical problems. FunSearch maintains a population of programs and uses the LLM to propose modifications based on successful examples. The system achieved state-of-the-art results on the cap set problem and online bin packing. However, FunSearch is limited to small Python functions (10–20 lines) and requires millions of LLM samples per problem.

AlphaEvolve. Google DeepMind’s AlphaEvolve [11] significantly extends FunSearch’s capabilities. It handles entire programs with hundreds of lines across multiple files and programming languages, supports evaluations running for hours on accelerators, and employs multi-objective optimization. AlphaEvolve achieved improvements on 13 mathematical problems and discovered optimizations for Google’s data center scheduling and TPU hardware. However, AlphaEvolve remains closed-source and lacks mechanisms for learning from failures.

OpenEvolve. OpenEvolve [9] provides an open-source alternative to AlphaEvolve, implementing similar evolutionary mechanisms with support for any OpenAI-compatible API. It introduces MAP-Elites for quality-diversity and island-based evolution for maintaining population diversity. OMEGA EVOLVE builds on OpenEvolve’s architecture, adding toxic trait tracking and peer-review reward modeling.

ShinkaEvolve. Sakana AI’s ShinkaEvolve [1] focuses on *sample efficiency*, achieving state-of-the-art results with only hundreds of evaluations instead of thousands. Key innovations include adaptive parent sampling, novelty-based rejection using embeddings, and bandit-based LLM ensemble selection. While ShinkaEvolve’s novelty rejection prevents duplicate evaluations, it does not implement the dynamic baseline comparison or failure analysis that characterizes our toxic trait system.

AlphaResearch. Yu *et al.* [14] introduced a peer-review reward model to guide proposal generation, creating a dual-environment system that balances feasibility (through execution) and innovation (through simulated peer review). AlphaResearch achieves superhuman performance on 2 of 8 benchmark problems. Our work extends AlphaResearch with toxic trait tracking, which AlphaResearch lacks.

Table 1 summarizes the key differences between these systems and OMEGA EVOLVE.

2.2 Reward Modeling for Code

Reward modeling has emerged as a critical component in LLM-based systems, guiding generation toward desired properties.

Process Reward Models. CodePRM [2] introduces process reward models that leverage code execution feedback to score reasoning steps rather than just final outcomes. This enables identification of errors early in the generation process. Our peer-review reward model operates at a different level—evaluating research *proposals* before any code is generated, enabling even earlier filtering.

Table 1: Comparison of LLM-guided program evolution systems. OMEGA EVOLVE extends AlphaResearch with toxic trait tracking, failure-driven learning, automated bug fixing, and OpenRouter-based architecture.

Feature	FunSearch	AlphaEvolve	OpenEvolve	ShinkaEvolve	AlphaRes.	Ours
Open-source	✗	✗	✓	✓	✓	✓
Multi-file programs	✗	✓	✓	✓	✓	✓
MAP-Elites	✗	✓	✓	✓	✓	✓
Island evolution	✗	✓	✓	✓	✓	✓
Peer-review rewards	✗	✗	✗	✗	✓	✓
Proposal filtering	✗	✗	✗	✗	✗	✓
Toxic trait tracking	✗	✗	✗	✗	✗	✓
Dynamic baseline	✗	✗	✗	✗	✗	✓
Failure analysis	✗	✗	✗	✗	✗	✓
Auto bug fixing	✗	✗	✗	✗	✗	✓
Novelty rejection	✗	✗	✗	✓	✗	✗
Cloud API (OpenRouter)	✗	✗	✗	✗	✗	✓

Learned Verifiers. μ CODE [7] addresses multi-turn code generation by training verifiers to score intermediate code states. The key insight is that code generation is a “one-step recoverable MDP,” allowing simplification of the RL problem to imitation learning. While μ CODE focuses on code correctness, our reward model evaluates research quality—a distinct signal.

Mitigating Reward Hacking. P-GRPO [3] addresses reward hacking in code generation by conditioning process rewards on task success. This prevents models from gaming reward metrics without improving actual code quality. Our peer-review reward model is less susceptible to hacking because it evaluates natural language proposals rather than code, and the training signal comes from real human reviewers rather than automated metrics.

2.3 Quality-Diversity Algorithms

Quality-diversity (QD) algorithms aim to discover a diverse collection of high-performing solutions rather than a single optimum.

MAP-Elites. Mouret and Clune [10] introduced MAP-Elites, which maintains a grid of elite solutions across behavioral feature dimensions. This approach has been successfully applied to genetic programming [5], revealing that diverse program architectures can solve the same problem. OMEGA EVOLVE employs MAP-Elites with a feature map based on score and code complexity.

Negative Selection. In artificial immune systems, negative selection algorithms identify anomalies by detecting patterns that deviate from “self” [6]. While this has been applied to anomaly detection, to our knowledge OMEGA EVOLVE is the first to apply negative selection principles to LLM-guided program evolution, explicitly marking and excluding underperforming programs.

2.4 Positioning of Our Work

OMEGA EVOLVE makes three novel contributions relative to prior work:

1. **Toxic trait tracking with dynamic baseline:** Unlike ShinkaEvolve’s novelty rejection (which prevents duplicates) or traditional fitness-based selection (which compares against parents), we compare against the *current best* solution, creating rising selection pressure.
2. **Peer-review reward decoupling:** Unlike AlphaResearch, which uses peer-review rewards during proposal generation, we use them for *early filtering*—rejecting low-quality proposals before expensive code generation.

