

PROGRAMMATIC CONTEXT AUGMENTATION FOR LLM-BASED SYMBOLIC REGRESSION

Anonymous authors

Paper under double-blind review

ABSTRACT

Symbolic regression (SR), the task of discovering mathematical expressions that best describe a given dataset, remains a fundamental challenge in scientific discovery. Traditional approaches, primarily based on genetic algorithms and related evolutionary methods, have proven useful but suffer from scalability and expressivity limitations. Recently, large language model (LLM)-based evolutionary search methods have been introduced into SR and show promise. However, existing LLM-based approaches typically rely on scalar evaluation metrics, such as mean squared error, as the sole source of feedback during the search process, thereby overlooking the rich information embedded in the dataset. To address this limitation, we propose a novel LLM-based evolutionary search framework that incorporates programmatic context augmentation. By enabling code-based interactions with the dataset, our method can actively perform data analysis and extract informative signals, beyond aggregated evaluation scores. We evaluate our framework on advanced benchmarks, such as LLM-SRBench, and demonstrate superior efficiency and accuracy compared to strong baselines.

1 INTRODUCTION

Symbolic regression (SR) (Kronberger et al., 2024; Makke & Chawla, 2024) is a central task across many scientific domains. Given a dataset collected from observations or experiments, the goal of SR is to discover concise and interpretable mathematical equations that can best explain the underlying structure of the data. Unlike purely predictive models, SR offers human-readable equations that provide interpretability and scientific insight, making it a powerful tool for discovery.

The key challenges of SR lie in efficiently searching the vast combinatorial space of candidate equations. Traditional approaches focus on using genetic programming (Kronberger et al., 2024), which iteratively generate populations of equations, evaluate them using fitness scores derived from the data, and refine them through mutation and crossover. Other approaches include neural-guided search (Cranmer et al., 2020; Shah et al., 2020) and reinforcement learning (Petersen et al., 2019). While these methods have shown success in small-scale settings, they remain limited in efficiency and scalability, largely due to their reliance on random mutations and their inability to leverage prior experience. To improve upon these limitations, neural network-based SR methods trained on datasets of equation-data pairs have been proposed (Kamienny et al., 2022).

More recently, large language models (LLMs) have been introduced into SR through evolutionary search frameworks (Shojaee et al., 2024). At each iteration, a textual prompt that combines the problem description with insights from past iterations is provided to the LLM, which then generates new candidate equations. This paradigm leverages the domain knowledge encoded in pretrained LLMs and has demonstrated promising results. Methods such as concept-library learning (Grayeli et al., 2024) have further enhanced LLMs’ ability to reason symbolically in SR.

Despite these advances, current LLM-based SR methods remain constrained by their feedback mechanisms. Specifically, they rely almost exclusively on scalar evaluation metrics—most commonly mean squared error (MSE)—as the sole interaction with the dataset during the search process. That is, while candidate equations are evaluated against the dataset, the only signal returned to the model is a single numerical score. This stands in stark contrast to how human scientists approach the same task: rather than relying only on scalar evaluations, they actively conduct exploratory data

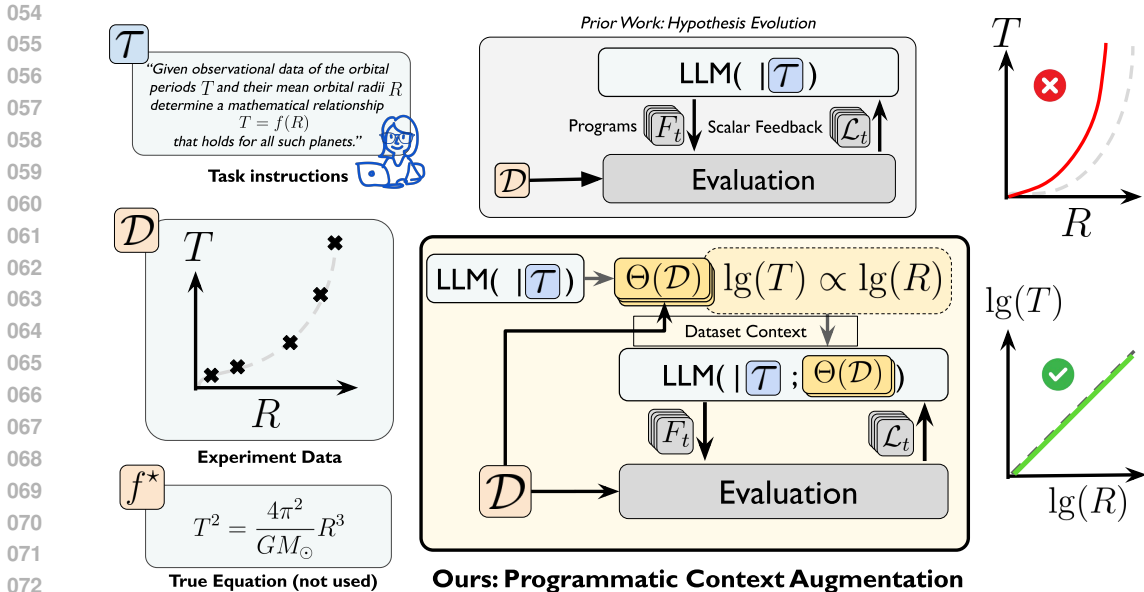


Figure 1: An overview of PROAUG, with τ being the task instructions, D being the dataset, Θ being the dataset context, F_t being the generated program and L_t being the scalar feedback. Prior work in LLM guided symbolic regression (Shojaee et al., 2024) (in gray) rely on simplistic scalar evaluation metrics, such as MSE as feedback, ignoring rich information contained in dataset. Instead, PROAUG enriches the model’s interaction with data by assigning the LLM a dual role—to generate data analysis code and extract informative signals for context augmentation. Considering the classical problem of deriving Kepler’s third law from planetary motion data as a example. The law states that the square of a planet’s orbital period (T) is proportional to the cube of its semi-major axis (R), i.e., $T^2 \propto R^3$. While standard method might struggle to identify this nonlinear relationship directly, PROAUG leverages the LLM’s ability to reason by data analysis. Specifically, the LLM generates code to apply a logarithmic transformation to both T and R , yielding $\log(T)$ and $\log(R)$. When plotted, these transformed variables exhibit a clear linear relationship: $\log(T) \approx \frac{3}{2} \log(R) + \text{constant}$. This linear trend immediately suggests a power-law relationship between the original variables. Our method demonstrates how data-driven statistics complement structural constraints in equation discovery.

analysis, examining relationships, distributions, and trends to guide the design of hypothesis equations. In this work, we bridge this gap by introducing **PROgrammatic context AUGmentation** (PROAUG). Our framework empowers the LLM not only to propose candidate equations but also to generate code for analyzing the underlying dataset. By executing such code, the model gains access to richer signals—such as statistical properties and variable correlations—that can be used to refine its evolutionary search.

Concretely, we propose a new LLM-based SR framework in which the model is tasked with two complementary objectives: (1) proposing potential equation candidates, and (2) generating data-analysis code that extracts informative signals from the dataset. To rigorously assess our method, we benchmark it on LLM-SRBench (Shojaee et al., 2025), a comprehensive and recently proposed evaluation suite that mitigates concerns about LLMs leveraging memorized knowledge. Across both zero-shot and supervised fine-tuning settings, our method consistently outperforms strong baselines in terms of both efficiency and accuracy.

Our contributions are summarized as follows:

- We introduce **programmatically context augmentation**, a novel mechanism that enables LLMs to actively interact with the dataset during evolutionary SR.
- We design a dual-task framework in which the LLM generates both candidate equations and code for data-driven analysis.
- We evaluate our approach on LLM-SRBench and demonstrate consistent improvements over state-of-the-art baselines across multiple settings. For instance, we see a 3x error reduction on LSR-Transform when using programmatically context augmentation with DeepSeek-V3.1 as the base model.

2 PRELIMINARIES

Symbolic Regression. Using the framework of Empirical Risk Minimization (ERM), the problem of symbolic regression can be formalized as follows (Kronberger et al., 2024; Vapnik, 1991).

Given a dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, where $\mathbf{x}_i \in \mathbb{R}^d$ denotes the input features and $y_i \in \mathbb{R}$ the scalar target, the objective of SR is to find a function $f^* \in \mathcal{F}$ that minimizes the empirical loss:

$$f^* = \arg \min_{f \in \mathcal{F}} \mathcal{L}(f),$$

where \mathcal{F} is the discrete hypothesis class of functions mapping \mathbb{R}^d to \mathbb{R} , and \mathcal{L} is a loss function measuring predictive accuracy. A common choice is the mean squared error (MSE):

$$\mathcal{L}(f) = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2.$$

Since the hypothesis class \mathcal{F} is constrained by Python code, one can also view this problem as an instance of program synthesis (Chaudhuri et al., 2021).

