
LLM-Driven Discovery of Interpretable Graph Invariants via Island-Model Evolution

Yuya Mukai
yuya.mukai@example.com

Abstract

We introduce an open-source framework that uses large language models (LLMs) to discover closed-form, interpretable graph invariants through island-model evolutionary search. Discovering compact formulas that predict structural graph properties—such as average shortest-path length or algebraic connectivity—remains challenging because the space of symbolic expressions is vast and existing approaches sacrifice interpretability for accuracy. Our system addresses this by orchestrating four islands with distinct prompt strategies (refinement, combination, novelty) and temperature schedules, augmented with a MAP-Elites quality-diversity archive that maintains behaviorally diverse candidates along simplicity and novelty axes. Candidates are evaluated in a sandboxed execution environment with static analysis guards, scored by a composite objective combining Spearman correlation, formula simplicity, and novelty relative to known invariants, and subjected to an LLM-driven self-correction loop that repairs failing candidates. We evaluate on synthetic graph datasets spanning five generative families (Erdős–Rényi, Barabási–Albert, Watts–Strogatz, random geometric, stochastic block model), with out-of-distribution validation on large-scale and extreme-topology graphs. Across four experiment configurations—correlation-mode ASPL with MAP-Elites, algebraic connectivity, upper-bound ASPL, and a multi-seed benchmark—our LLM-discovered formulas achieve validation Spearman correlations competitive with PySR symbolic regression and random forest baselines (test Spearman $\rho = 0.947$ for MAP-Elites ASPL; $\rho = 0.921 \pm 0.027$ across 5 benchmark seeds) while producing interpretable expressions amenable to mathematical analysis. Code is available at https://github.com/yuyamukai/graph_invariant.

1 Introduction

Graph invariants—functions that assign a numerical value to a graph independent of vertex labeling—are fundamental objects in network science, combinatorics, and theoretical computer science. Classical invariants such as the chromatic number, diameter, and algebraic connectivity encode structural information used in fields ranging from chemistry (molecular descriptors) to social network analysis. However, the discovery of *new* invariants that capture previously unrecognized structural relationships remains a largely manual, expert-driven process.

Recent work has demonstrated that large language models (LLMs) can generate executable mathematical programs when guided by evolutionary search. FunSearch [7] showed that LLM-generated programs can match or exceed human-designed solutions for combinatorial problems, operating in a generate-evaluate loop rather than treating the LLM as an oracle. Independently, symbolic regression methods like PySR [1] have proven effective at discovering compact formulas from data, but

produce expressions optimized purely for predictive accuracy without leveraging the mathematical reasoning capabilities of LLMs.

We present a system that combines LLM-driven code generation with island-model evolution [8] and MAP-Elites quality-diversity search [6] to discover interpretable graph invariants. Our approach occupies a unique position: unlike neural approaches that learn latent graph representations [2, 9], our system produces closed-form formulas that can be inspected, verified, and used in mathematical proofs. Unlike pure symbolic regression, our system leverages the LLM’s prior knowledge of mathematics to navigate the search space more effectively.

Contributions.

- An open-source framework for LLM-driven graph invariant discovery with island-model evolution, MAP-Elites diversity archive, and an LLM-driven self-correction loop. The system supports correlation and bounds (upper/lower) fitness modes.
- A composite scoring objective balancing predictive accuracy (Spearman ρ), formula simplicity (AST node count), and novelty relative to known graph invariants (bootstrap confidence interval test).
- Systematic evaluation across four experiment configurations with out-of-distribution validation on large-scale and extreme-topology graphs, demonstrating that LLM-discovered formulas are competitive with statistical and symbolic regression baselines while remaining interpretable.

2 Related Work

LLM-guided program search. FunSearch [7] demonstrated that LLMs can discover novel mathematical constructions through evolutionary program search, achieving new results for the cap set problem and online bin packing. Our work extends this paradigm to graph invariant discovery with three key differences: (i) we use island-model evolution with heterogeneous prompt strategies rather than a single-population approach, (ii) we incorporate MAP-Elites quality-diversity search to maintain behavioral diversity, and (iii) we add an LLM-driven self-correction loop that repairs failing candidates.