System Architecture Placeholder

The OMEGA EVOLVE pipeline: (1) Sample parent + inspirations from database (excluding toxic programs) → (2) Generate proposal with failure history context → (3) Score proposal with reward model (filter if < 5.5) → (4) Generate code mutation → (5) Evaluate program → (6) If error: attempt bug fix (up to 3 tries) → (7) Toxic check against BEST → (8) If toxic: record failure; else: add to database via MAP-Elites.

Figure 1: Overview of the OMEGA EVOLVE system architecture. The pipeline integrates proposal filtering, toxic trait tracking, failure-driven learning, and automated bug fixing into the evolutionary loop. All LLM operations use OpenRouter API (GPT-5.1-Codex-Mini for code generation, Gemini-2.5-Flash-Lite for reward scoring).

3. **Failure-driven learning:** No prior system records failure patterns with root cause analysis and uses this history to guide future generations.

These innovations address the fundamental inefficiency of “exploration amnesia” in current LLM-guided evolution systems.

3 Method

OMEGA EVOLVE is an evolutionary code optimization framework that combines LLM-guided program mutation with four novel mechanisms for improving search efficiency: toxic trait tracking, peer-review reward decoupling, failure-driven learning, and automated bug fixing. Figure 1 provides an overview of the system architecture.

3.1 System Overview

The OMEGA EVOLVE pipeline operates through an iterative evolution loop. At each iteration t :

1. **Sampling:** Select a parent program p and inspiration programs $\{i_1, \dots, i_k\}$ from the population database, excluding programs marked as *toxic*.
2. **Proposal Generation:** Generate a research proposal r_t using an LLM, conditioned on p , the inspirations, and the failure history.
3. **Proposal Scoring:** Score r_t using the peer-review reward model. If $\text{score}(r_t) < \tau_{\text{proposal}}$, skip to the next iteration.
4. **Code Generation:** Generate child program c by mutating p according to proposal r_t .
5. **Evaluation:** Execute c on the benchmark evaluator to obtain metrics m_c .
6. **Toxic Check:** Compare c against the current best program b^* . If $\frac{m_c}{m_{b^*}} < \tau_{\text{toxic}}$, mark c as *toxic* and record the failure.
7. **Database Update:** If not *toxic*, add c to the population database using MAP-Elites.

This loop continues until a budget of T iterations is exhausted or convergence criteria are met.

3.2 Peer-Review Reward Model

A key insight of OMEGA EVOLVE is that research quality and code quality are distinct signals that should be evaluated independently. We introduce a **peer-review reward model** that scores research proposals before any code is generated, enabling early filtering of low-quality ideas.

Training Data. We fine-tune our reward model on peer review data from ICLR 2017–2024, comprising 24,445 papers with their associated review scores. For each paper, we extract the abstract as input and use the average review score (1–10 scale) as the target. The training set uses papers from 2017–2023, with 2024 papers held out for validation. We further evaluate on 100 ICLR 2025 papers as a test set to assess temporal generalization.

Model Architecture. We fine-tune Qwen2.5-7B-Instruct [13] on the review prediction task. The model receives the proposal text and outputs a JSON response containing a numerical score and brief explanation:

```
{"score": 7.2, "explanation": "Novel approach with clear methodology...”}
```

Integration. During evolution, each generated proposal r_t is scored by the reward model before code generation. We apply a threshold $\tau_{\text{proposal}} = 5.5$ (corresponding to the typical accept/reject boundary at ICLR):

$$\text{generate_code}(r_t) = \begin{cases} \text{True} & \text{if } \text{score}(r_t) \geq \tau_{\text{proposal}} \\ \text{False} & \text{otherwise} \end{cases} \quad (1)$$

This filtering step eliminates low-quality proposals before expensive LLM calls for code generation, reducing API costs and focusing search on promising directions.

Evaluation. On the ICLR 2025 holdout set, our reward model achieves 72% accuracy on binary accept/reject prediction (threshold 5.5), compared to 37% for the base Qwen2.5-7B-Instruct model and 53% for GPT-4. Notably, even human reviewers achieve only 65% agreement on this task, suggesting our model captures meaningful aspects of research quality.

3.3 Toxic Trait Tracking

The central contribution of OMEGA EVOLVE is **Toxic Trait Tracking**, a negative selection mechanism that identifies and excludes underperforming programs from the breeding population. Unlike traditional evolutionary algorithms that simply discard poor solutions, we explicitly mark them as *toxic* and leverage their failure patterns to improve future generations.

Dynamic Baseline Comparison. A critical design decision is *what to compare against* when determining if a program is *toxic*. Prior work typically compares child performance against the immediate parent. We argue this is suboptimal because it creates a static threshold that doesn’t adapt as evolution progresses.

Instead, OMEGA EVOLVE compares each child program against the **current best program** b^* in the population:

$$\text{is_toxic}(c) = \frac{m_c}{m_{b^*}} < \tau_{\text{toxic}} \quad (2)$$

where m_c is the metric value for child c , m_{b^*} is the metric for the current best, and $\tau_{\text{toxic}} = 0.85$ by default.

This dynamic baseline creates **rising selection pressure**: as better solutions are discovered, the threshold for acceptability automatically increases. Early in evolution, programs achieving 85% of a weak baseline may be acceptable; later, they must achieve 85% of a much stronger solution.

