
Decompose, Adapt, and Evolve: Towards Efficient Scientific Equation Discovery with Large Language Models

Pouya Behzadifar^{1*} Parshin Shojaee^{2*} Sanchit Kabra^{2*}
Kazem Meidani³ Chandan K. Reddy²

¹Sharif University of Technology ²Virginia Tech ³Capital One

Abstract

Finding mathematical relations underlying natural phenomena and scientific systems has been one of the fundamental tasks in the history of scientific discovery. Recent advancements in evolutionary search with Large Language Models (LLMs), with their embedded scientific knowledge, have shown great promise for this task. However, discovering such mathematical models governing scientific observations still remains significantly challenging, as it requires navigating vast combinatorial hypothesis spaces with an explosion of possible relations. Existing LLM-based approaches overlook the impact of data on the structure of mathematical relations, and treat LLMs as a static hypothesis generator unaware of the observed scientific system. This leads to inefficient exploration of the hypothesis space with over-reliance on LLMs’ internal priors. To bridge this gap, we introduce *Decompose, Adapt, Evolve (DecAEvolve)*, a framework that leverages granular feedback from symbolic term decomposition and LLM refinement through reinforcement learning (RL) fine-tuning to enhance both robustness and efficiency of evolutionary discovery frameworks. Our experiments across diverse datasets demonstrate that DecAEvolve significantly improves the accuracy of discovered equations and the efficiency of the discovery process compared to the state-of-the-art baseline.

1 Introduction

The emergence of Large Language Models (LLMs) has fundamentally transformed automated problem-solving across diverse domains. Beyond their well-established capabilities in natural language understanding and programming [1, 2], LLMs have recently demonstrated remarkable reasoning abilities that enable them to tackle complex optimization and discovery tasks. Their capacity to leverage embedded domain knowledge, interpolate between them, generate structured hypotheses and engage in iterative refinement, positions LLMs as powerful engines for systematic exploration of complex solution spaces towards discovery goals [3, 4, 5]. This potential extends naturally to scientific discovery tasks, where the combination of domain expertise and systematic search/exploration in the hypothesis space can unlock new approaches to longstanding challenges of scientific inquiry [6]. Scientific equation discovery—the process of uncovering compact and interpretable mathematical models that govern natural phenomena—represents one of the fundamental tasks in automated scientific discovery, with applications across many fields of science such as physics, biology, and material science [7]. Traditional approaches in equation discovery rely on genetic programming and evolutionary strategies [8, 9]; however, these approaches often struggle with scalability limitations and inefficient exploration of the vast combinatorial hypothesis space [10]. More recent advances

*Equal contribution.

To address these limitations, we introduce **DecAEvolve** (Decompose, Adapt, Evolve), a novel framework that enhances the effectiveness and efficiency of LLM-based equation discovery through several synergistic contributions:

- We develop a systematic methodology for providing LLMs with interpretable directional feedback about which components of their generated hypothesis prove effective. Through structured hypothesis decomposition and evaluations, the contributions of individual terms are quantified and provided as feedback. This enables LLMs to understand not just which hypotheses succeed, but *why* specific mathematical building blocks are effective.
- We employ reinforcement learning with Group-Relative Policy Optimization (GRPO) to implicitly encode the data distribution into the model’s parameters for better hypothesis generation process. This test-time adaptation/training approach allows the LLM to implicitly learn from successful equation discoveries, progressively aligning its hypothesis generation with the underlying symbolic relationships through reward-weighted gradient updates.
- We demonstrate that these contributions improve search effectiveness and efficiency, requiring significantly fewer iterations to discover accurate symbolic models. Our comprehensive evaluation across multiple benchmarks shows better performance compared to state-of-the-art baselines in both in-domain and out-of-domain settings.

2 Method

We propose **DecAEvolve** (Decompose, Adapt, Evolve), shown in Figure 1, a framework that shifts the evolutionary search of equation discovery towards guided discovery, achieved through granular and directional feedback as well as test-time adaptation with reinforcement learning fine-tuning of the backbone LLM to the observed scientific system. The core premise of our approach is that effective symbolic discovery requires the generator to learn not only *what* works, but also *why* it works and *how* to search. We implement this via two main components:

Directional Feedback with Term-Level Contribution. At the core of our framework is an iterative discovery process with LLM-based evolutionary search where the LLM generates candidate symbolic equations guided by structured, interpretable feedback. Unlike prior approaches that rely solely on coarse performance metrics, we introduce a fine-grained contribution analysis mechanism that quantifies the importance of individual terms and their interactions within discovered equations. In the contribution analysis, we parse the generated Python function of equation program skeleton into an Abstract Syntax Tree (AST) and decompose it into constituent terms. Addition and subtraction operations serve as natural term boundaries, while multiplicative structures, powers, and unary function calls (e.g., $\sin(x)$) are preserved as atomic units. Following decomposition, we conduct ablated contribution analysis by removing terms and measuring the resulting performance discrepancy. These ablated scores reveal contribution of symbolic components to the performance with respect to data. These contribution annotations are saved and passed to LLMs alongside the corresponding python equation program hypothesis that gets stored in the experience buffer. When the LLM encounters these annotated equations with directional feedback from term decomposition, it can see which terms drive better performance, and build upon them in the hypothesis generation of subsequent iterations.

Test-Time Adaptation with GRPO. To further enhance the LLM’s hypothesis generation capabilities, we incorporate a test-time adaptation approach using reinforcement learning fine-tuning with Group-Relative Policy Optimization (GRPO). This allows us to adapt the model to the specific scientific system observed data by learning from the distribution of successful equation discoveries. After each iteration of hypothesis generation, we collect a dataset of prompts paired with candidate equations and their corresponding rewards. Each equation is evaluated using negative MSE on the training data, which we transform to a bounded reward between 0 and 1 via $r = \exp(-\text{MSE})$. Failed or invalid completions receive a floor reward of 0.01. For each prompt x with k candidate hypotheses $\{y_i\}_{i=1}^k$, we compute group-relative advantages $A_i = r_i - b(x)$ where $b(x) = \frac{1}{k} \sum_i r_i$ serves as the baseline. The training objective balances reward maximization with a KL regularization term to prevent the model from drifting too far from its initial policy: $\mathcal{L}(\theta) = -\mathbb{E}_{x, \{y_i\}} \left[\sum_{i=1}^k A_i \log \pi_\theta(y_i|x) \right] + \beta \cdot \text{KL}(\pi_\theta || \pi_{\text{ref}})$. This GRPO training serves as an implicit mechanism for incorporating the underlying data distribution into the model’s hypothesis generation process to go beyond its internal priors.

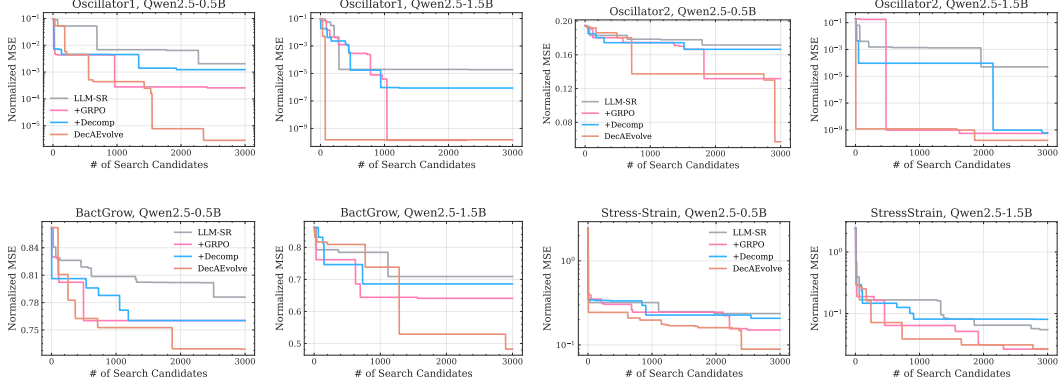


Figure 2: Best-score trajectories of DecAEvolve and its variants against the LLM-SR baseline across benchmark problems (lower is better). Each curve shows the average across five runs.

3 Experiments

As noted in [6], LLMs have significant memorization issues with the well-known physics equations of the commonly used equation discovery benchmarks such as Feynman [11]. So, instead of evaluating on this standard benchmark, we evaluate DecAEvolve on the benchmark datasets from [6] which are designed for LLM-based discovery methods to address this memorization issue, covering different scientific domains including physics, biology, and material science. We compare DecAEvolve with the state-of-the-art baseline LLM-SR [6] under same configurations: $3k$ LLM calls per problem with sampling temperature $\tau = 0.8$. As per [6], equation parameters are optimized with the BFGS solver from Scipy library and a 30s timeout is used for the execution of each hypothesis. In the GRPO adaptation phase, we use batch size of 16 per device, gradient accumulation 4, learning rate 10^{-6} , and KL coefficient $\beta = 0.05$. For fine-tuning, we use LoRA adapters with $r = 16$. Decomposition analysis is also limited to 7 terms and their pairwise interactions per program hypothesis. We conduct experiments on two open-source Qwen model variants (Qwen2.5-0.5B and Qwen2.5-1.5B). For the analysis, we use the normalized mean squared error (NMSE) as in [19] on both in-domain (ID) and out-of-domain (OOD) test settings. NMSE normalizes errors by scale of dataset variance, enabling comparison across datasets.