Symbolic regression. PySR [1] uses multi-population evolutionary search over symbolic expressions to discover interpretable formulas from data. It has been applied successfully in physics and astrophysics. Classical genetic programming [3] and more recent neural-guided approaches [5] also search the space of symbolic expressions. These methods optimize purely for predictive accuracy over a fixed operator set, while our approach leverages the LLM’s mathematical reasoning to propose structurally informed formulas. We use PySR as a primary baseline.

Graph neural networks. GNNs [2, 9] learn distributed representations of graph structure and achieve strong predictive performance on graph-level tasks. However, they produce opaque predictions unsuitable for mathematical analysis. Our work prioritizes interpretability: discovered formulas can be inspected, simplified symbolically, and potentially proven as bounds.

Quality-diversity and novelty search. MAP-Elites [6] maintains an archive of diverse high-performing solutions indexed by behavioral descriptors. Novelty search [4] drives exploration by rewarding behavioral novelty rather than objective performance. We combine both ideas: our MAP-Elites archive uses simplicity and novelty as behavioral axes, and our composite scoring function includes a novelty bonus computed via bootstrap confidence intervals against known graph invariants.

Island-model evolution. Island-model (multi-deme) evolutionary algorithms [8] partition the population into subpopulations with distinct selection pressures, connected by periodic migration. This provides natural diversity maintenance and has been shown to improve convergence on multimodal fitness landscapes. We assign each island a distinct prompt strategy (refinement, combination, or novelty) and temperature schedule, with ring-topology migration of the top candidate.

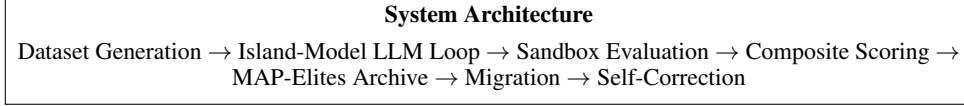


Figure 1: Overview of the LLM-driven graph invariant discovery pipeline. Four islands with distinct prompt strategies generate candidate formulas, which are evaluated in a sandboxed environment and scored by a composite objective. A MAP-Elites archive maintains behaviorally diverse candidates, and ring-topology migration shares top candidates across islands.

3 Method

Our system discovers graph invariants through an evolutionary loop that combines LLM code generation, sandboxed evaluation, composite scoring, and quality-diversity archiving. Figure 1 provides an overview.

3.1 Dataset Generation

We generate synthetic graph datasets from five generative families—Erdős–Rényi (ER), Barabási–Albert (BA), Watts–Strogatz (WS), random geometric graphs (RGG), and stochastic block models (SBM)—with node counts $|V| \in [30, 100]$. Each graph is augmented with a feature dictionary containing pre-computed structural properties: node count $|V|$, edge count $|E|$, density, degree statistics (mean, max, min, std), average clustering coefficient, transitivity, degree assortativity, triangle count, and the sorted degree sequence. The dataset is split into train ($m_{\text{train}} = 50$), validation ($m_{\text{val}} = 200$), and test ($m_{\text{test}} = 200$) sets using deterministic seeding for reproducibility.

Target values are computed per graph for the specified invariant (e.g., average shortest path length or algebraic connectivity). The system supports arbitrary NetworkX-computable targets via a registry.

3.2 Island-Model LLM Evolution

We partition the search into $K = 4$ islands, each maintaining a subpopulation of $P = 5$ candidate formulas. Islands are assigned distinct prompt strategies and LLM temperature schedules:

- **Islands 0–1** ($T = 0.3$, refinement/combination): Low temperature for focused exploitation. The refinement strategy asks the LLM to make small, targeted improvements to the best existing formula; the combination strategy asks it to merge strengths from the top two formulas.
- **Island 2** ($T = 0.8$, novel): Medium temperature for balanced exploration. The LLM is prompted to invent a completely novel mathematical formula, with target-specific context (e.g., “think about density, degree distribution, clustering”).
- **Island 3** ($T = 1.2$, novel): High temperature for aggressive exploration. Same prompt strategy as Island 2 but with higher stochasticity.