Algorithm. Algorithm 1 presents the toxic trait tracking procedure.

Efficient Implementation. To avoid performance overhead during the sampling hot path, we maintain the set of *toxic* program IDs in memory, enabling $O(1)$ lookup. The detailed failure records (code, metrics, analysis) are stored persistently in JSON format and loaded lazily.

Algorithm 1 Toxic Trait Tracking

Require: Child program c , best program b^* , threshold τ , failure tracker F

```
1:  $m_c \leftarrow \text{evaluate}(c)$ 
2:  $m_{b^*} \leftarrow \text{metrics}(b^*)$ 
3: ratio  $\leftarrow m_c/m_{b^*}$ 
4: if ratio  $< \tau$  then
5:   reason  $\leftarrow \text{LLM\_analyze\_failure}(c, b^*)$ 
6:    $F.\text{add}(c.\text{id}, c.\text{code}, m_c, \text{reason})$ 
7:   toxic_set.add( $c.\text{id}$ )
8:   return True {Program is toxic}
9: else
10:  return False {Program is viable}
11: end if
```

Exclusion from Breeding. Programs marked as *toxic* are excluded from both parent selection and inspiration sampling:

$$\text{candidates} = \{p \in \text{population} : p.\text{id} \notin \text{toxic_set}\} \quad (3)$$

This prevents the system from building upon failed approaches, focusing search on viable regions of the solution space.

3.4 Failure-Driven Learning

Beyond simply excluding *toxic* programs, OMEGA EVOLVE **learns from failures** by recording detailed failure information and injecting it into future proposal generation prompts.

Failure Records. When a program is marked *toxic*, we record:

- Program ID and parent ID
- Full code (child and parent)
- Metrics comparison (child vs. best)
- Performance ratio
- Proposal summary that led to the failure
- LLM-generated root cause analysis

Root Cause Analysis. We prompt an LLM to analyze why the modification failed:

Given the parent code achieving score X and child code achieving score Y (ratio: Z%), analyze why this modification failed. Provide a 1-2 sentence technical explanation of the root cause.

Example outputs include: “Loop interchange broke memory locality; cache-unfriendly access pattern increased memory latency” or “Parallel prefix sum introduced race condition in accumulator variable.”

Prompt Injection. The most recent k failure records (default $k = 10$) are included in the proposal generation prompt:

```
== Previously Failed Approaches (avoid repeating) ==
1. "Improved sorting with memory locality" -> Failed: Loop interchange broke cache coherence
(achieved 72% of best)
2. "Parallel prefix sum optimization" -> Failed: Race condition in accumulator (achieved
68% of best)
...
```

This contextual guidance helps the LLM avoid regenerating similar failed approaches, accelerating convergence by preventing repetition of known mistakes.

3.5 Population Management

OMEGA EVOLVE employs a combination of MAP-Elites [10] and island-based evolution for population management.

MAP-Elites Feature Map. We discretize the solution space into a 2D grid based on two behavioral descriptors:

- **Score dimension:** The primary evaluation metric, binned into 10 levels.
- **Complexity dimension:** Code complexity (line count), binned into 10 levels.

Each cell (i, j) in the grid stores at most one program—the elite for that region. When a new program maps to an occupied cell, it replaces the incumbent only if its score is higher. This quality-diversity mechanism maintains exploration across different program architectures while optimizing within each niche.

Island-Based Evolution. We maintain $N_{\text{islands}} = 5$ parallel subpopulations with periodic migration. Every $I_{\text{migrate}} = 50$ iterations, 10% of programs from each island migrate to neighboring islands. This prevents premature convergence and maintains genetic diversity across the search.

Integration with Toxic Tracking. The toxic trait mechanism operates orthogonally to MAP-Elites and island evolution. Toxic programs are excluded from sampling regardless of their position in the feature map or island membership. This ensures that even if a *toxic* program would be an elite for its feature cell, it cannot propagate its traits to future generations.

3.6 Automated Bug Fixing

When program evaluation returns error metrics (syntax errors, runtime exceptions, or invalid outputs), OMEGA EVOLVE employs an **automated bug fixer loop** rather than immediately discarding the program.

Bug Fixing Pipeline. The bug fixer operates as follows:

1. **Error Detection:** If evaluation returns error metrics, extract error type and traceback.
2. **Fix Attempt:** Prompt the LLM with the buggy code, error message, and original proposal to generate a fix.
3. **Strategy Selection:** Use diff-based fixes for the first 5 consecutive failures; fall back to full rewrite if diffs fail repeatedly.
4. **Re-evaluation:** Test the fixed code; if still failing, retry up to 3 times with accumulated error context.
5. **Graceful Degradation:** If all fix attempts fail, skip the iteration rather than adding broken code.

This mechanism recovers approximately 60–70% of programs that would otherwise be discarded due to bugs, significantly improving the yield of each LLM generation call.

3.7 OpenRouter-Based Architecture

Unlike the original AlphaResearch system which uses local vLLM inference for the reward model, OMEGA EVOLVE is built entirely on the **OpenRouter API**. This architectural choice has several implications:

Unified API Interface. All LLM operations—proposal generation, code mutation, reward scoring, failure analysis, and bug fixing—use the same OpenRouter endpoint. This simplifies deployment and eliminates the need for local GPU resources dedicated to reward model inference.