Results Figure 2 presents the discovery trajectories showing best-achieved NMSE scores across search iterations for DecAEvolve and its ablated variants compared to the state-of-the-art LLM-SR baseline. Each curve shows the average across five runs. The results demonstrate that both core components contribute meaningfully to performance: **Adaptation** (+GRPO) and **Decomposition** (+Decomp) consistently achieve lower discovery errors and converge faster than the LLM-SR baseline across all benchmark datasets. Notably, the full **DecAEvolve** framework, which integrates both components, delivers best performance in terms of both final accuracy (lower terminal NMSE) and search efficiency (faster convergence), establishing new state-of-the-art results across all scientific discovery tasks.

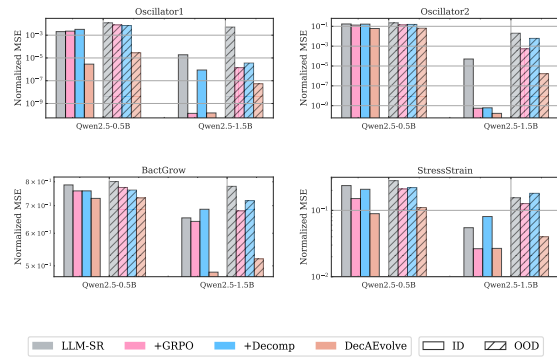


Figure 3: In-domain (ID) and out-of-domain (OOD) performance of DecAEvolve and its variants compared to LLM-SR (lower is better).

To evaluate the generalizability of discovered equations—a fundamental prerequisite for scientific equations and laws—we assess all methods on out-of-distribution (OOD) test data from [6] beyond their training domains. Figure 3 compares in-domain (ID) versus out-of-domain (OOD) NMSE perfor-

mance across all model variants and benchmark datasets. While all methods exhibit expected performance degradation on OOD data, **DecAEvolve** consistently achieves the lowest NMSE in both settings. DecAEvolve’s strong OOD performance indicates that our framework discovers equations with better inherent generalizability rather than merely fitting to the observed training distributions, a critical distinction for scientific discovery applications where extrapolation beyond observed data is essential.

Lastly, Figure 4 shows consistent reward improvement during GRPO adaptation across both model scales and all datasets, validating our reinforcement learning fine-tuning approach as test-time adaptation for equation discovery. Notably, we observe some interesting domain-dependent scaling behaviors: the smaller model (Qwen2.5-0.5B) converges faster on oscillator datasets, while the larger model (Qwen2.5-1.5B) shows better progression on bacterial growth and stress-strain tasks. We think that this pattern correlates with problem complexity for LLM, as evidenced by higher initial rewards on oscillator problems (0.6-0.7) compared to the more challenging biological and material systems (0.2-0.3). Interestingly, in all these cases, the smaller model eventually matches larger model performance even on complex datasets, suggesting that targeted adaptation through GRPO can help to effectively bridge the capability gap between model scales for the purpose of scientific discovery.

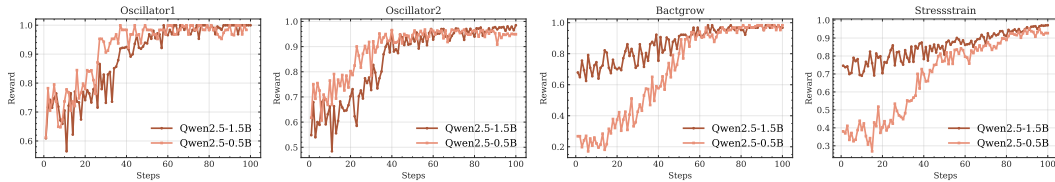


Figure 4: Reward improvement over steps during the GRPO adaptation phase across benchmark datasets.