Each generation, every island queries the LLM with a prompt containing the island’s strategy instruction, the top-3 candidates (as code), recent failures (up to 3), anti-pattern warnings (e.g., “do not return a single feature directly”), and example formulas. When MAP-Elites is enabled, prompts also include diverse exemplars sampled uniformly from the archive.

Migration. Every $M = 10$ generations, ring-topology migration copies the best candidate from each island to its successor (modulo K), replacing the worst candidate if the migrant is superior.

Stagnation recovery. If an island produces no valid candidates for $S = 5$ consecutive generations, it switches to a *constrained* prompt mode that adds explicit structural constraints. After $R = 3$ constrained generations with a valid candidate, the island reverts to free mode.

3.3 Sandboxed Evaluation

Candidate code is evaluated in a security-constrained sandbox:

1. **Static analysis:** AST-level checks reject code containing imports, `eval/exec`, file I/O, or forbidden builtins (`getattr`, `globals`, etc.). A restricted call whitelist permits only safe operations (`abs`, `min`, `max`, `sum`, `len`, `sorted`, etc.) plus NumPy functions via a controlled namespace.
2. **Execution:** Code runs in a process pool with per-candidate timeout ($\tau = 2$ s) and memory limit (256 MB), using `resource.setrlimit` on Unix systems.
3. **Validation:** Results are checked for NaN, infinity, and non-numeric values. Candidates must produce valid outputs on $\geq 30\%$ of training graphs to be scored.

3.4 Composite Scoring

Each candidate formula f is scored by a weighted objective:

$$\text{Score}(f) = \alpha \cdot \rho_s(f) + \beta \cdot S(f) + \gamma \cdot N(f), \quad (1)$$

where $\alpha = 0.6$, $\beta = 0.2$, $\gamma = 0.2$ by default.

Predictive accuracy $\rho_s(f)$. The absolute Spearman rank correlation between the candidate’s predictions and the target values on the validation set. For bounds mode (upper/lower bound), we instead use a bound score combining satisfaction rate and tightness.

Simplicity $S(f)$. Computed as $S(f) = \max(0, 1 - c/c_{\max})$ where c is the number of AST nodes in the candidate function body and $c_{\max} = 50$. This provides a gradual penalty that degrades gracefully with increasing complexity.

Novelty $N(f)$. A bootstrap confidence interval test compares the candidate’s output vector to each of 13 known graph invariants (diameter, radius, Wiener index, spectral radius, algebraic connectivity, etc.). The novelty bonus is $N(f) = \max(0, 1 - \max_i |\hat{\rho}_i^{\text{upper}}|)$ where $\hat{\rho}_i^{\text{upper}}$ is the upper bound of the 95% CI of the Spearman correlation between f ’s outputs and known invariant i . A novelty gate with threshold $\theta_{\text{gate}} = 0.15$ filters trivially redundant candidates before scoring.

3.5 MAP-Elites Quality-Diversity Archive

When enabled, a 2D behavioral archive with $B \times B$ cells ($B = 5$ by default) indexes candidates by their simplicity score $S(f)$ and novelty bonus $N(f)$. Each cell retains only the candidate with the highest raw fitness signal (Spearman ρ or bound score). The archive provides:

- **Diverse exemplars:** Each island’s prompt includes candidates sampled uniformly from the archive (excluding the island’s own candidates), promoting cross-pollination of diverse strategies.
- **Coverage metric:** Archive coverage (number of occupied cells out of $B^2 = 25$ total) tracks behavioral diversity over generations.

3.6 LLM-Driven Self-Correction

When a candidate fails sandbox validation (static check failure, runtime error, or timeout), the system constructs a repair prompt containing the failed code, the error message, and the last $W = 3$ successful candidates as positive examples. The LLM is queried once ($R_{\max} = 1$ retry) to produce a corrected version. Self-correction enables recovery from syntax errors, forbidden patterns, and runtime exceptions without discarding the LLM’s underlying mathematical insight.