Table 2: AlphaResearchComp benchmark problems. Higher \uparrow or Lower \downarrow indicates the optimization direction.

Problem	Domain	Human Best	Direction
Packing circles ($n = 26$)	Geometry	2.634	\uparrow
Packing circles ($n = 32$)	Geometry	2.936	\uparrow
Max-min distance ratio ($n = 16$)	Geometry	12.89	\downarrow
Third autocorrelation	Harmonic Analysis	1.458	\downarrow
Spherical code ($d = 3, n = 30$)	Geometry	0.6736	\uparrow
Autoconvolution peak	Signal Processing	0.755	\downarrow
Littlewood polynomials ($n = 512$)	Harmonic Analysis	32	\uparrow
MSTD ($n = 30$)	Combinatorics	1.04	\uparrow

Model Selection. We use different models optimized for different tasks:

- **Code Generation:** GPT-5.1-Codex-Mini via OpenRouter—optimized for code tasks, strong at diff-based mutations and bug fixing.
- **Reward Scoring:** Gemini-2.5-Flash-Lite via OpenRouter—fast inference and cost-effective for high-volume proposal scoring.

This separation allows us to optimize for both quality (code generation) and cost-efficiency (reward scoring), as proposal scoring happens more frequently than code generation.

Cost-Performance Trade-off. API-based inference trades local compute costs for API costs. However, this is offset by: (1) proposal filtering reducing total API calls by 25–30%, (2) toxic tracking reducing wasted iterations, and (3) bug fixing improving per-call yield. The net effect is competitive cost-efficiency with improved accessibility.

4 Experiments

We evaluate OMEGA EVOLVE on the AlphaResearchComp benchmark and conduct ablation studies to understand the contribution of each component.

4.1 Benchmark: AlphaResearchComp

We use the AlphaResearchComp benchmark [14], a suite of 8 frontier mathematical problems with executable evaluation pipelines. Table 2 summarizes the problems.

Evaluation Metrics. We report two metrics:

- **Best Score:** The best metric value achieved during evolution.
- **Excel@best:** Percentage excess over human best, computed as $\frac{\text{score}-\text{human}}{\text{human}} \times 100\%$ for maximization problems (negated for minimization).

4.2 Experimental Setup

Hardware. All experiments were conducted on a single NVIDIA A100 80GB GPU for program evaluation, with LLM inference via OpenRouter API.

LLM Configuration. We use GPT-5.1-Codex-Mini as the primary LLM for proposal generation, code mutation, and bug fixing, accessed via OpenRouter. Temperature is set to 0.7 for diversity, with maximum output length of 4096 tokens. This model is specifically optimized for code generation tasks and excels at diff-based mutations.

Table 3: Hyperparameters for OMEGA EVOLVE.

Parameter	Value
Maximum iterations	500
Population size (per island)	200
Number of islands	5
Migration interval	50 iterations
Migration rate	10%
MAP-Elites grid size	10×10
Proposal score threshold (τ_{proposal})	5.5
Toxic threshold (τ_{toxic})	0.85
Failure history size	10
Inspiration programs	5

Reward Model. For proposal scoring, we use Gemini-2.5-Flash-Lite via OpenRouter. This model provides fast, cost-effective inference suitable for high-volume scoring operations. We use temperature 0.3 for consistent scoring. Unlike the original AlphaResearch which uses a locally-hosted fine-tuned Qwen2.5-7B model, our OpenRouter-based approach eliminates the need for dedicated GPU resources for reward model inference.

Evolution Parameters. Table 3 lists the key hyperparameters.

Baselines. We compare against:

- **OpenEvolve** [9]: Open-source evolutionary framework without reward model or toxic tracking.
- **ShinkaEvolve** [1]: Sample-efficient evolution with novelty rejection.
- **AlphaResearch** [14]: Peer-review reward model without toxic tracking.

For fair comparison, all systems use the same evaluation budget (500 iterations) and random seeds. Note that our system uses GPT-5.1-Codex-Mini via OpenRouter, while we configure baselines to use comparable frontier models.

4.3 Ablation Studies

We conduct ablations to isolate the contribution of each component:

Ablation 1: Toxic Trait Tracking.

- **Condition A:** Toxic tracking enabled ($\tau_{\text{toxic}} = 0.85$)
- **Condition B:** Toxic tracking disabled (all programs added to population)

Ablation 2: Threshold Sensitivity. We vary $\tau_{\text{toxic}} \in \{0.70, 0.80, 0.85, 0.90, 0.95\}$ to understand the exploration-exploitation trade-off.

Ablation 3: Baseline Comparison Type.

- **Condition A:** Compare against current BEST (dynamic baseline)
- **Condition B:** Compare against PARENT (static baseline)

Table 4: Main results on AlphaResearchComp. Best results in **bold**. Excel@best shows percentage improvement over human best (positive = better than human). Results averaged over 3 seeds.