LLM evolutionary search for SR. Recent work has integrated large language models (LLMs) into the framework of evolutionary algorithms, yielding LLM agents that iteratively propose solutions to tasks such as coding or scientific discovery. We focus on their application to SR, as exemplified by LLM-SR (Shojaee et al., 2024), which serves as a main baseline in this work. The framework consists of three key components:

(i) *Hypothesis generation.* At each iteration, the LLM is prompted with a structured input that includes: task instructions, problem specifications (e.g., the physical meaning of variables), the optimization and evaluation function, and demonstrations of promising prior hypotheses and their improvement trajectories from an experience buffer. Conditioned on this prompt, the LLM generates a candidate equation in the form of a Python function skeleton with learnable parameters.

(ii) *Hypothesis optimization and assessment.* The free parameters of the generated equation skeletons are optimized using external solvers (e.g., NumPy with BFGS). The resulting hypothesis is then evaluated on the dataset to produce a fitness score. In LLM-SR, this score is defined as the negative MSE on the training set. However, we observed that this definition can lead to generalization mismatches between training and test performance. To mitigate this, we split the training set into a *train-train* (tr-tr) and *train-val* (tr-val) subset. Parameters are optimized on the train-train set, and the fitness score is computed on the train-val set using the negative normalized MSE.

(iii) *Experience management.* Inspired by mechanisms in evolutionary algorithms (e.g., island models (Cranmer, 2023; Romera-Paredes et al., 2024)), LLM-SR maintains an experience buffer of past hypotheses. Candidate hypotheses are sampled from this buffer according to a distribution weighted by fitness scores, and incorporated into future prompts as demonstrations. This balances diversity and efficiency during the evolutionary search process.

3 OVERALL FRAMEWORK

Although LLM-based methods have demonstrated effectiveness in symbolic regression tasks, they still exhibit key limitations. As discussed earlier, existing approaches primarily rely on scalar evaluation metrics derived from raw data samples, overlooking the rich statistical characteristics of the dataset that can be extracted through program. These characteristics can provide valuable inductive biases, helping to constrain the search space and guide the discovery of correct symbolic expressions. This naturally raises the question: *Can statistical information enhance LLM-based SR?* In this section, we address this question through controlled analysis and then introduce our proposed framework, PROAUG.

3.1 CAN STATISTICAL INFORMATION ENHANCE LLM-BASED SR?

Analysis Setup. To systematically examine the role of statistical information in LLM-based SR, we conduct a controlled case study on a representative instance, II.6.15b_1_0, from the LLM-

SRBench benchmark. The ground-truth expression for this case is: $A_{vec} = \frac{j \cdot m}{q \cdot \rho \cdot c_0}$, which corresponds to a physical relationship involving the vector potential, current density, charge carrier density, elementary charge, and mass of a charge carrier. To eliminate confounding effects from prior knowledge, we anonymize all variables and remove the original physical background in the experimental conditions here. Specifically, inputs and outputs are replaced with x_i and y respectively. We then compare the following three settings:

- **Zero-shot:** The model receives no data samples or additional information.
- **Few-shot:** The model is provided with random 12 raw data samples without any additional background or statistical information.
- **Statistical Hint:** In addition to the 12 raw samples, the model is provided with a structured dictionary that summarizes dataset-level statistics. This dictionary includes (i) basic statistical measures for each variable and the output—such as mean, standard deviation, minimum, and maximum, and (ii) feature-level correlations, such as R^2 values between the logarithm of the output and simple transformations of the inputs (e.g., log, exp, sin, cos).

Prompts used here for all settings are provided in Appendix E. We use DeepSeek-V3.1 (Liu et al., 2024) as the base model within the LLM-SR framework, conducting a total of 100 evolutionary steps for each experiment.

Results. Figure 2 shows the best normalized mean squared error (NMSE) score trajectory under all settings. Both the zero-shot and few-shot settings perform poorly. The 12 randomly sampled data points provided in the few-shot condition are insufficient for the model to infer the true underlying expression. As a result, the NMSE decreases only slowly over the 100 evolutionary iterations. In contrast, the inclusion of statistical hints significantly accelerates convergence, ultimately yielding an NMSE on the order of 10^{-13} . To understand the reason for this improvement, we look into the model’s behavior under the statistical hint setting. We find that the high R^2 value between $\log y$ and $\log x$ provided leads the model to infer that y is likely proportional to certain powers of the four input variables—suggesting a functional form of $y = x_0^a \cdot x_1^b \cdot x_2^c \cdot x_3^d$, where a, b, c, d are parameters can be optimized. This inferred form aligns perfectly with the ground-truth expression. These results demonstrate that even basic statistical information can provide substantial guidance in the LLM-SR process for the specific case. By constraining the search space and suggesting plausible functional forms, statistical cues can improve both the efficiency and the accuracy of symbolic regression. This mirrors how human scientists often practice exploratory data analysis before formulating hypotheses in the process of scientific discovery.

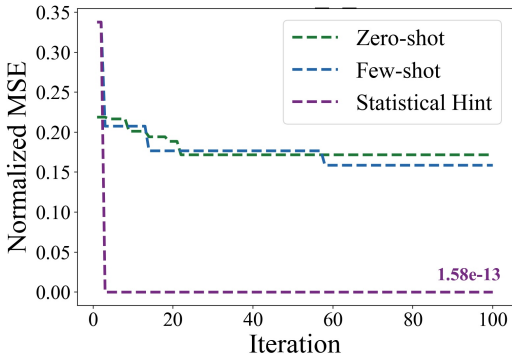


Figure 2: NMSE trajectories under three settings with varying amounts of background information.

Limitations. Despite its utility, the design of the statistical context here has certain limitations. The correlation metrics are manually designed and confined to a small predefined set of transformations (e.g., exp, sin, cos). Consequently, complex nonlinear or multivariate interactions may remain undetected. While effective for highlighting simple algebraic relationships, this approach may fail with more intricate structures.

3.2 PROGRAMMATIC CONTEXT AUGMENTATION (PROAUG)

To fully leverage statistical context, it is essential to move beyond manually engineered heuristics. In the next section, we introduce **Programmatic Context Augmentation (PROAUG)**, in which LLMs—guided by physical knowledge—automatically generate dataset-analysis code. This enables

Algorithm 1: Pseudocode for PROAUG.

```

216 Input: Pretrained LLM  $\mathcal{M}$ , dataset  $\mathcal{D}$ , number of iterations  $T$ 
217 Output: Best symbolic expression  $f^*$ 
218 for  $t = 0, \dots, T - 1$  do
219   // Programmatic context augmentation
220   Generate dataset-analysis prompt and sample analysis code
221   Execute code on  $\mathcal{D}$  to obtain informative signals
222   // Context construction
223   Construct enriched prompt using task instruction, problem specification, experience demonstrations,
224   and extracted signals
225   // Hypothesis generation
226   Sample candidate equation from  $\mathcal{M}$  given enriched prompt
227   // Optimization and evaluation
228   Optimize free parameters on  $\mathcal{D}_{\text{tr-tr}}$ 
229   Evaluate fitness score on  $\mathcal{D}_{\text{tr-val}}$ 
230   // Experience management
231   Update population buffer with candidate and score
232   Update best solution  $f^*$ 
233 end

```

the model to extract richer and more adaptive signals, preserving the benefits of inductive bias while broadening coverage.

The goal of PROAUG is to empower the LLM to actively interact with the dataset during evolutionary search by first proposing *data analysis programs*. These programs are executed to extract informative signals from the dataset, which are then incorporated into the prompt context. The enriched context prompt is subsequently used by the LLM to generate candidate symbolic expressions. In contrast, the original LLM-SR pipeline (see Section 2) bypasses this step: it directly prompts the LLM to generate hypotheses based only on task instructions, problem specifications, and past experience, without any dataset-level analysis.

Formally, PROAUG extends the LLM-SR pipeline by introducing an additional *dataset analysis phase* before hypothesis generation. The new pipeline proceeds as follows: 1. In the *data analysis phase*, The LLM is instructed with a structured prompt to generate Python code that extracts statistical or structural information from the dataset. 2. The generated code is executed, and the extracted signals (e.g., summary statistics, correlations, transformation relationships) are fed back as part of an enriched context prompt. 3. Using this enriched prompt, the LLM generates candidate equations within an evolutionary search framework similar as LLM-SR.

The structured prompt in the *data analysis phase* includes both a task description and a lightweight code template of basic analysis tools, which includes computations such as (i) descriptive statistics of each variable (mean, variance, range), (ii) linear correlations, and R^2 tests between transformed features (e.g., log, exp, sin, cos) or feature combinations and the target variable or the transformation of it.

A brief overview of PROAUG is shown in Algorithm 1, along with a pipeline illustration and an example of a PROAUG prompt, is provided in Figure 1 and Appendix E.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Datasets Benchmarking LLMs in symbolic regression is challenging due to concerns that pre-trained models may exploit memorized knowledge rather than demonstrating genuine reasoning and search capabilities. We evaluate PROAUG on LLM-SRBENCH (Shojaee et al., 2025), a recently proposed benchmark designed to mitigate this issue. LLM-SRBench consists of two primary subsets: **LSR-Transform** and **LSR-Synth**.