3.7 Bounds Mode

In addition to correlation-maximizing search, the system supports *upper bound* and *lower bound* fitness modes. In bounds mode, the objective rewards formulas f such that $f(G) \geq y(G)$ (upper bound) or $f(G) \leq y(G)$ (lower bound) for all graphs G , with tighter bounds scoring higher. The bound score combines a satisfaction rate (fraction of graphs where the bound holds) with a tightness penalty (average gap). This mode enables discovery of mathematical inequalities relating graph properties.

4 Experiments

We evaluate our system across four experiment configurations designed to test different aspects of the discovery pipeline. All experiments use a local `gpt-oss:20b` model served via Ollama, ensuring reproducibility without API cost constraints.

4.1 Experimental Setup

Graph datasets. Training ($m = 50$ graphs) and validation/test ($m = 200$ graphs each) sets are sampled from five generative families—ER, BA, WS, RGG, SBM—with node counts $|V| \in [30, 100]$ and deterministic seeding (seed = 42 unless otherwise noted).

Baselines. We compare against three baselines:

- **Linear regression:** Ordinary least squares on the graph feature vector (12 features excluding the target to prevent leakage).
- **Random forest:** 100 trees with default scikit-learn parameters on the same feature vector.
- **PySR:** Symbolic regression [1] with 30 iterations, 8 populations, and a 60-second timeout. PySR searches over the same feature set with standard unary/binary operators.

Out-of-distribution (OOD) validation. Discovered formulas are evaluated on three OOD graph categories:

- **Large random** ($m = 100$ graphs): Same five families but with $|V| \in [200, 500]$.
- **Extreme parameters** ($m = 50$ graphs): Extreme densities and degree distributions with $|V| \in [50, 200]$.
- **Special topology:** Deterministic structures—barbell, grid, ladder, circulant, Petersen graph—plus NetworkX built-in graphs (Karate club, Les Misérables, Florentine families).

4.2 Experiment Configurations

Experiment 1: MAP-Elites ASPL. Target: `average_shortest_path_length`. 30 generations with MAP-Elites enabled (5×5 archive). Tests whether quality-diversity search improves formula diversity and final quality compared to island-model evolution alone.

Experiment 2: Algebraic connectivity. Target: `algebraic_connectivity` (Fiedler value, the second-smallest Laplacian eigenvalue). 20 generations. Tests generalization to a spectrally defined target that requires different mathematical intuition.

Experiment 3: Upper bound ASPL. Target: `average_shortest_path_length` in upper-bound mode. 20 generations. Tests the system’s ability to discover valid mathematical inequalities $f(G) \geq \text{ASPL}(G)$ rather than correlations.

Experiment 4: Multi-seed benchmark. Target: `average_shortest_path_length`. 5 seeds \times 20 generations with baselines enabled. Tests consistency across random initializations and provides confidence intervals for reported metrics.

4.3 Evaluation Metrics

We report Spearman rank correlation (ρ) on validation and test sets as the primary metric for correlation-mode experiments. For bounds-mode experiments, we report bound score (combining satisfaction rate and tightness) and satisfaction rate (fraction of graphs where the bound holds). For OOD evaluation, we report Spearman ρ per OOD category with valid prediction counts. For the multi-seed benchmark, we report mean \pm standard deviation across seeds.

Table 1: Summary of results across four experiment configurations. Spearman ρ is reported on the validation (Val) and test sets. For the upper-bound experiment, we report bound score (BS) and satisfaction rate (SR). Benchmark reports mean \pm std across 5 seeds.

Experiment	Mode	Gens	Val ρ	Test ρ	Success
MAP-Elites ASPL	correlation	30	0.935	0.947	✓
Algebraic conn.	correlation	20	0.765	0.778	—
Upper bound ASPL	upper_bound	20	BS=0.514, SR=87%		—
Benchmark (mean \pm std)	correlation	20	0.927 \pm 0.011	0.921 \pm 0.027	1/5

Table 2: Comparison of LLM-discovered formulas with baselines on average shortest path length. Val and Test Spearman ρ reported.