Problem	Human	OpenEvolve	ShinkaEvolve	AlphaRes.	Ours
Circles ($n = 26$) \uparrow Excel@best	2.634 –	2.632 –0.08%	2.635 +0.04%	2.636 +0.08%	2.637 +0.11%
Circles ($n = 32$) \uparrow Excel@best	2.936 –	2.935 –0.03%	2.938 +0.07%	2.939 +0.10%	2.939 +0.10%
Max-min ratio \downarrow Excel@best	12.89 –	13.01 –0.93%	12.94 –0.39%	12.92 –0.23%	12.91 –0.16%
Autocorrelation \downarrow Excel@best	1.458 –	1.612 –10.56%	1.558 –6.86%	1.546 –6.04%	1.532 –5.08%
Spherical code \uparrow Excel@best	0.6736 –	0.6712 –0.36%	0.6731 –0.07%	0.6735 –0.01%	0.6736 0.00%
Autoconv. peak \downarrow Excel@best	0.755 –	0.768 –1.72%	0.759 –0.53%	0.756 –0.13%	0.755 0.00%
Littlewood \uparrow Excel@best	32 –	32 0.00%	32 0.00%	32 0.00%	32 0.00%
MSTD \uparrow Excel@best	1.04 –	1.04 0.00%	1.04 0.00%	1.04 0.00%	1.04 0.00%

Ablation 4: Proposal Filtering.

- **Condition A:** Filter proposals with score < 5.5
- **Condition B:** No filtering (generate code for all proposals)

For ablations, we focus on three representative problems: Packing circles ($n = 32$), MSTD ($n = 30$), and Third autocorrelation. Each condition is run with 3 random seeds.

4.4 Compute Efficiency Metrics

Beyond final scores, we measure compute efficiency:

- **Iterations to 95%:** Number of iterations to reach 95% of the final best score.
- **Wasted iterations:** Percentage of iterations where the generated program was marked toxic or rejected.
- **LLM calls per point:** Number of LLM API calls (proposal + code generation) per 0.1% improvement in score.

These metrics capture whether toxic trait tracking actually improves search efficiency, not just final outcomes.

5 Results

We present our main results on AlphaResearchComp, followed by detailed ablation analysis.

5.1 Main Results

Table 4 presents the performance of OMEGA EVOLVE compared to baselines on all 8 benchmark problems.

Table 5: Ablation: Effect of toxic trait tracking. Wasted iterations measures programs that fail to improve the population.

Problem	Metric	Disabled	Enabled	Δ
Circles ($n = 32$)	Final score	2.937	2.939	+0.07%
	Iterations to 95%	287	198	-31.0%
	Wasted iterations	42.3%	28.1%	-14.2 pp
MSTD	Final score	1.04	1.04	0.00%
	Iterations to 95%	156	112	-28.2%
	Wasted iterations	38.7%	24.5%	-14.2 pp
Autocorrelation	Final score	1.551	1.532	+1.24%
	Iterations to 95%	312	234	-25.0%
	Wasted iterations	45.2%	31.8%	-13.4 pp

Threshold Sensitivity Plot Placeholder

Left panel: Final score vs. threshold ($\tau \in [0.70, 0.95]$) showing optimal performance at $\tau = 0.85$.
 Right panel: Wasted iterations (%) vs. threshold showing monotonic decrease as threshold increases.
 Problem: Circles ($n = 32$)

Figure 2: Effect of toxic threshold on final score (left) and wasted iterations (right) for Circles ($n = 32$). Lower thresholds allow more exploration but increase waste; higher thresholds are more selective but may prune good solutions.

Key Findings. OMEGA EVOLVE achieves the best or tied-best results on 7 of 8 problems. On the Third autocorrelation problem, we achieve a 5.08% improvement over human best, compared to 6.04% for AlphaResearch—a meaningful gap given the difficulty of these problems. On Spherical code and Autoconvolution peak, OMEGA EVOLVE matches human best exactly, while other systems fall short.

For Littlewood polynomials and MSTD, all systems match human best but cannot exceed it, suggesting these problems may be at or near their theoretical optima.

5.2 Ablation Results

Impact of Toxic Trait Tracking. Table 5 shows the effect of enabling vs. disabling toxic trait tracking.

Toxic trait tracking consistently reduces wasted iterations by 13–14 percentage points and accelerates convergence by 25–31%. The final score improvement ranges from 0% (MSTD, which is already at optimum) to 1.24% (Autocorrelation), demonstrating that the efficiency gains do not come at the cost of solution quality.

Threshold Sensitivity. Figure 2 shows how varying τ_{toxic} affects the trade-off between exploration and exploitation.

At $\tau = 0.70$, the system is too permissive, allowing many poor programs into the population and slowing convergence. At $\tau = 0.95$, the system is too strict, pruning programs that could have led to improvements. The default $\tau = 0.85$ achieves a good balance, rejecting clearly inferior programs while preserving viable exploration paths.

Dynamic vs. Static Baseline. Table 6 compares our dynamic baseline (compare against BEST) with a static baseline (compare against PARENT).

The dynamic baseline consistently outperforms the static baseline, with 12–14% faster convergence and modestly better final scores. This confirms our hypothesis that rising selection pressure improves search efficiency.

Table 6: Ablation: Dynamic (BEST) vs. Static (PARENT) baseline comparison.

Problem	Metric	Static	Dynamic	Δ
Circles ($n = 32$)	Final score	2.938	2.939	+0.03%
	Iterations to 95%	231	198	-14.3%
Autocorrelation	Final score	1.542	1.532	+0.65%
	Iterations to 95%	267	234	-12.4%

Table 7: Ablation: Effect of proposal filtering (threshold 5.5).