LSR-Transform includes 111 target functions adapted from Feynman equations (Udrescu & Tegmark, 2020), further modified through variable substitutions and symbolic rewrites to increase

Table 1: Comparison of different methods on LLM-SRBENCH, evaluated using normalized mean squared error (NMSE). The benchmark subsets contain 17 tasks from LSR-Transform, and 5, 3, 5, and 3 tasks from Chemistry, Biology, Physics, and Materials Science, respectively. Each entry reports the average NMSE across all tasks within the corresponding subset. The final column (Average) reports the mean NMSE over all LSR-Synth tasks (16 instances in total). Bold values indicate the best performance for each model, while underlined values denote the best overall performance across models.

Models	LSR-Transform	LSR-Syn				Average
		Chemistry	Biology	Physics	Material Science	
Qwen3-4B-Instruct-2507 (Yang et al., 2025)						
LLM-SR	0.258	<u>7.55e-6</u>	2.48	0.048	1.60e-6	0.480
Statistical Hint	0.225	0.032	0.095	0.22	0.012	0.099
PROAUG	0.145	4.73e-3	0.036	0.060	9.62e-6	0.027
Qwen3-8B (Yang et al., 2025)						
LLM-SR	0.199	0.012	0.041	0.077	2.05e-4	0.036
Statistical Hint	0.279	4.11e-3	0.019	<u>0.019</u>	0.034	0.017
PROAUG	0.148	3.10e-3	<u>1.08e-3</u>	0.021	1.67e-7	<u>7.75e-3</u>
Deepseek-V3.1 (Liu et al., 2024)						
LLM-SR	0.230	0.080	0.060	0.110	8.47e-7	0.071
Statistical Hint	0.188	0.022	1.112	0.034	0.333	0.288
PROAUG	<u>0.067</u>	1.02e-05	0.094	0.058	<u>1.82e-9</u>	0.036

structural complexity. Due to computational constraints, we randomly sample 15% of the tasks (17 instances) for evaluation. LSR-Synth consists of symbolic regression tasks that integrate canonical scientific terms with synthetically constructed expressions, spanning four domains—chemistry, physics, biology, and materials science. From this subset, we randomly select 16 tasks for testing. The complete list of tasks used in our experiments is provided in the Appendix B.

Baselines We adopt LLM-SR (Shojaee et al., 2024) as our primary baseline and evaluate three backbone models: Qwen3-4B-Instruct-2507 (Yang et al., 2025), Qwen3-8B (Yang et al., 2025), and DeepSeek-V3.1 (Liu et al., 2024), covering a range of model scales. In addition to LLM-SR, we also compare against a hand-crafted statistical hint baseline, which incorporates both the problem specifications (physical meaning of variables) as well as the statistical hint. All methods are evolved for 150 iterations per task, with 2 samples generated per prompt. To reduce the effect of randomness, each experiment is repeated three times with different random seed, and we report the mean results.

Evaluation Metrics Follow prior work, we primarily evaluate performance using the normalized mean squared error (NMSE). NMSE quantifies prediction error normalized by target variance:
$$\text{NMSE} = \frac{\sum_{i=1}^{N_{\text{test}}} (f(x_i) - y_i)^2}{\sum_{i=1}^{N_{\text{test}}} (y_i - \bar{y})^2},$$
 where \bar{y} denotes the mean of the true values. By emphasizing large deviations, NMSE is well-suited for assessing numeric precision in symbolic regression.

4.2 RESULTS

As shown in Table 1, our method achieves lower NMSE on most datasets, consistently outperforming the LLM-SR baseline. For instance, in the LSR-Transform subset with DeepSeek-V3.1 as the backbone, our method attains an NMSE of 0.067 compared to 0.23 for LLM-SR, which is a 3x reduction in error. Moreover, we observe that the performance advantage of our method over LLM-SR becomes more pronounced as the capability of the underlying model improves. In particular, with DeepSeek models, our approach yields significant improvements across all tasks. We attribute this trend to the nature of programmatic context augmentation, which depends on the model’s semantic understanding: weaker models are more likely to be overwhelmed by the extended context, while stronger models can effectively leverage it. Therefore, as base models continue to grow stronger

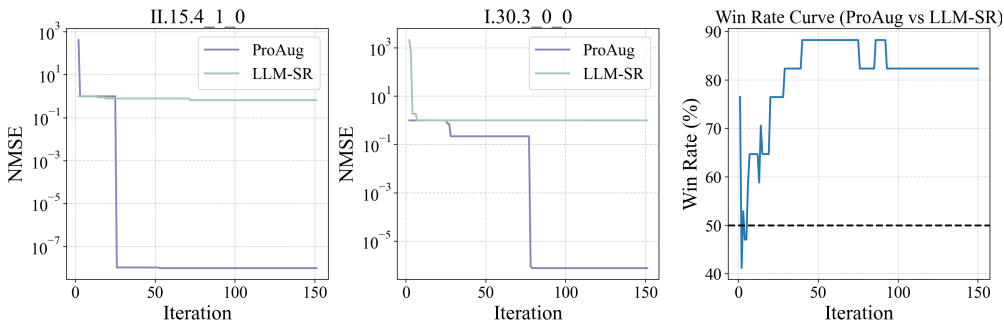


Figure 3: Analysis of convergence efficiency on LSR-Transform for DeepSeek-V3.1.

Table 2: SFT results using NMSE on LSR-Transform. best performance for each model in bold.

LSR-Transform	Qwen3-4B-Instruct-2507	Qwen3-8B
LLM-SR	0.258	0.199
LLM-SR-SFT	0.110	0.131
PROAUG	0.145	0.148
PROAUG-SFT	0.107	0.106

in semantic understanding and context management, we expect our method to become even more effective in the near future.

4.3 FINETUNING

Setting Furthermore, we conduct supervised fine-tuning (SFT) for the Qwen3-4B-Instruct-2507 and Qwen3-8B models using samples drawn from the remaining portion of LSR-Transform. The full training set consists of 77 problems, and we employ DeepSeek-V3.1 for data distillation. The distillation data generation procedure is divided into two stages. In the first stage, the model is provided with the ground-truth expression and prompted to generate a data analysis program based on background knowledge and the correct answer. The generated program is then executed, and if execution fails, the generation process is retried up to ten attempts. In the second stage, we concatenate the results of the data analysis and prompt the model to produce the final expression together with its reasoning steps. The SFT model was trained on two NVIDIA A800 GPUs with a batch size of 16, a learning rate of $5e-6$, and ZeRO Stage 2 for parallel acceleration.

We observe that the model occasionally generates non-executable code. To mitigate this issue, we augment the training data with additional examples where the data analysis code fails, but the model is still required to generate the correct expression. For comparison, we also distill outputs from DeepSeek-V3.1 for the LLM-SR method as an additional baseline. Finally, we omit SFT on LSR-Synth, as in this dataset subset, distinct problems share identical input prompts, leading to a one-to-many mapping issue that makes SFT ill-posed and counterproductive.

Results Table 2 reports the SFT results on LSR-Transform. We observe that supervised fine-tuning consistently improves performance for both LLM-SR and PROAUG, indicating that even standard fine-tuning provides clear benefits in this setting. For example, with Qwen3-4B-Instruct-2507, NMSE decreases from 0.258 to 0.110 for LLM-SR after SFT. Similarly, PROAUG also benefits from fine-tuning, with NMSE reduced from 0.145 to 0.107. Importantly, PROAUG-SFT achieves the lowest errors overall across both backbones, highlighting that our method remains effective and complementary to supervised fine-tuning.

4.4 METHOD BEHAVIOR ANALYSIS

4.4.1 CONVERGENCE EFFICIENCY

We evaluated the convergence efficiency of PROAUG relative to LLM-SR on the LSR-Transform subset from the DeepSeek-V3.1 model, with results shown in Figure 3. The two plots on the left

provide representative examples, illustrating that our method achieves a more rapid reduction in NMSE on train-val set during the iterative search process and ultimately converges to a lower error. The plot on the right displays the win rate of PROAUG over LLM-SR across 17 LSR-Transform test cases as a function of the number of iterations. We observe that PROAUG starts with a high win rate (above 50%) early in the process, indicating fast convergence of scores in the initial stages. As the iterations progress, the win rate continues to increase, eventually exceeding 80%. These results demonstrate that PROAUG maintains superior performance while also converges efficiently.

4.4.2 VARIANCE ANALYSIS

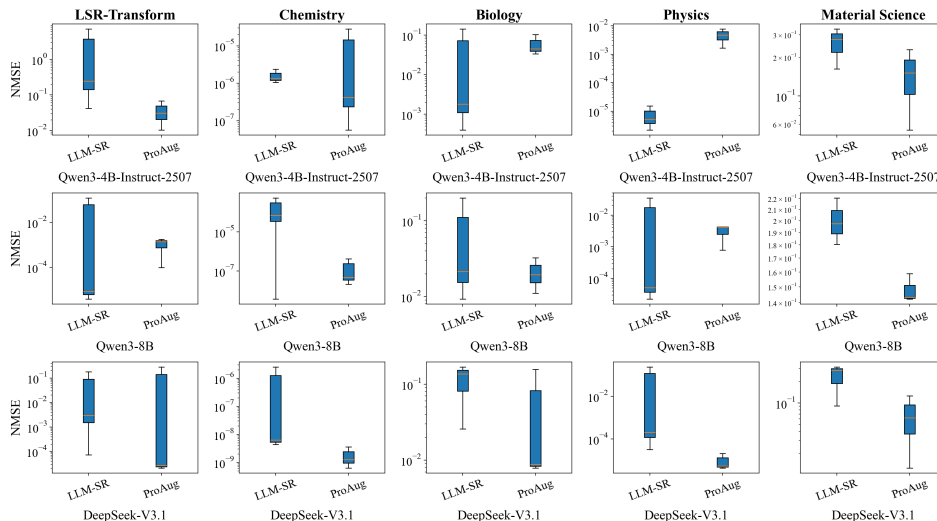


Figure 4: Boxplot comparison of tri-run variance for LLM-SR vs. PROAUG.