4 Conclusion

We introduce DecAEvolve, a framework that enhances LLM-based equation discovery through granular term-level feedbacks, test-time adaptation via GRPO and, evolutionary search with LLMs. Our approach transforms static hypothesis generation into adaptive learning, enabling LLMs to progressively align with nuances of underlying observed scientific systems through reinforcement learning model adaptation and interpretable feedback mechanisms. Experimental results across diverse benchmark datasets demonstrate that DecAEvolve consistently outperforms state-of-the-art baselines in both discovery accuracy and search efficiency, while maintaining strong out-of-domain generalization. The success of smaller models through targeted test-time adaptation suggests promising directions for democratizing scientific discovery tools without requiring large, resource-intensive models. Future work could extend our simple decomposition mechanisms to more complex structures and explore better optimization strategies for the evolutionary process. The term-level feedback approach developed here may also prove valuable for broader program synthesis tasks requiring iterative refinement in the symbolic space of programs based on component-level understanding.

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [2] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023.
- [3] 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. *Nat.*, 625(7995):468–475, January 2024.

- [4] Alexander Novikov, Ngân Vũ, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco J. R. Ruiz, Abbas Mehrabian, M. Pawan Kumar, Abigail See, Swarat Chaudhuri, George Holland, Alex Davies, Sebastian Nowozin, Pushmeet Kohli, and Matej Balog. Alphaevolve: A coding agent for scientific and algorithmic discovery, 2025.
- [5] Anja Surina, Amin Mansouri, Lars Quaedvlieg, Amal Seddas, Maryna Viazovska, Emmanuel Abbe, and Caglar Gulcehre. Algorithm discovery with llms: Evolutionary search meets reinforcement learning. *arXiv preprint arXiv:2504.05108*, 2025.
- [6] Parshin Shojaee, Kazem Meidani, Shashank Gupta, Amir Barati Farimani, and Chandan K Reddy. Llm-sr: Scientific equation discovery via programming with large language models, 2025.
- [7] Nour Makke and Sanjay Chawla. Interpretable scientific discovery with symbolic regression: a review. *Artificial Intelligence Review*, 57(1):2, 2024.
- [8] John R. Koza. Genetic programming as a means for programming computers by natural selection. *Statistics and Computing*, 4(2):87–112, 1994.
- [9] William La Cava, Patryk Orzechowski, Bogdan Burlacu, Fabrício Olivetti de França, Marco Virgolin, Ying Jin, Michael Kommenda, and Jason H. Moore. Contemporary symbolic regression methods and their relative performance, 2021.
- [10] Marco Virgolin and Solon P. Pissis. Symbolic regression is np-hard, 2022.
- [11] Silviu-Marian Udrescu and Max Tegmark. Ai feynman: a physics-inspired method for symbolic regression, 2020.
- [12] J.-P. Bruneton. Enhancing symbolic regression with quality-diversity and physics-inspired constraints (qdsr). *arXiv preprint*, 2025.
- [13] Author Kamienny et al. End-to-end transformer-based equation generation for symbolic regression. In *NeurIPS*, 2022.
- [14] Parshin Shojaee, Kazem Meidani, Amir Barati Farimani, and Chandan Reddy. Transformer-based planning for symbolic regression. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 45907–45919. Curran Associates, Inc., 2023.
- [15] Kazem Meidani, Parshin Shojaee, Chandan K. Reddy, and Amir Barati Farimani. SNIP: Bridging mathematical symbolic and numeric realms with unified pre-training. In *The Twelfth International Conference on Learning Representations*, 2024.
- [16] Arya Grayeli, Atharva Sehgal, Omar Costilla-Reyes, Miles Cranmer, and Swarat Chaudhuri. Symbolic regression with a learned concept library. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- [17] Miles Cranmer. Interpretable machine learning for science with pysr and symbolicregression. jl. *arXiv preprint arXiv:2305.01582*, 2023.
- [18] Pingchuan Ma, Tsun-Hsuan Wang, Minghao Guo, Zhiqing Sun, Joshua B. Tenenbaum, Daniela Rus, Chuang Gan, and Wojciech Matusik. LLM and simulation as bilevel optimizers: A new paradigm to advance physical scientific discovery. In *Forty-first International Conference on Machine Learning*, 2024.
- [19] Parshin Shojaee, Ngoc-Hieu Nguyen, Kazem Meidani, Amir Barati Farimani, Khoa D Doan, and Chandan K. Reddy. LLM-SRBench: A new benchmark for scientific equation discovery with large language models. In *Forty-second International Conference on Machine Learning*, 2025.
- [20] Chandan K Reddy and Parshin Shojaee. Towards scientific discovery with generative ai: Progress, opportunities, and challenges. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 28601–28609, 2025.