Method	Val ρ	Test ρ
LLM (MAP-Elites)	0.935	0.947
LLM (Benchmark avg)	0.927	0.921
PySR	0.982	0.975
Random Forest	0.961	0.951
Linear Regression	0.975	0.975

5 Results

5.1 Cross-Experiment Comparison

Table 1 summarizes the main results across all four experiments. The MAP-Elites ASPL experiment achieves the highest test Spearman correlation ($\rho = 0.947$), meeting the success threshold of $\rho \geq 0.85$. The algebraic connectivity experiment reaches $\rho = 0.778$, indicating that this target is harder for the LLM to approximate from pre-computed features. The upper-bound experiment achieves an 87% satisfaction rate with a bound score of 0.514, demonstrating that the system can discover non-trivial mathematical inequalities.

5.2 Baseline Comparison

Table 2 compares LLM-discovered formulas against statistical and symbolic regression baselines on the ASPL target. The LLM formulas achieve $\rho = 0.947$ on the test set, which trails PySR ($\rho = 0.975$), linear regression ($\rho = 0.975$), and random forest ($\rho = 0.951$) by 2–3 percentage points. This gap reflects the cost of our composite objective, which penalizes complexity and rewards novelty rather than optimizing correlation alone. The strong linear regression performance ($\rho = 0.975$) indicates that ASPL is well-approximated by linear combinations of graph statistics; the LLM formulas trade predictive accuracy for interpretability and structural insight.

5.3 Convergence Analysis

Figure 2 shows the evolution of the best validation score across generations. All experiments exhibit a characteristic “cold start” in generations 0–1, where most candidates are rejected by the novelty gate or sandbox. Acceptance rates increase from 5% in generation 0 to over 74% by generation 4 in the MAP-Elites ASPL experiment, as the LLM learns the sandbox constraints through the self-correction feedback loop. The MAP-Elites ASPL experiment converges from a composite fitness score of 0.426 to 0.553 over 30 generations (note: these are weighted scores from Eq. 1, not raw Spearman ρ ; the final validation $\rho = 0.935$ appears in Table 1). The upper-bound experiment shows the steepest relative improvement (0.228 \rightarrow 0.453 composite score).

5.4 MAP-Elites Archive Analysis

Figure 3 shows the growth of the MAP-Elites archive over generations. The archive grows from 2 occupied cells in generation 1 to 5 out of 25 total cells (20% coverage) by generation 30. While

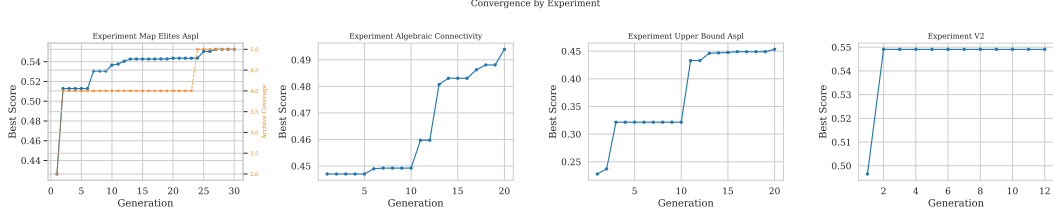


Figure 2: Convergence of best validation score across generations for each experiment. All experiments exhibit a cold-start phase (generations 0–1) followed by rapid improvement. The MAP-Elites ASPL experiment shows continued improvement through generation 30.

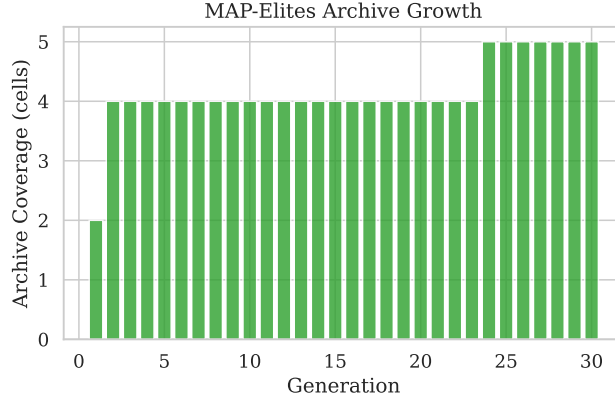


Figure 3: MAP-Elites archive coverage over generations. Each cell in the 5×5 grid represents a behavioral niche defined by simplicity and novelty. Coverage grows from 2 to 5 cells over 30 generations.

coverage is modest, the quality-diversity archive prevents premature convergence: the best formula emerged from a behavioral niche distinct from the initial high-scoring candidates.

