
README: Rapid Equation Discovery with Multimodal Encoders

Gregory Kang Ruey Lau^{*12} Yue Ran Kang^{*1} Zi-Yu Khoo¹ Apivich Hemachandra¹ Ruth Wan Theng Chew¹
Bryan Kian Hsiang Low¹

Abstract

Discovering scientific laws or interpretable symbolic equations from data rapidly is important in many setting, such as decision-making in time-sensitive high-stake scenarios or applications involving interactive or iterative experimentation such as in scientific or machine learning workflows. However, existing methods, generally known as symbolic regression (SR), typically require long computational time to achieve good performance and have to run from scratch for each dataset. Recent methods that use pre-training SR foundation models for faster inference also suffer from performance limitations and require large training datasets. In this work, we propose README, a framework for rapid equation discovery that can generate performant, interpretable equations from limited, noisy data in just a few seconds, and requires significantly less training data compared to past SR foundation model approaches. We achieve this by being the first to (1) work with image representations of datasets to efficiently capture their key properties, (2) combine the capabilities of open-sourced pre-trained text and image encoders to produce an informative SR embedding space, and (3) develop a novel Grey Wolf Optimizer with Bayesian Optimization (GWBO) algorithm to rapidly optimize for the best symbolic expression within seconds. We empirically show that README outperforms benchmarks on a wide range of realistic datasets, including real experimental data from various domains and noisy video-extracted dynamics.

1. Introduction

In many scientific and industrial settings, obtaining interpretable symbolic expressions that describe systems accurately is a critical objective. For example, symbolic representation of physical phenomena in areas such as climate science (Grundner et al., 2024; AL NAJAR et al., 2023), material science (Wang et al., 2022; 2019), and robotics (Zhang & Chen, 2023; Mor, 2011) are important in building scientific understanding, and interpretable symbolic expressions describing industrial processes and systems can help in high-stakes decision-making scenarios (Rudin, 2019) and applications in aerospace engineering (Brunton et al., 2021), electrical systems (Andelić et al., 2024) and healthcare (Wahlquist et al., 2024; Fitzsimmons & Moscato, 2018), where verifiability and human oversight are often required.

Symbolic regression (SR) methods aim to achieve automated discovery of the symbolic expressions that best approximate a given dataset, which is a challenging problem given the large search space of possible expressions (Virgolin & Pissis, 2022). However, while existing methods such as genetic programming-based algorithms (Makke & Chawla, 2024) can generate good approximations, they are typically computationally intensive and slow to converge (Biggio et al., 2021), and suffer from high sensitivity to the choice of hyperparameters and the right basis functions (Petersen et al., 2020). To address these issues, some works have proposed pre-training transformer-based SR models on large corpora of data, so as to amortize computational cost and enable faster inference (Valipour et al., 2021; Kamienny et al., 2022). These include approaches using CLIP-based (Radford et al., 2021a) multi-modal architectures trained on symbolic expressions and numerical data that could be used for candidate generation together with genetic programming-based SR methods (Meidani et al., 2023; Liu et al., 2023; Shojaee et al., 2024). However, these works require large datasets and computational time to train the models from scratch.

Importantly, most of the past works have not emphasized low-latency requirement settings where the time constraint for accurate symbolic expression is *in seconds*, not minutes or hours. While there are a few methods for fast SR (McConaghy, 2011), their performance tend to degrade for more

^{*}Equal contribution ¹Department of Computer Science, National University of Singapore ²CNRS@CREATE, 1 Create Way, #08-01 Create Tower, Singapore 138602. Correspondence to: <{greglau, lowkh}@comp.nus.edu.sg, kyueran@u.nus.edu>.

realistic, noisy data. As a result, SR remains challenging for use in interactive, real-time or iterative scenarios, potentially limiting its utility in applications such as adaptive scientific experimentation and close to real-time decision making in high-stakes environment.

In this work, we have identified three insights to achieve performant rapid equation discovery. First, rather than working with raw numerical data, *image representations of numerical data* can aid SR. Humans use plots to quickly extract key trends and identify candidate equations. A similar approach might be adapted for multi-modal large language models. Images, even complex human-uninterpretable plots, can efficiently summarize mathematical trends with many variables while remaining readable to well-trained models, enabling better candidate equations generation.