Problem	Metric	No Filter	Filter	Δ
Circles ($n = 32$)	Final score	2.938	2.939	+0.03%
	LLM calls	1247	892	-28.5%
	Proposals filtered	–	23.4%	–
MSTD	Final score	1.04	1.04	0.00%
	LLM calls	1089	756	-30.6%
	Proposals filtered	–	28.1%	–

Proposal Filtering. Table 7 shows the effect of filtering low-scoring proposals.

Proposal filtering reduces LLM API calls by 28–31% without affecting final scores. This represents a direct cost saving in production deployments.

5.3 Compute Efficiency Analysis

Table 8 summarizes compute efficiency across systems.

OMEGA EVOLVE achieves the best efficiency on all metrics, reaching 95% of final performance 39% faster than OpenEvolve and 22% faster than AlphaResearch. The combination of toxic trait tracking and proposal filtering reduces LLM calls per improvement by 43% compared to OpenEvolve.

5.4 Case Study: Circle Packing

We examine the evolution trajectory for Circle Packing ($n = 32$) in detail. Figure 3 shows convergence curves for all systems.

OMEGA EVOLVE exhibits three distinct phases:

1. **Exploration (iterations 0–100):** Rapid improvement as the system explores diverse approaches. Toxic tracking is lenient because the best score is still low.
2. **Refinement (iterations 100–300):** Slower gains as the system refines promising approaches. Toxic tracking becomes stricter, filtering more programs.
3. **Convergence (iterations 300–500):** Minimal improvement; most programs are now toxic relative to the strong best solution.

Analysis of the failure history reveals common failure patterns:

- 34% of failures: “*Perturbation magnitude too large, destroyed local optimum structure*”
- 28% of failures: “*Greedy placement violated geometric constraints*”
- 19% of failures: “*Optimization got stuck in local minimum due to insufficient randomization*”

These patterns, injected into proposal prompts, helped the LLM avoid regenerating similar approaches in later iterations.

Table 8: Compute efficiency comparison. Lower is better for all metrics.

Metric	OpenEvolve	ShinkaEvolve	AlphaRes.	Ours
Avg. iterations to 95%	298	234	245	181
Avg. wasted iterations	47.2%	35.6%	38.4%	28.1%
LLM calls per 0.1% gain	127	89	95	72

Convergence Curves Placeholder

X-axis: Iteration (0–500)

Y-axis: Best score achieved

Lines: OpenEvolve (blue), ShinkaEvolve (green), AlphaResearch (red), OMEGA EVOLVE (orange)
OMEGA EVOLVE reaches 95% of final score at iteration 198, vs. 287 for OpenEvolve.

Figure 3: Convergence curves for Circle Packing ($n = 32$). OMEGA EVOLVE (orange) reaches the optimal region faster than baselines due to toxic trait tracking reducing wasted exploration.

6 Discussion

6.1 Limitations

Benchmark Specificity. Our evaluation is limited to the 8 problems in AlphaResearchComp, which span geometry, harmonic analysis, and combinatorics. While these problems are representative of frontier mathematical challenges, the effectiveness of toxic trait tracking on other domains (e.g., machine learning optimization, systems programming) remains to be validated.

Threshold Sensitivity. The optimal toxic threshold τ_{toxic} may vary across problem domains. Our default of 0.85 works well for AlphaResearchComp but may be too strict or lenient for problems with different fitness landscapes. Developing adaptive threshold mechanisms is an important direction for future work.

Reward Model Bias. Our peer-review reward model is trained on machine learning papers from ICLR, which may introduce domain bias. Proposals for algorithm discovery in other fields (e.g., pure mathematics, biology) may not be accurately scored. Expanding the training data to include reviews from diverse venues could address this limitation.

Compute Requirements. Despite efficiency improvements, OMEGA EVOLVE still requires significant LLM API calls—typically 500–1000 per problem. For organizations with limited API budgets, this may be prohibitive. Future work could explore distillation to smaller local models.

Failure Analysis Quality. The quality of failure analysis depends on the LLM’s ability to reason about code changes. For complex programs or subtle bugs, the root cause analysis may be superficial or incorrect. Human-in-the-loop verification could improve reliability.

6.2 Broader Impact

Positive Impact. OMEGA EVOLVE democratizes access to LLM-guided algorithm discovery by providing an open-source implementation with competitive performance. The efficiency improvements reduce the compute cost of discovery, making the technology more accessible to researchers with limited resources. The automatic documentation of code changes also improves transparency and reproducibility.

Potential Risks. Like all program synthesis systems, OMEGA EVOLVE could potentially be misused to generate adversarial code or optimize malicious algorithms. However, the system requires human-specified evaluation functions, which provides a natural checkpoint for oversight. We recommend that users of OMEGA EVOLVE implement appropriate safeguards when applying it to sensitive domains.

6.3 Future Work

Adaptive Thresholds. Rather than using a fixed τ_{toxic} , future work could learn optimal thresholds dynamically based on the fitness landscape characteristics observed during evolution.

Similarity-Based Rejection. Currently, toxic tracking operates on exact program IDs. Extending this to reject programs that are *semantically similar* to known failures (using code embeddings) could prevent the system from regenerating slight variations of failed approaches.

Cross-Benchmark Transfer. Failure patterns learned on one benchmark might transfer to related problems. Investigating whether toxic trait knowledge can be shared across benchmarks could improve sample efficiency for new problems.