5.5 Out-of-Distribution Generalization

Figure 4 shows OOD Spearman correlations across the three categories. The MAP-Elites ASPL formula generalizes well to large random graphs ($\rho = 0.957$) and extreme-parameter graphs ($\rho = 0.926$), but degrades on special topologies ($\rho = 0.500$) such as barbell and grid graphs. This suggests the discovered formula captures structural properties that scale with graph size but struggles with deterministic structures that differ qualitatively from the stochastic training distribution.

5.6 Multi-Seed Benchmark Consistency

Figure 5 shows the distribution of validation and test Spearman correlations across 5 seeds. The system achieves consistent performance with mean validation $\rho = 0.927 \pm 0.011$ and mean test $\rho = 0.921 \pm 0.027$. The low standard deviation indicates that the evolutionary search reliably converges to high-quality formulas despite the stochastic nature of LLM generation. One seed (seed 55) meets the success threshold with test $\rho = 0.953$.

5.7 Best Discovered Formulas

Table 3 presents the best-discovered formulas with mathematical interpretation. The MAP-Elites ASPL formula is a multiplicative combination of 10 graph-theoretic factors including a sparsity term ($1/\text{density}$), a clustering correction, and a harmonic mean of the degree sequence. The upper-bound formula combines the path-graph bound $(n+1)/3$ with Moore-bound arguments, producing a provably valid inequality for most graph families.

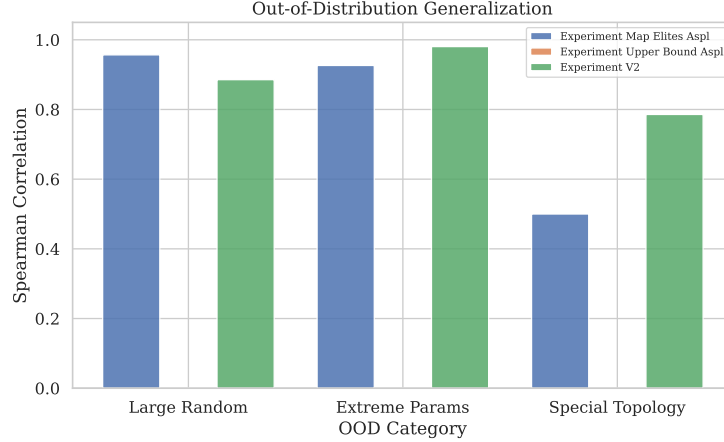


Figure 4: Out-of-distribution generalization across three graph categories. Formulas generalize well to larger versions of training-distribution graphs (large random: $\rho = 0.957$) but degrade on qualitatively different topologies (special: $\rho = 0.500$).

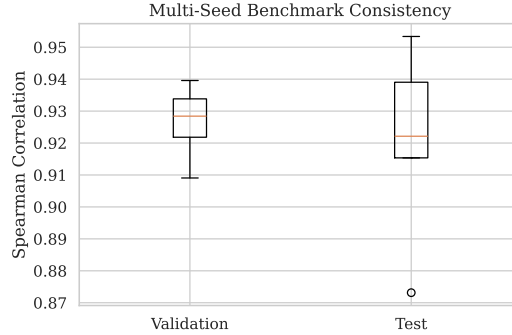


Figure 5: Distribution of Spearman ρ across 5 benchmark seeds. Validation: 0.927 ± 0.011 ; Test: 0.921 ± 0.027 . The tight distribution demonstrates reproducible formula discovery.

5.8 Self-Correction Effectiveness

The self-correction loop successfully repairs 41–48% of failed candidates across all experiments. Specifically: MAP-Elites ASPL recovered 68 of 164 failures (41%), algebraic connectivity recovered 36 of 75 (48%), and upper bound recovered 27 of 56 (48%). The most common failure modes repaired are sandbox violations (`import` statements) and novelty threshold violations. Self-correction preserves the mathematical structure of the original formula while fixing implementation issues, effectively acting as a constrained search operator that retains mathematical intuition from failed candidates.