To rigorously evaluate the stability and reproducibility of LLM-SR and PROAUG, we conducted multiple independent runs under controlled hyperparameters and random seeds. Figure 4 presents boxplots summarizing the distribution of results, including the median and interquartile range (IQR) for each method.

Overall, PROAUG exhibits substantially lower variance across runs. LLM-SR shows high sensitivity to randomness, with outcomes occasionally differing by several orders of magnitude. For example, on the Chemistry task with Qwen3-8B, the NMSE varied from 10^{-4} to 10^{-9} across runs. Such findings highlight a critical gap in current LLM-based SR research, where multi-run evaluations are rarely reported. We therefore recommend that future studies adopt repeated trials to mitigate randomness and enhance reliability.

5 RELATED WORK

Symbolic Regression. Symbolic regression (SR) is core to scientific discovery. Interest began in the 1970s (Gerwin, 1974; Langley, 1977) and SR has recently become prominent in AI-for-science (Makke & Chawla, 2024; Merler et al., 2024; Romera-Paredes et al., 2024). SR methods follow three algorithmic themes:

Search based methods. SR is a challenging combinatorial optimization task. These methods search the space of possible equations while minimizing an objective function, relying on custom heuristics and parallelization to efficiently hypothesize and evaluate candidates (Cranmer, 2023; Petersen et al., 2019; Stephens, 2024; Udrescu & Tegmark, 2020). Notably, PySR (Cranmer, 2023) presents a multi-population evolutionary algorithm that has found practical utility in cosmology (Davis & Jin, 2023), international economics (Verstyuk & Douglas, 2022), and climate modeling (Grundner et al., 2024). While such approaches work well out-of-the-box, they are difficult to steer towards specific class of equations. PROAUG extends such approaches by providing two flexible steering mechanisms:

432 scientists can either manually steer search through natural language priors or such priors can be
 433 automatically extracted using dataset statistics.

434 *Learning based methods.* Another approach uses neural networks to accelerate SR (Tenachi et al.,
 435 2023; Biggio et al., 2021; Landajuela et al., 2022). These methods train networks on experimental
 436 data to either directly induce target expressions (Merler et al., 2024; Romera-Paredes et al., 2024;
 437 Meyerson et al., 2024; Biggio et al., 2021) or to accelerate the search for target expressions (Tenachi
 438 et al., 2023; Romera-Paredes et al., 2024; Zhang et al., 2025; Shah et al., 2020). However, neu-
 439 ral networks require retraining when datasets change significantly, necessitating scientists to mon-
 440 itor data distribution shifts and repeatedly retrain models – both computationally expensive tasks.
 441 PROAUG follows prior work in primarily being a data-driven SR algorithm, but avoids the com-
 442 putational overheads by replacing learned neural representations with programmatically extracted
 443 insights, greatly increasing practical utility. In particular, compared with RAG-SR (Zhang et al.,
 444 2025), PROAUG focuses on improving LLM-based SR methods that leverage the physical reasoning
 445 and world knowledge capabilities of large language models as agents (e.g., LLM-SR, LaSR). In con-
 446 trast, RAG-SR employs transformer-based neural networks without LLM-scale world knowledge or
 447 physical reasoning components, representing a fundamentally different methodological paradigm.
 448 Our methodology of leveraging programmatic interaction with the dataset to extract complex sta-
 449 tistical insights for enhancing reasoning and search is also fundamentally different from RAG-SR.
 450 RAG-SR reconciles with our motivation that scalar feedback can be limited, while complementarily
 451 work on different lines of methods.

452 *LLM based methods.* Recent works leverage large language models for symbolic regression
 453 (Romera-Paredes et al., 2024; Novikov et al., 2025; Shojaee et al., 2024; Grayeli et al., 2024; Ma
 454 et al., 2024). These methods rely on LLMs’ world knowledge to propose and mutate programs con-
 455 ditioned on natural language instructions and execution feedback. Specifically, Funsearch (Romera-
 456 Paredes et al., 2024) and AlphaEvolve (Novikov et al., 2025) demonstrate that LLMs as mutation
 457 operators within evolutionary algorithms can discover novel heuristics for long-standing problems
 458 in mathematics, hardware design, and algorithms. Many exciting works augment the evolutionary
 459 process with sketching (Shojaee et al., 2024), library learning (Grayeli et al., 2024), and tool use (Ma
 460 et al., 2024). PROAUG leverages pretrained LLMs similarly but treats the scientific dataset as a first-
 461 class citizen in the optimization process rather than merely a scoring mechanism. Consecutively,
 462 PROAUG’s primary contribution is *complementary* to that of these methods, and it is technically
 463 possible to employ programmatic context augmentation to enhance other SR algorithms.

464 **Scientific Discovery with Foundation Models.** Modern scientific breakthroughs increasingly de-
 465 mand both mastery of individual fields and trans-disciplinary insights (Gottweis et al., 2025; Jinek
 466 et al., 2012). As such, many LLM frameworks have emerged that assist scientists in various stages
 467 of the scientific process – from hypothesis falsification (Huang et al., 2025) and idea generation
 468 (Radensky et al., 2025) to data-driven discovery (Majumder et al., 2024), heuristic design (Romera-
 469 Paredes et al., 2024; Novikov et al., 2025), and multi-purpose systems spanning hypothesis gener-
 470 ation to experiment design (Gottweis et al., 2025; Lu et al., 2024). While these systems primarily
 471 interface with knowledge sources and scientists through natural language, PROAUG (like (Romera-
 472 Paredes et al., 2024; Novikov et al., 2025)) interfaces with these sources through code – leveraging
 473 the precise semantics and deterministic execution of code to rigorously analyze scientific data, in
 474 addition to using textual reasoning.

475 476 477 6 CONCLUSION

479 In this work, we address a key limitation of existing LLM-based symbolic regression methods:
 480 their reliance on simplistic scalar feedback such as mean squared error. We propose PROAUG,
 481 a novel framework that enriches the model’s interaction with data by assigning the LLM a dual
 482 role—generating both candidate equations and data-analysis code. This design fosters a more ac-
 483 tive, human-like discovery process that leverages rich statistical signals and variable relationships.
 484 Experimental results on the challenging LLM-SRBENCH benchmark highlight the clear advantages
 485 of our approach. We believe this work represents a step toward more intelligent and interactive AI
 systems with stronger capabilities for scientific reasoning and discovery.

486 **Limitations and future work** Despite its promising results, our approach is not without limita-
487 tions. The current framework, while leveraging richer statistical feedback than MSE-based methods,
488 still relies on a predefined set of data analysis operations. This may restrict its ability to discover
489 equations that depend on highly complex or non-standard data relationships not captured by our
490 initial code-generation template. Future work could explore enabling the LLM to propose entirely
491 novel analytical procedures.