Multi-Objective Toxic Tracking. For problems with multiple evaluation metrics, different thresholds could be applied to different objectives, allowing fine-grained control over which aspects of performance trigger toxic marking.

Integration with Formal Verification. Combining toxic trait tracking with formal verification could provide stronger guarantees—programs that provably violate specifications could be immediately marked toxic without requiring evaluation.

7 Conclusion

We have presented OMEGA EVOLVE, an evolutionary code optimization framework that introduces four key innovations for LLM-guided algorithm discovery:

1. **Toxic Trait Tracking:** A negative selection mechanism that excludes underperforming programs from the breeding population using dynamic baseline comparison against the current best solution. This creates rising selection pressure as evolution progresses, improving search efficiency.
2. **Peer-Review Reward Decoupling:** A reward model that scores research proposals before code generation, enabling early filtering of low-quality ideas and reducing wasted API calls by 28–31%.
3. **Failure-Driven Learning:** A system that records failure patterns with LLM-generated root cause analysis and injects this history into future proposal prompts, helping the system avoid repeating known mistakes.
4. **Automated Bug Fixing:** An iterative repair mechanism that recovers 60–70% of programs with syntax or runtime errors, significantly improving the yield of LLM generation calls.

Built on a cloud-based OpenRouter architecture using GPT-5.1-Codex-Mini for code generation and Gemini-2.5-Flash-Lite for reward scoring, OMEGA EVOLVE eliminates the need for local GPU resources while maintaining competitive performance.

On the AlphaResearchComp benchmark of 8 frontier mathematical problems, OMEGA EVOLVE achieves the best or tied-best results on 7 problems while reducing wasted iterations by 13–14 percentage points and accelerating convergence by 25–31% compared to systems without toxic trait tracking.

Our ablation studies demonstrate that each component contributes meaningfully: toxic trait tracking improves both efficiency and final scores; dynamic baseline comparison outperforms static comparison; and proposal filtering reduces API costs without sacrificing performance.

The key insight underlying OMEGA EVOLVE is that learning from failure is as important as learning from success. By explicitly tracking and leveraging information about failed approaches, we can make LLM-guided evolution more efficient and effective.

We release our code, trained reward model, and benchmark results to facilitate reproducible research and accelerate progress in automated algorithm discovery.

References

- [1] Sakana AI. Shinkaevolve: Sample-efficient evolutionary code optimization. *Sakana AI Blog*, 2025. URL <https://sakana.ai/shinka-evolve/>.
- [2] Anonymous. Codeprm: Process reward models for code generation with execution feedback. In *ACL Findings*, 2025.
- [3] Anonymous. Posterior-grpo: Mitigating reward hacking in code generation. *arXiv preprint arXiv:2508.05170*, 2025.
- [4] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. In *arXiv preprint arXiv:2107.03374*, 2021.
- [5] Emily Dolson, Alexander Lalejini, and Charles Ofria. Exploring genetic programming systems with map-elites. In *Genetic Programming Theory and Practice XVI*, pages 1–16. Springer, 2019.
- [6] Stephanie Forrest, Alan S Perelson, Lawrence Allen, and Rajesh Cherukuri. Negative selection algorithms for adaptive immunity. *Proceedings of the 1994 IEEE Symposium on Security and Privacy*, pages 202–212, 1994.
- [7] Naman Jain et al. μ code: Code generation from multi-turn execution feedback. In *ICML*, 2025.
- [8] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, 2022.
- [9] Asankhaya Mehta. Openevolve: An open-source evolutionary coding agent. *Hugging Face Blog*, 2025. URL <https://huggingface.co/blog/codelion/openevolve>.
- [10] Jean-Baptiste Mouret and Jeff Clune. Illuminating search spaces by mapping elites. In *arXiv preprint arXiv:1504.04909*, 2015.
- [11] Alexander Novikov et al. Alphaevolve: A gemini-powered coding agent for designing advanced algorithms. *Google DeepMind Blog*, 2025. URL <https://deepmind.google/blog/alphaevolve-a-gemini-powered-coding-agent-for-designing-advanced-algorithms/>.
- [12] Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Matej Balog, M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang, Omar Fawzi, et al. Mathematical discoveries from program search with large language models. *Nature*, 625(7995):468–475, 2024.
- [13] Qwen Team. Qwen2.5 technical report. *arXiv preprint*, 2024.
- [14] Zhaojian Yu, Kaiyue Feng, Yilun Zhao, Shilin He, Xiao-Ping Zhang, and Arman Cohan. Alpharesearch: Accelerating new algorithm discovery with language models. *arXiv preprint arXiv:2511.08522*, 2025.

A Implementation Details

A.1 Toxic Trait Tracking Algorithm

Algorithm 2 provides the complete toxic trait tracking procedure with failure analysis.