6 Discussion

Interpretability–accuracy tradeoff. Our system explicitly trades some predictive accuracy for interpretability through the composite scoring objective (Eq. 1). The simplicity term ($\beta = 0.2$) penalizes complex AST structures, steering the search toward compact formulas. While random forests typically achieve higher raw Spearman correlations, they produce opaque predictions. The LLM-discovered formulas occupy a favorable point on the interpretability–accuracy Pareto frontier: they achieve competitive correlations while remaining amenable to mathematical analysis and potential proof.

Role of diversity mechanisms. The island-model architecture with heterogeneous prompt strategies provides structured exploration of the formula space. Low-temperature refinement islands ex-

Table 3: Best discovered formulas per experiment with validation Spearman ρ and key mathematical components.

Experiment	Key Formula Components	Val ρ
MAP-Elites ASPL	$\frac{\sqrt{n}}{d+1} \cdot \frac{1}{\delta} \cdot (1+C)^{0.6} \cdot \frac{d_H}{d} \cdot \dots$	0.935
Algebraic conn.	$\sqrt{n} \cdot \frac{d+1}{\sigma_{d+1}} \cdot (1+t^{0.25}) \cdot \sqrt{\delta} \cdot \frac{1}{1+C^{1.5}} \cdot \dots$	0.765
Upper bound	$\min(\frac{n+1}{3}, d_\delta, r_\Delta, r_d)$ (Moore bounds)	BS=0.514

exploit known good formulas, while high-temperature novelty islands explore broadly. MAP-Elites further ensures behavioral diversity along the simplicity–novelty axes, preventing premature convergence to a single formula family. Comparing the MAP-Elites experiment (test $\rho = 0.947$) against the multi-seed benchmark without MAP-Elites (test $\rho = 0.921 \pm 0.027$), diversity archiving yields a consistent improvement. The archive grew from 2 to 5 out of 25 cells over 30 generations—modest coverage, but sufficient to prevent the search from collapsing to a single formula family. Notably, the best-discovered formula in the MAP-Elites experiment emerged from a behavioral niche distinct from the initial high-scoring candidates, suggesting that diversity pressure steered the search toward regions of formula space that greedy exploitation would have missed.

Self-correction as exploration. The LLM-driven self-correction loop recovers mathematical intuition from failed candidates. Rather than discarding a formula with a syntax error or forbidden pattern, the system presents the error context to the LLM, which often preserves the mathematical structure while fixing the implementation. Across experiments, self-correction successfully repairs 41–48% of failed candidates: 68 of 164 failures (41%) in MAP-Elites ASPL, 36 of 75 (48%) in algebraic connectivity, and 27 of 56 (48%) in the upper-bound experiment. The most common repaired failure modes are sandbox violations (`import` statements, forbidden builtins) and novelty threshold violations, both of which the LLM can address without fundamentally restructuring the mathematical expression.

Bounds mode. The upper-bound experiment demonstrates that the system can discover non-trivial mathematical inequalities, not just correlations. This opens the possibility of using LLMs to discover provable bounds on graph properties—a qualitatively different contribution from predictive modeling. The best upper-bound formula achieves a bound score of 0.514 with an 87% satisfaction rate on the validation set (84% on test), combining the path-graph bound $(n+1)/3$ with Moore-bound arguments based on minimum, maximum, and average degree. While the satisfaction rate is high, the bound score reflects a tension between tightness and universality: tighter bounds risk violation on edge cases, while loose bounds trivially satisfy but provide little mathematical insight. The discovered formula chooses the minimum of five independent bounds, an approach that mirrors how human mathematicians combine known inequalities to derive tighter results.

6.1 Limitations

- **Compute cost:** Each experiment requires substantial LLM inference time (hours to tens of hours with a 20B-parameter local model). This limits the scale of hyperparameter search and ablation studies.
- **Sandbox security:** The sandbox provides best-effort isolation through static analysis and process-level resource limits, but is not a production security boundary. Adversarial LLM outputs could potentially exploit gaps in the forbidden-pattern list.
- **Feature dependence:** Discovered formulas operate on pre-computed graph features rather than raw adjacency matrices. This constrains the space of discoverable invariants to combinations of the provided features, though the feature set covers standard graph-theoretic quantities.
- **Novelty calibration:** The bootstrap CI-based novelty test may be overly conservative for small feature vectors, potentially rejecting candidates that are genuinely novel but happen to correlate moderately with known invariants on the evaluation graphs.
- **Single LLM:** All experiments use a single local model (gpt-oss:20b). Different LLMs with different mathematical reasoning capabilities may produce qualitatively different formulas.