492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

REFERENCES

- 540
541
542 Luca Biggio, Tommaso Bendinelli, Alexander Neitz, Aurelien Lucchi, and Giambattista Parascandolo. Neural symbolic regression that scales, 2021. URL <https://arxiv.org/abs/2106.06427>.
543
544
- 545 Swarat Chaudhuri, Kevin Ellis, Oleksandr Polozov, Rishabh Singh, Armando Solar-Lezama, Yisong Yue, et al. Neurosymbolic programming. *Foundations and Trends® in Programming Languages*, 7(3):158–243, 2021.
546
547
548
- 549 Miles Cranmer. Interpretable machine learning for science with pysr and symbolicregression. *jl. arXiv preprint arXiv:2305.01582*, 2023.
550
551
- 552 Miles Cranmer, Alvaro Sanchez Gonzalez, Peter Battaglia, Rui Xu, Kyle Cranmer, David Spergel, and Shirley Ho. Discovering symbolic models from deep learning with inductive biases. *Advances in neural information processing systems*, 33:17429–17442, 2020.
553
554
- 555 Benjamin L Davis and Zehao Jin. Discovery of a planar black hole mass scaling relation for spiral galaxies. *The Astrophysical Journal Letters*, 956(1):L22, 2023.
556
557
- 558 Donald Gerwin. Information processing, data inferences, and scientific generalization. *Behavioral Science*, 19(5):314–325, 1974.
559
560
- 561 Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, Khaled Saab, Dan Popovici, Jacob Blum, Fan Zhang, Katherine Chou, Avinatan Hassidim, Burak Gokturk, Amin Vahdat, Pushmeet Kohli, Yossi Matias, Andrew Carroll, Kavita Kulkarni, Nenad Tomasev, Yuan Guan, Vikram Dhillon, Eeshit Dhaval Vaishnav, Byron Lee, Tiago R D Costa, José R Penadés, Gary Peltz, Yunhan Xu, Annalisa Pawlosky, Alan Karthikesalingam, and Vivek Natarajan. Towards an ai co-scientist, 2025. URL <https://arxiv.org/abs/2502.18864>.
562
563
564
565
566
- 567 Arya Grayeli, Atharva Sehgal, Omar Costilla Reyes, Miles Cranmer, and Swarat Chaudhuri. Symbolic regression with a learned concept library. *Advances in Neural Information Processing Systems*, 37:44678–44709, 2024.
568
569
570
- 571 Arthur Grundner, Tom Beucler, Pierre Gentine, and Veronika Eyring. Data-driven equation discovery of a cloud cover parameterization. *Journal of Advances in Modeling Earth Systems*, 16(3):e2023MS003763, 2024.
572
573
- 574 Kexin Huang, Ying Jin, Ryan Li, Michael Y. Li, Emmanuel Candès, and Jure Leskovec. Automated hypothesis validation with agentic sequential falsifications, 2025. URL <https://arxiv.org/abs/2502.09858>.
575
576
577
- 578 Martin Jinek, Krzysztof Chylinski, Ines Fonfara, Michael Hauer, Jennifer A Doudna, and Emmanuelle Charpentier. A programmable dual-rna-guided dna endonuclease in adaptive bacterial immunity. *science*, 337(6096):816–821, 2012.
579
580
- 581 Pierre-Alexandre Kamienny, Stéphane d’Ascoli, Guillaume Lample, and François Charton. End-to-end symbolic regression with transformers. *Advances in Neural Information Processing Systems*, 35:10269–10281, 2022.
582
583
584
- 585 Gabriel Kronberger, Bogdan Burlacu, Michael Kommenda, Stephan M Winkler, and Michael Afenzeller. *Symbolic regression*. Chapman and Hall/CRC, 2024.
586
587
- 588 Mikel Landajuela, Chak Shing Lee, Jiachen Yang, Ruben Glatt, Claudio P Santiago, Ignacio Aravena, Terrell Mundhenk, Garrett Mulcahy, and Brenden K Petersen. A unified framework for deep symbolic regression. *Advances in Neural Information Processing Systems*, 35:33985–33998, 2022.
589
590
591
- 592 Pat Langley. Bacon: A production system that discovers empirical laws. In *International Joint Conference on Artificial Intelligence*, 1977. URL <https://api.semanticscholar.org/CorpusID:2320342>.
593

- 594 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
595 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint*
596 *arXiv:2412.19437*, 2024.
- 597 Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist:
598 Towards fully automated open-ended scientific discovery, 2024. URL <https://arxiv.org/abs/2408.06292>.
- 601 Pingchuan Ma, Tsun-Hsuan Wang, Minghao Guo, Zhiqing Sun, Joshua B. Tenenbaum, Daniela
602 Rus, Chuang Gan, and Wojciech Matusik. Llm and simulation as bilevel optimizers: A new
603 paradigm to advance physical scientific discovery, 2024. URL <https://arxiv.org/abs/2405.09783>.
- 605 Bodhisattwa Prasad Majumder, Harshit Surana, Dhruv Agarwal, Bhavana Dalvi Mishra, Abhi-
606 jeetsingh Meena, Aryan Prakhhar, Tirth Vora, Tushar Khot, Ashish Sabharwal, and Peter Clark.
607 Discoverybench: Towards data-driven discovery with large language models, 2024. URL
608 <https://arxiv.org/abs/2407.01725>.
- 609 Nour Makke and Sanjay Chawla. Interpretable scientific discovery with symbolic regression: a
610 review. *Artificial Intelligence Review*, 57(1):2, 2024.
- 612 Matteo Merler, Nicola Dainese, and Katsiaryna Haitsiukevich. In-context symbolic regression:
613 Leveraging language models for function discovery. *arXiv preprint arXiv:2404.19094*, 2024.
- 614 Elliot Meyerson, Mark J. Nelson, Herbie Bradley, Adam Gaier, Arash Moradi, Amy K. Hoover, and
615 Joel Lehman. Language model crossover: Variation through few-shot prompting, 2024. URL
616 <https://arxiv.org/abs/2302.12170>.
- 617 Alexander Novikov, Ngân Vũ, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt
618 Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco JR Ruiz, Abbas Mehrabian,
619 et al. Alphaevolve: A coding agent for scientific and algorithmic discovery. *arXiv preprint*
620 *arXiv:2506.13131*, 2025.
- 622 Brenden K Petersen, Mikel Landajuela, T Nathan Mundhenk, Claudio P Santiago, Soo K Kim, and
623 Joanne T Kim. Deep symbolic regression: Recovering mathematical expressions from data via
624 risk-seeking policy gradients. *arXiv preprint arXiv:1912.04871*, 2019.
- 625 Marissa Radensky, Simra Shahid, Raymond Fok, Pao Siangliulue, Tom Hope, and Daniel S. Weld.
626 Scideator: Human-llm scientific idea generation grounded in research-paper facet recombination,
627 2025. URL <https://arxiv.org/abs/2409.14634>.
- 628 Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Matej Balog,
629 M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang,
630 Omar Fawzi, et al. Mathematical discoveries from program search with large language models.
631 *Nature*, 625(7995):468–475, 2024.
- 633 Ameesh Shah, Eric Zhan, Jennifer Sun, Abhinav Verma, Yisong Yue, and Swarat Chaudhuri. Learn-
634 ing differentiable programs with admissible neural heuristics. *Advances in neural information*
635 *processing systems*, 33:4940–4952, 2020.
- 636 Parshin Shojaee, Kazem Meidani, Shashank Gupta, Amir Barati Farimani, and Chandan K Reddy.
637 Llm-sr: Scientific equation discovery via programming with large language models. *arXiv*
638 *preprint arXiv:2404.18400*, 2024.
- 639 Parshin Shojaee, Ngoc-Hieu Nguyen, Kazem Meidani, Amir Barati Farimani, Khoa D Doan, and
640 Chandan K Reddy. Llm-srbench: A new benchmark for scientific equation discovery with large
641 language models. *arXiv preprint arXiv:2504.10415*, 2025.
- 643 Trevor Stephens. gplearn: Genetic programming in python, with a scikit-learn inspired api, 2024.
644 URL <https://github.com/trevorstevens/gplearn>. Accessed: 2024-05-22.
- 645 Wassim Tenachi, Rodrigo Ibata, and Foivos I. Diakogiannis. Deep Symbolic Regression for Physics
646 Guided by Units Constraints: Toward the Automated Discovery of Physical Laws. *ApJ*, 959(2):
647 99, December 2023. doi: 10.3847/1538-4357/ad014c.

648 Silviu-Marian Udrescu and Max Tegmark. Ai feynman: A physics-inspired method for symbolic
649 regression. *Science advances*, 6(16):eaay2631, 2020.
650

651 Vladimir Vapnik. Principles of risk minimization for learning theory. *Advances in neural informa-*
652 *tion processing systems*, 4, 1991.

653 Sergiy Verstyuk and Michael R Douglas. Machine learning the gravity equation for international
654 trade. *Available at SSRN 4053795*, 2022.
655

656 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,
657 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*
658 *arXiv:2505.09388*, 2025.

659 Hengzhe Zhang, Qi Chen, Bing Xue, Wolfgang Banzhaf, and Mengjie Zhang. Rag-sr: Retrieval-
660 augmented generation for neural symbolic regression. In *International Conference on Learning*
661 *Representations*, 2025.
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A USE OF LLMs

In the course of this work, we made limited use of LLMs as auxiliary tools. Specifically, we used an LLM to assist in generating small utility Python scripts for tasks such as data preprocessing, formatting, and visualization. These scripts were only supportive in nature and did not contribute to the conceptual development, design of experiments, or the core scientific results of the paper. All key research ideas, experimental design, and writing of the main paper were carried out independently by the authors. The LLM was not used to generate or edit the scientific content of the manuscript, nor was it involved in producing research hypotheses or interpretations. The authors take full responsibility for the correctness, originality, and integrity of all content in the paper.

B PROBLEM DETAILS

Details of the problems we used in our experiments are shown in Table 3.

Table 3: Test data we use from LLM-SRBench.

Discipline	Name
Chemistry	CRK13, CRK12, CRK35, CRK31, CRK8
Biology	BPG16, BPG8, BPG9
Physics	PO12, PO37, PO40, PO24, PO8
Material Science	MatSci9, MatSci18, MatSci8
LSR-Transform	I.29.16_1_0, II.6.15b_1_0, II.11.27_1_0, I.34.1_2_0, I.30.3_0_0, II.24.17_1_1, III.15.27_2_0, I.12.2_3_0, III.10.19_2_1, I.37.4_0_1, II.11.17_4_0, II.6.15b_3_0, II.15.4_1_0, II.13.23_1_0, II.34.29b_3_0, I.11.19_2_0, I.32.5_1_1

Table 4: Results of $Acc_{0.1}$.