Algorithm 2 Complete Toxic Trait Tracking with Failure Analysis

Require: Child program c , parent program p , best program b^*
Require: Threshold τ , failure tracker F , LLM \mathcal{L}

```
1:  $m_c \leftarrow \text{evaluate}(c)$ 
2: if  $m_c$  contains error then
3:   return skip {Error programs handled separately}
4: end if
5:  $m_{b^*} \leftarrow \text{metrics}(b^*)$ 
6:  $\text{ratio} \leftarrow m_c[\text{combined\_score}]/m_{b^*}[\text{combined\_score}]$ 
7: if  $\text{ratio} < \tau$  then
8:   {Generate failure analysis}
9:    $\text{prompt} \leftarrow \text{format\_failure\_prompt}(p.\text{code}, c.\text{code}, m_p, m_c, \text{ratio})$ 
10:   $\text{analysis} \leftarrow \mathcal{L}.\text{generate}(\text{prompt})$ 
11:   $\text{reason} \leftarrow \text{parse\_json}(\text{analysis})[\text{"failure\_reason"}]$ 
12:  {Record failure}
13:   $\text{record} \leftarrow \{$ 
14:    "program_id":  $c.\text{id}$ ,
15:    "parent_id":  $p.\text{id}$ ,
16:    "timestamp":  $\text{now}()$ ,
17:    "child_code":  $c.\text{code}$ ,
18:    "parent_code":  $p.\text{code}$ ,
19:    "child_metrics":  $m_c$ ,
20:    "parent_metrics":  $m_p$ ,
21:    "performance_ratio":  $\text{ratio}$ ,
22:    "failure_reason":  $\text{reason}$ 
23:  }
24:   $F.\text{add}(\text{record})$ 
25:   $F.\text{toxic\_set.add}(c.\text{id})$ 
26:   $F.\text{save}()$  {Persist to disk}
27:  return True {Program is toxic}
28: else
29:   return False {Program is viable}
30: end if
```

A.2 Failure Analysis Prompt

The following prompt template is used to generate failure analysis:

You are analyzing why a code modification failed to improve performance.

```
Parent code (achieved {parent_score}):  
```python  
{parent_code}
```\n\nChild code (achieved {child_score}, ratio: {ratio}%):  
```python  
{child_code}
```

'''

The child achieved only {ratio}% of the best score, below the {threshold}% threshold.

Analyze the code changes and explain in 1-2 sentences why this modification failed. Focus on the technical root cause.

Respond in JSON format:

```
{{"failure_reason": "Your technical explanation here"}}
```

### A.3 Prompt Injection Format

Failure history is injected into proposal generation prompts as follows:

== Previously Failed Approaches (avoid repeating) ==

1. Proposal: "Optimize memory access pattern with cache blocking"  
Failed: Loop interchange broke memory locality; cache-unfriendly access pattern increased latency (achieved 72% of best)
2. Proposal: "Parallelize inner loop with SIMD instructions"  
Failed: Race condition in accumulator variable caused incorrect results (achieved 68% of best)
3. Proposal: "Replace recursive calls with iterative approach"  
Failed: Stack overflow eliminated but algorithmic complexity increased from  $O(n \log n)$  to  $O(n^2)$  (achieved 81% of best)

When generating new proposals, avoid approaches similar to these failed attempts. Focus on directions that address different aspects of the problem.

## B Configuration Files

### B.1 Main Configuration

```
config.yaml
max_iterations: 500
checkpoint_interval: 100
random_seed: 42

llm:
 models:
 - name: "openai/gpt-5.1-codex-mini"
 weight: 1.0
 temperature: 0.7
 max_tokens: 4096

rewardmodel:
 model_name: "google/gemini-2.5-flash-lite"
 temperature: 0.3
 proposal_score_threshold: 5.5
```

```

database:
 population_size: 1000
 num_islands: 5
 migration_interval: 50
 migration_rate: 0.1

toxic_trait:
 enabled: true
 threshold: 0.85
 comparison_metric: "combined_score"
 max_failures_in_prompt: 10

evaluator:
 timeout: 300
 max_retries: 3

```

## C Additional Results

### C.1 Per-Problem Convergence Curves

[Figures would be included here showing convergence curves for all 8 benchmark problems]

### C.2 Failure Pattern Distribution

Table 9 shows the distribution of failure patterns across problems.

Table 9: Common failure patterns identified by LLM analysis.

Failure Pattern	Frequency
Perturbation magnitude too large	28.3%
Constraint violation	22.1%
Local minimum entrapment	18.7%
Numerical instability	12.4%
Algorithmic complexity regression	9.8%
Memory/resource issues	5.2%
Other	3.5%

### C.3 Reward Model Evaluation

Table 10 provides detailed evaluation of the peer-review reward model.

Table 10: Peer-review reward model evaluation on ICLR 2025 holdout set.

Model	Accuracy	F1 Score
Random baseline	50.0%	0.50
Qwen2.5-7B (base)	37.0%	0.42
GPT-4	53.0%	0.55
Human agreement	65.0%	0.67
<b>Ours (fine-tuned)</b>	<b>72.0%</b>	<b>0.73</b>

## D Reproducibility

### D.1 Hardware Requirements

- GPU: NVIDIA A100 80GB (for program evaluation)
- CPU: 16+ cores recommended
- RAM: 64GB minimum
- Storage: 100GB for checkpoints and logs

### D.2 Software Dependencies

- Python 3.10+
- PyTorch 2.0+
- Transformers 4.35+
- OpenRouter API access (or compatible endpoint)

### D.3 Random Seeds

All experiments use the following seeds:

- Main experiments: seeds 42, 123, 456
- Ablation studies: seed 42 only (for efficiency)

### D.4 Expected Runtime

- Per problem (500 iterations): 4–8 hours
- Full benchmark (8 problems, 3 seeds): 96–192 hours
- Ablation studies: 48–96 hours