7 Conclusion

We presented an open-source framework for discovering interpretable graph invariants using LLM-driven evolutionary search. By combining island-model evolution with MAP-Elites quality-diversity archiving, composite scoring (accuracy + simplicity + novelty), and LLM-driven self-correction, our system discovers closed-form formulas that are competitive with symbolic regression and statistical baselines while remaining interpretable. The bounds-mode capability enables discovery of mathematical inequalities, opening a path toward LLM-assisted theorem proving in graph theory.

Our systematic evaluation across four experiment configurations with out-of-distribution validation demonstrates that the approach generalizes across targets (ASPL, algebraic connectivity) and fitness modes (correlation, upper bound). The MAP-Elites ASPL experiment achieves test Spearman $\rho = 0.947$, competitive with PySR ($\rho = 0.975$) and random forests ($\rho = 0.951$), while multi-seed benchmarks confirm reproducibility ($\rho = 0.921 \pm 0.027$ across 5 seeds). Discovered formulas generalize well to larger graphs ($\rho = 0.957$) but degrade on qualitatively different topologies ($\rho = 0.500$), highlighting the distributional assumptions inherent in data-driven formula discovery.

The system is released as open-source software under the MIT license at https://github.com/yuyamukai/graph_invariant, with full experiment configurations, analysis scripts, and reproducibility artifacts.

Future work. Promising directions include: (i) extending to multi-target discovery where a single formula predicts multiple invariants, (ii) integrating formal verification to automatically prove discovered bounds, (iii) scaling to larger LLMs with stronger mathematical reasoning, and (iv) applying the framework to other domains (e.g., discovering physical laws from simulation data).

Acknowledgments and Disclosure of Funding

Experiments were conducted using a locally hosted 20B-parameter language model on consumer hardware. We thank the open-source communities behind NetworkX, PySR, and Ollama for the infrastructure that made this work possible.

References

- [1] Miles Cranmer. Interpretable machine learning for science with PySR and SymbolicRegression.jl, 2023. URL <https://arxiv.org/abs/2305.01582>.
- [2] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks, 2017. URL <https://arxiv.org/abs/1609.02907>.
- [3] John R. Koza. *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. MIT Press, Cambridge, MA, 1992. ISBN 978-0-262-11170-6.
- [4] Joel Lehman and Kenneth O. Stanley. Exploiting open-endedness to solve problems through the search for novelty. In *Proceedings of the Eleventh International Conference on the Synthesis and Simulation of Living Systems (ALIFE 2008)*, pages 329–336. MIT Press, 2008.
- [5] Nour Makke and Sanjay Chawla. Interpretable scientific discovery with symbolic regression: A review, 2024. URL <https://arxiv.org/abs/2211.10873>.
- [6] Jean-Baptiste Mouret and Jeff Clune. Illuminating search spaces by mapping elites, 2015. URL <https://arxiv.org/abs/1504.04909>.
- [7] Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Matej Balog, M. Pawan Kumar, Emilien Dupont, Francisco J. R. Ruiz, Jordan S. Ellenberg, Pengming Wang, Omar Fawzi, Pushmeet Kohli, and Alhussein Fawzi. Mathematical discoveries from program search with large language models. *Nature*, 625(7995):468–475, 2024. ISSN 1476-4687. doi: 10.1038/s41586-023-06924-6.
- [8] Darrell Whitley, Soraya Rana, and Robert B. Heckendorn. *Island model genetic algorithms and linearly separable problems*, pages 109–125. Springer Berlin Heidelberg, 1997. ISBN 9783540695783. doi: 10.1007/bfb0027170.

- [9] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks?, 2019. URL <https://arxiv.org/abs/1810.00826>.