Models	LSR-Transform	LSR-Synth				
		Chemistry	Biology	Physics	Material Science	Average
<i>Qwen3-4B-Instruct-2507</i> (Yang et al., 2025)						
LLM-SR	59.9±4.2	72.5±0.8	44.2±26.3	82.6±15.5	98.0±1.3	75.1±7.3
PROAUG	65.0±3.2	71.3±3.8	36.4±1.8	58.0±10.8	99.2±0.3	65.8±4.6
<i>Qwen3-8B</i> (Yang et al., 2025)						
LLM-SR	68.0±1.4	72.1±2.3	62.2±23.7	63.8±8.9	92.9±7.1	71.6±3.6
PROAUG	57.5±0.9	73.2±14.4	39.0±12.2	65.3±11.3	98.7±1.5	69.1±4.5
<i>Deepseek-V3.1</i> (Liu et al., 2024)						
LLM-SR	70.1±4.1	65.2±10.5	68.0±13.7	56.3±5.5	99.5±0.7	69.4±7.0
PROAUG	78.3±5.2	69.5±1.7	68.8±12.8	70.4±8.5	99.9±0.1	75.3±5.1

C OTHER EXPERIMENT RESULTS

We provide results of Accuracy under error tolerance (ACC_τ) which measures numeric precision. Specifically ACC_τ is defined as the proportion of test points with error below a threshold τ : $ACC_\tau = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} \mathbf{1} \left(\left| \frac{f(x_i) - y_i}{y_i} \right| \leq \tau \right)$, where $\mathbf{1}(\cdot)$ is the indicator function. ACC_τ captures “good enough” performance under bounded error, aligning with real-world tolerance requirements. We show results in Table 4.

810 D EQUATION EXAMPLES

811 We present several representative equations discovered by PROAUG and demonstrate how LLM-
812 generated statistical information enhances algorithmic performance.

813 D.1 EQUATION EXAMPLES

814 The following examples illustrate equations discovered by PROAUG.

815 Example 1: I.32.5.1.1

```
816 Problem background: radiation from accelerating charges in classical
817 ↪ electrodynamics.
818 Ground truth:  $\sqrt{6} \cdot \sqrt{\pi} \cdot \sqrt{Pwr \cdot c^3 \cdot \epsilon_0} / q$ 
819 NMSE: 3.15e-14
820 Symbol meaning: [Target: 'the acceleration of the charged particle',
821 ↪ Feature: 'the power of an electromagnetic wave', 'the charge of a
822 ↪ particle', 'the electric constant or permittivity of free space',
823 ↪ 'the speed of light']
824
825 ... Discovered equation:
826 def equation(Pwr: np.ndarray, q: np.ndarray, epsilon: np.ndarray, c:
827 ↪ np.ndarray, params: np.ndarray) -> np.ndarray:
828     term1 = np.power(Pwr, params[0])
829     term2 = np.power(q, params[1])
830     term3 = np.power(epsilon, params[2])
831     term4 = np.power(c, params[3])
832     output = params[4] * term1 * term2 * term3 * term4 + params[5]
833     return output
834 ...
```

835 Example 2: I.37.4.0.1

```
836 Problem background: interference and intensity relations for two
837 ↪ coherent wave sources.
838 Ground truth:  $2 \cdot I_2 \cdot \cos(\delta)^2 + I_2 + I_1 + 2 \cdot \sqrt{I_2 \cdot (I_2 \cdot \cos(\delta)^2 + I_2 + I_1)} \cdot \cos(\delta)$ 
839 NMSE: 2.68e-14
840 Symbol meaning: [Target: 'the intensity of the first wave source',
841 ↪ Feature: 'the resultant intensity of two wave sources', 'the
842 ↪ intensity of the second wave source', 'the phase difference
843 ↪ between the two wave sources']
844
845 ... Discovered equation:
846 def equation(Int: np.ndarray, I2: np.ndarray, delta: np.ndarray,
847 ↪ params: np.ndarray) -> np.ndarray:
848     cos_term = np.cos(delta + params[2])
849     sqrt_I2 = np.sqrt(np.maximum(I2, 1e-8))
850     a = 1.0
851     b = params[1] * cos_term * sqrt_I2
852     c = params[0] * I2 - Int
853     discriminant = b**2 - 4*a*c
854     valid_mask = discriminant >= 0
855     sqrt_I1 = np.zeros_like(Int)
856     sqrt_I1[valid_mask] = (-b[valid_mask] +
857 ↪ np.sqrt(discriminant[valid_mask])) / (2*a)
858     I1_result = np.maximum(sqrt_I1**2 + params[3] + params[4] * I2 *
859 ↪ cos_term, 0)
860     return I1_result
861 ...
```

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

Example 3: III.15.27.2_0

```

Problem background: quantized wave modes in periodic physical
↪ systems.
Ground truth: 2*pi*alpha/(k*n)
NMSE: 3.54e-15
Symbol meaning: [Target: 'the diameter or a characteristic length',
↪ Feature: 'the spring constant or a proportionality constant', 'a
↪ material-dependent constant or a dimensionless parameter', 'the
↪ number of turns or a dimensionless quantity']
... Discovered equation:
def equation(k: np.ndarray, alpha: np.ndarray, n: np.ndarray, params:
↪ np.ndarray) -> np.ndarray:
    output = params[0] * alpha / (k * n)
    return output
...

```

These examples demonstrate PROAUG’s ability to discover accurate equations across diverse problems with low NMSE that closely match the ground truth.

D.2 LLM-GENERATED STATISTICAL INFORMATION

We provide a concrete example demonstrating how GPT-generated statistical information improves performance on MatSci9 from the Material Science category in LSR-Synth. The problem is about finding the equation that represents Stress (σ), given data on Strain (ϵ), and Temperature (T).

PROAUG generates code to extract statistical information including correlation and regression r^2 values for LLM-proposed features including $\epsilon * T$, $\epsilon ** 3$, $\epsilon ** 2$, $\epsilon - T$, etc.

Given statistical context, the LLM reasoned that

“From the dataset and the extracted information, we can observe several important trends and correlations. The output ‘ σ ’ (Stress) shows a strong correlation with both ‘ ϵ ’ (Strain) and ‘ T ’ (Temperature), with a correlation coefficient of 0.693 for both. The R^2 values also support this, indicating that both variables contribute significantly to predicting Stress. The strongest correlation with ‘ σ ’ is observed with ‘ $\epsilon_squared$ ’, followed by ‘ ϵ_cubed ’, and ‘ ϵ_minus_T ’. Additionally, the R^2 values for ‘ $\sigma_vs_epsilon_squared$ ’ and ‘ $\sigma_vs_T_squared$ ’ suggest that non-linear relationships may be more appropriate than linear ones. Given the high correlation with ‘ $\epsilon_squared$ ’ and the fact that Stress is often related to the square of Strain in physical models (such as in Hooke’s Law or stress-strain relationships), we can infer that the function might be a quadratic or cubic function of Strain, possibly combined with Temperature.”,

which leads to the discovered equation ‘ $output = params[0] * \epsilon + params[1] * T + params[2] * \epsilon ** 3 + params[3] * \epsilon ** 3 * T$ ’ after several search iterations.

The discovered equation, with a NMSE of $8.029e-14$, closely matches the ground truth, which is “ $4.255830478818137 * \epsilon ** 3 - 0.39029568781474044 * (T - 308.8646781296164) + \epsilon ** 3 * 8.63802815432207 * (T - 308.8646781296164)$ ”, demonstrating how statistical context guides the LLM toward physically meaningful functional forms.

918 E PROMPT DESIGN

919 E.1 PROAUG PROMPT DESIGN

920 Aside from the statistical information augmentation component, we adopt the same prompt struc-
 921 ture as the original LLM-SR/LLM-SRBench framework. The prompt template below illustrates the
 922 structure of the prompt used in PROAUG.

925 PROAUG Prompt Template

926 You are a helpful assistant tasked with discovering mathematical function structures for sci-
 927 entific data science tasks. Complete the 'equation' function below, considering the dataset
 928 information.

929 """
 930 Find the mathematical function skeleton that represents \$OUTPUT_VAR_DESC, given data on
 931 \$INPUT_VAR_DESC[0], ..., \$INPUT_VAR_DESC[N].
 932 """

```
933
934 import numpy as np
935
936 #Initialize parameters
937 MAX_NPARAMS = 10
938 params = [1.0]*MAX_NPARAMS
939
940 def equation_v0($INPUT_VAR[0], ..., $INPUT_VAR[N], params):
941     """
942     Args:
943     $INPUT_VAR[0]: A numpy array representing observations of
944     ↪ {$INPUT_VAR_DESC[0]}.
945     ...
946     $INPUT_VAR[N]: A numpy array representing observations of
947     ↪ {$INPUT_VAR_DESC[N]}.
948     params: Array of numeric constants or parameters to be
949     ↪ optimized
950
951     Return:
952     A numpy array representing {$OUTPUT_VAR_DESC} as the result
953     ↪ of applying the mathematical function to the inputs.
954     """
955     # Equation example 1 logic as function body
956     ...
957
958 def equation_v1($INPUT_VAR[0], ..., $INPUT_VAR[N], params):
959     """Improved version of `equation_v0`."""
```

960 You are to write Python code to extract and analyze relevant information from the dataset, with
 961 the goal of uncovering the underlying scientific function structure. The code will be executed
 962 in a sandbox environment that has direct access to the data. The output will be returned to
 963 assist your reasoning and support you in deriving the final answer.

964 You can do, but not limited to, investigating common statistics, use physical reasoning to
 965 propose candidate engineered features (inverse, cos/sin, log, exp, ratio, products, powers, etc)
 966 as well as compositions of these terms, and use tools including but not limited to correlation
 967 analysis, log-log fit, log-linear fit, multivariate linear fits, periodicity tests, and additive models
 968 in log space, etc, to identify whether the system follows additive, multiplicative, exponential,
 969 or power-law relationships.

```
970 def extract_dataset_info($OUTPUT_VAR, $INPUT_VAR[0], ...,
971 ↪ $INPUT_VAR[N]):
972     """
973     Args:
974     $OUTPUT_VAR[0]: A numpy array representing observations of
975     ↪ {$OUTPUT_VAR_DESC}.
```

```

972
973     $INPUT_VAR[0]: A numpy array representing observations of
974     ↪ {$INPUT_VAR_DESC}.
975     ...
976     $INPUT_VAR[N]: A numpy array representing observations of
977     ↪ {$INPUT_VAR_DESC[N]}.
978
979     Return:
980     """ A dict summarizing the extracted dataset information.
981     """
982     import numpy as np
983     from scipy.stats import linregress
984     def r2(x, y):
985         slope, intercept, r_value, _, _ = linregress(x, y)
986         return float(np.round(r_value ** 2, 3))
987
988     # Basic statistics
989     """ mean / std of certain variables """
990     info['stats'] = {}
991
992     # Feature transform
993     """ use physical reasoning to propose a variety of candidate
994     ↪ engineered features as well as compositions of these terms
995     ↪ for [$OUTPUT_VAR, $INPUT_VAR[0], ..., $INPUT_VAR[N]] """
996     features = {}
997
998     # Linear correlation with $OUTPUT_VAR
999     info['correlations_with_$OUTPUT_VAR'] = {}
1000
1001     # R2 values from regression fits
1002     """ regression fits between $OUTPUT_VAR or transforms of
1003     ↪ $OUTPUT_VAR and [$INPUT_VAR[0], ..., $INPUT_VAR[N]] and
1004     ↪ features of [$INPUT_VAR[0], ..., $INPUT_VAR[N]] """
1005
1006     info['regression_r2'] = {}
1007     return info
1008
1009     info = {'stats': {}, 'correlations_with_Int_0': {}, 'regression_r2':
1010     ↪ {}}
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

```

Given this template, the LLM first generates the statistical analysis code function `extract_dataset_info`, which is executed in a sandbox environment with direct access to the data. The output is then returned to the LLM to generate the final equation.

Computational cost: Compared to LLM-SR, The main additional overhead of PROAUG is that requires two LLM calls per iteration (one for statistical code generation, one for equation generation) versus one call for original LLM-SR. However, this can be viewed as a single tool-use process, where the LLM generates analysis tools (statistical code), obtains tool-call results (extracted statistics), and performs the final action (equation generation). As LLM inference techniques continue to advance and become more efficient, we expect this overhead to diminish further.

For comparison, the LLM-SR prompt template is shown as follows:

```

1016 LLM-SR Prompt Template
1017
1018 You are a helpful assistant tasked with discovering mathematical function structures for sci-
1019 entific data science tasks. Complete the 'equation' function below, considering the physical
1020 meaning and relationships of inputs.
1021 """
1022 Find the mathematical function skeleton that represents $OUTPUT_VAR_DESC, given data on
1023 $INPUT_VAR_DESC[0], ..., $INPUT_VAR_DESC[N].
1024 """
1025
1026 import numpy as np

```

```

1026
1027
1028 #Initialize parameters
1029 MAX_NPARAMS = 10
1030 params = [1.0]*MAX_NPARAMS
1031
1032 def equation_v0($INPUT_VAR[0], ..., $INPUT_VAR[N], params):
1033     """
1034     Args:
1035     $INPUT_VAR[0]: A numpy array representing observations of
1036     ↪ {$INPUT_VAR_DESC[0]}.
1037     ...
1038     $INPUT_VAR[N]: A numpy array representing observations of
1039     ↪ {$INPUT_VAR_DESC[N]}.
1040     params: Array of numeric constants or parameters to be
1041     ↪ optimized
1042
1043     Return:
1044     A numpy array representing {$OUTPUT_VAR_DESC} as the result
1045     ↪ of applying the mathematical function to the inputs.
1046     """
1047     # Equation example 1 logic as function body
1048     ...
1049
1050 def equation_v1($INPUT_VAR[0], ..., $INPUT_VAR[N], params):
1051     """Improved version of `equation_v0`."""
1052
1053
1054
1055
1056
1057
1058

```

As shown, PROAUG shares the same prompt structure as LLM-SR, aside from the statistical information augmentation component. Both LLM-SR and PROAUG use a linear regression model as the seed hypothesis at the beginning of the search process. For both LLM-SR and PROAUG, the prompts for different problem instances are generated by replacing the input and output variable names `$OUTPUT_VAR`, `$INPUT_VAR[0]`, ..., `$INPUT_VAR[N]`, and short descriptions `$OUTPUT_VAR_DESC`, `$INPUT_VAR_DESC[0]`, ..., `$INPUT_VAR_DESC[N]` in the template with instance-specific values. Finally, we provide a concrete example using a linear regression equation for instance II.6.15b.1.0 to illustrate what these variable names and descriptions look like.

```

1059 Equation example for II.6.15b.1.0
1060
1061 """
1062 Find the mathematical function skeleton that represents the dipole moment, given data on the
1063 electric field, the electric constant or permittivity of the medium, the angle between the dipole
1064 axis and the position vector, and the distance from the dipole to the point where the electric
1065 field is being measured.
1066 """
1067
1068 import numpy as np
1069
1070 #Initialize parameters
1071 MAX_NPARAMS = 10
1072 params = [1.0]*MAX_NPARAMS
1073
1074 def equation_v0(Ef, epsilon, theta, r, params):
1075     """
1076     Args:
1077     Ef: A numpy array representing observations of the electric
1078     ↪ field.
1079     epsilon: A numpy array representing observations of the
1080     ↪ electric constant or permittivity of the medium.
1081     theta: A numpy array representing observations of the angle
1082     ↪ between the dipole axis and the position vector.

```

```
1080
1081     r: A numpy array representing observations of the distance
1082     ↪ from the dipole to the point where the electric field is
1083     ↪ being measured.
1084     params: Array of numeric constants or parameters to be
1085     ↪ optimized
1086
1087     Return:
1088     A numpy array representing the dipole moment as the result of
1089     ↪ applying the mathematical function to the inputs.
1090     """
1091     output = params[0] * Ef + params[1] * epsilon + params[2] * theta
1092     ↪ + params[3] * r + params[4]
1093     return output
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
```

1134 E.2 PROMPTS USED IN SECTION 3.1
1135

1136 We present the Zero-shot, Few-shot, and Statistical hint prompts used in Section 3.1 for instance
1137 II.6.15b.1_0 as follows.
1138

1139 **Zero-Shot Prompt**

1140 You are a helpful assistant tasked with discovering mathematical function structures for sci-
1141 entific data science tasks. Complete the 'equation' function below, considering the dataset
1142 information.
1143

```
1144 """
1145 ### Your Task:
1146 FIRST, provide BRIEF reasoning inside a <thought>...</thought>
1147 ↪ block.
1148
1149 THEN, AFTER THE </thought> BLOCK, output your evolved mathematical
1150 ↪ function skeleton in Python. The function should begin with a
1151 ↪ 'def' line and follow standard Python formatting.
```

1151 Note: DO NOT use more than 10 params
1152 """

```
1153
1154
1155 import numpy as np
1156
1157 #Initialize parameters
1158 MAX_NPARAMS = 10
1159 params = [1.0]*MAX_NPARAMS
1160
1161 def equation_v0(x_0: np.ndarray, x_1: np.ndarray, x_2: np.ndarray,
1162 ↪ x_3: np.ndarray, params: np.ndarray) -> np.ndarray:
1163     """ Mathematical function
1164
1165     Args:
1166         x_0: A numpy array.
1167         x_1: A numpy array.
1168         x_2: A numpy array.
1169         x_3: A numpy array.
1170
1171         params: Array of numeric constants or parameters to be
1172 ↪ optimized
1173
1174     Return:
1175         A numpy array as the result of applying the mathematical
1176 ↪ function to the inputs.
1177     """
1178     output = params[0] * x_0 + params[1] * x_1 + params[2] * x_2 +
1179     ↪ params[3] * x_3 + params[4]
1180     return output
1181
1182 def equation_v1(x_0: np.ndarray, x_1: np.ndarray, x_2: np.ndarray,
1183 ↪ x_3: np.ndarray, params: np.ndarray) -> np.ndarray:
1184     """Improved version of `equation_v0`."""
```

1183 **Few-shot Prompt**

1184 You are a helpful assistant tasked with discovering mathematical function structures for sci-
1185 entific data science tasks. Complete the 'equation' function below, considering the dataset
1186 information.
1187

```

1188 """
1189 Find the mathematical function skeleton that represents y, given data
1190 ↪ on x_0, x_1, x_2, x_3.
1191
1192 ### 12 Random Samples (X, Y) (Sorted by Y from small to large):
1193 X[0] = [-32.193, 4.357, 4.774, 1.672], Y[0] = 2.588
1194 X[1] = [-8.799, 2.543, 2.522, 2.127], Y[1] = 2.918
1195 X[2] = [-45.200, 4.476, 3.700, 2.379], Y[2] = 6.494
1196 X[3] = [-64.941, 3.223, 3.793, 1.486], Y[3] = 7.892
1197 X[4] = [-46.401, 3.874, 4.408, 4.287], Y[4] = 11.650
1198 X[5] = [-76.242, 2.469, 2.298, 1.248], Y[5] = 16.766
1199 X[6] = [-84.225, 4.124, 1.593, 1.798], Y[6] = 23.048
1200 X[7] = [-72.211, 3.124, 3.184, 3.350], Y[7] = 24.319
1201 X[8] = [-43.300, 1.639, 4.213, 4.297], Y[8] = 26.954
1202 X[9] = [-41.009, 1.516, 3.111, 3.784], Y[9] = 32.907
1203 X[10] = [-83.370, 3.689, 2.829, 4.502], Y[10] = 35.962
1204 X[11] = [-43.438, 1.281, 2.523, 3.718], Y[11] = 49.963
1205
1206 ### Your Task:
1207 FIRST, provide BRIEF reasoning inside a <thought>...</thought>
1208 ↪ block.
1209
1210 THEN, AFTER THE </thought> BLOCK, output your evolved mathematical
1211 ↪ function skeleton in Python. The function should begin with a
1212 ↪ 'def' line and follow standard Python formatting.
1213
1214 Note: DO NOT use more than 10 params
1215 """
1216
1217 import numpy as np
1218
1219 #Initialize parameters
1220 MAX_NPARAMS = 10
1221 params = [1.0]*MAX_NPARAMS
1222
1223 def equation_v0(x_0: np.ndarray, x_1: np.ndarray, x_2: np.ndarray,
1224 ↪ x_3: np.ndarray, params: np.ndarray) -> np.ndarray:
1225     """ Mathematical function
1226
1227     Args:
1228         x_0: A numpy array.
1229         x_1: A numpy array.
1230         x_2: A numpy array.
1231         x_3: A numpy array.
1232
1233         params: Array of numeric constants or parameters to be
1234 ↪ optimized
1235
1236     Return:
1237         A numpy array as the result of applying the mathematical
1238 ↪ function to the inputs.
1239     """
1240     output = params[0] * x_0 + params[1] * x_1 + params[2] * x_2 +
1241 ↪ params[3] * x_3 + params[4]
1242     return output
1243
1244 def equation_v1(x_0: np.ndarray, x_1: np.ndarray, x_2: np.ndarray,
1245 ↪ x_3: np.ndarray, params: np.ndarray) -> np.ndarray:
1246     """Improved version of `equation_v0`."""
1247

```

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

Statistical Hint Prompt

You are a helpful assistant tasked with discovering mathematical function structures for scientific data science tasks. Complete the 'equation' function below, considering the dataset information.

```

"""
Find the mathematical function skeleton that represents y, given data
↳ on x_0, x_1, x_2, x_3.
The information of (X, Y) dataset including random sample points,
↳ smoothness estimation and simple basis fit scores are as follows:
### 12 Random Samples (X, Y) (Sorted by Y from small to large):
X[0] = [-32.193, 4.357, 4.774, 1.672], Y[0] = 2.588
X[1] = [-8.799, 2.543, 2.522, 2.127], Y[1] = 2.918
X[2] = [-45.200, 4.476, 3.700, 2.379], Y[2] = 6.494
X[3] = [-64.941, 3.223, 3.793, 1.486], Y[3] = 7.892
X[4] = [-46.401, 3.874, 4.408, 4.287], Y[4] = 11.650
X[5] = [-76.242, 2.469, 2.298, 1.248], Y[5] = 16.766
X[6] = [-84.225, 4.124, 1.593, 1.798], Y[6] = 23.048
X[7] = [-72.211, 3.124, 3.184, 3.350], Y[7] = 24.319
X[8] = [-43.300, 1.639, 4.213, 4.297], Y[8] = 26.954
X[9] = [-41.009, 1.516, 3.111, 3.784], Y[9] = 32.907
X[10] = [-83.370, 3.689, 2.829, 4.502], Y[10] = 35.962
X[11] = [-43.438, 1.281, 2.523, 3.718], Y[11] = 49.963
Statistics: {'mean_Y': 21.331, 'std_Y': 24.473, 'min_Y': 0.03,
↳ 'max_Y': 335.661, 'mean_X_0': -44.025, 'std_X_0': 25.099,
↳ 'min_X_0': -87.583, 'max_X_0': -0.436, 'mean_X_1': 3.011,
↳ 'std_X_1': 1.151, 'min_X_1': 1.0, 'max_X_1': 5.0, 'mean_X_2':
↳ 2.995, 'std_X_2': 1.156, 'min_X_2': 1.0, 'max_X_2': 5.0,
↳ 'mean_X_3': 3.007, 'std_X_3': 1.156, 'min_X_3': 1.0, 'max_X_3':
↳ 5.0, 'r2_Y_X_0': 0.245, 'r2_Y_X_1': 0.151, 'r2_Y_X_2': 0.15,
↳ 'r2_Y_X_3': 0.111, 'r2_log(Y)_log(X_0)': 0.595,
↳ 'r2_log(Y)_log(X_1)': 0.133, 'r2_log(Y)_log(X_2)': 0.133,
↳ 'r2_log(Y)_log(X_3)': 0.134, 'r2_log(Y)_log(sin(x_0))': 0.0,
↳ 'r2_log(Y)_log(cos(x_0))': 0.0, 'r2_log(Y)_log(exp(x_0))': 0.433,
↳ 'r2_log(Y)_log(sqrt(x_0))': 0.595, 'r2_log(Y)_log(sin(x_1))':
↳ 0.003, 'r2_log(Y)_log(cos(x_1))': 0.001,
↳ 'r2_log(Y)_log(exp(x_1))': 0.128, 'r2_log(Y)_log(sqrt(x_1))':
↳ 0.133, 'r2_log(Y)_log(sin(x_2))': 0.004,
↳ 'r2_log(Y)_log(cos(x_2))': 0.001, 'r2_log(Y)_log(exp(x_2))':
↳ 0.128, 'r2_log(Y)_log(sqrt(x_2))': 0.133,
↳ 'r2_log(Y)_log(sin(x_3))': 0.003, 'r2_log(Y)_log(cos(x_3))':
↳ 0.001, 'r2_log(Y)_log(exp(x_3))': 0.129,
↳ 'r2_log(Y)_log(sqrt(x_3))': 0.134}

### Your Task:
Based on the statistics above, analyze and evolve the mathematical
↳ function skeleton.

IMPORTANT:
FIRST, provide **BRIEF** reasoning inside a <thought>...</thought>
↳ block about how you interpret the data and how it guides your
↳ function structure decisions.

THEN, AFTER THE </thought> BLOCK, output your evolved mathematical
↳ function skeleton in Python. The function should begin with a
↳ 'def' line and follow standard Python formatting.

```

```
1296
1297 Note: DO NOT use more than 10 params
1298 """
1299
1300
1301 import numpy as np
1302
1303 #Initialize parameters
1304 MAX_NPARAMS = 10
1305 params = [1.0]*MAX_NPARAMS
1306
1307 def equation_v0(x_0: np.ndarray, x_1: np.ndarray, x_2: np.ndarray,
1308 ↪ x_3: np.ndarray, params: np.ndarray) -> np.ndarray:
1309     """ Mathematical function
1310
1311     Args:
1312         x_0: A numpy array.
1313         x_1: A numpy array.
1314         x_2: A numpy array.
1315         x_3: A numpy array.
1316
1317         params: Array of numeric constants or parameters to be
1318         ↪ optimized
1319
1320     Return:
1321         A numpy array as the result of applying the mathematical
1322         ↪ function to the inputs.
1323     """
1324     output = params[0] * x_0 + params[1] * x_1 + params[2] * x_2 +
1325     ↪ params[3] * x_3 + params[4]
1326     return output
1327
1328 def equation_v1(x_0: np.ndarray, x_1: np.ndarray, x_2: np.ndarray,
1329 ↪ x_3: np.ndarray, params: np.ndarray) -> np.ndarray:
1330     """Improved version of `equation_v0`."""
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
```