Second, *existing pre-trained image and text encoders/models can be leveraged to efficiently build foundation models for SR*, given plots and symbolic equations. Building on rapidly improving open-sourced image and text encoders rather than separately training SR-specific models from scratch, can produce better models with less data and computational resources. These encoders may have been pre-trained to extract relevant features that are transferable to the SR task (e.g., shape features in the plots for the image encoder or relationships among math operations for the text encoder), and hence only require fine-tuning with a small amount of data to become effective.

Third, to rapidly optimize for the symbolic expression that best approximates a dataset with the desired properties (e.g., complexity), we consider whether *query-efficient approaches such as Bayesian Optimization (BO) could be used to significantly speed up* the search process, when combined with SR foundation model approaches. BO methods (Garnett, 2022) allow for optimization while reducing calls to expensive fit procedures, which is a natural combination with population-based algorithms such as Grey Wolf Optimizer (GWO) (Mirjalili et al., 2014) to enable rapid equation discovery.

Combining insights from these analysis, we propose README (Rapid Equation Discovery with Multimodal Encoders), a framework for SR that uses (1) an informative, compressed image representation of numerical data (Section 3.1); (2) an efficiently-trained transformer-based model built on top of pre-trained image and text encoders ($\sim 60\times$ less training data compared to past works) (Section 3.2); and (3) a novel combination of BO and GWO for a rapid, effective optimization process (Section 3.3), to (4) achieve state-of-the-art and robust SR results for challenging settings with realistic, noisy settings and tight time constraints ($\leq 10s$) (Section 4).

2. Problem formulation

SR inference phase. Consider a target system that is governed by an underlying equation $y = f(x)$, where $y \in \mathbb{R}$, $x \in \mathbb{R}^n$, and $f(x)$ is a function that can be symbolically expressed as a composition of math operators. Given an inference dataset \mathcal{D} consisting of a set of noisy m observations $\{(x_i, \tilde{y}_i)\}_{i=1}^m$ where $\tilde{y}_i = y_i + \epsilon_i$, the SR task is to obtain a symbolic expression for the underlying function $f(x)$ that is the most accurate while prioritizing parsimonious expressions. Specifically, our accuracy goal is to find a symbolic expression $G^* \in \mathbb{G}$, where \mathbb{G} is the space of all valid symbolic expressions consisting of symbolic representations of input variables and math operators for the task under consideration, that represents a function $g(x)$ with the maximum R^2 value¹ over a test dataset \mathcal{D}_t generated from the same underlying phenomenon as \mathcal{D} .

In practice, during the inference phase we aim to achieve the best SR expression subjected to two additional desiderata. First, we prefer equations that are more parsimonious (i.e., a symbolic representation G that is less complex as evaluated by the number of nodes in its expression tree, details in Section 4.3.1) as they tend to be more interpretable, though there is typically a trade-off between the achievable accuracy and parsimony of the symbolic representation G . Hence, we will evaluate methods on the parsimony of their proposed expressions, and examine accuracy-parsimony Pareto plots to analyze how well the methods balance this trade-off. Second, the SR methods should also have low *inference runtime*, as many practical scenarios may have strict time budgets. Hence, we evaluate these methods based on fixed, short time budgets (e.g., 10 – 30s) in our experiments (Section 4).

Training phase for foundation models. We consider the realistic setting where we can generate or have access to synthetic training datasets independently generated from any inference training or test data, i.e., a set $\{\mathcal{D}_i, F_i\}_i$ where F_i are ground truth symbolic representations of the dataset \mathcal{D}_i , which we can use to train SR foundation models. Given the cost of high-quality data generation, an additional desiderata is for SR methods involving the training of foundation models to be as data-efficient as possible. Hence, we also evaluate the amount of training data needed for foundation-based SR methods to achieve good inference data (Table 1).

3. Method Overview

README framework consists of three key components:

1. **Data processing.** [Section 3.1] For both inference and

¹Note that G^* is not unique if only the accuracy desiderata is considered (e.g., superfluous terms could be added to any expression to represent the same function $g(x)$), but the parsimony desiderata will mitigate this.

training, we have a data processing step \mathcal{P} that converts each datasets \mathcal{D}_i from raw numerical data to a single image I , i.e., $\mathcal{P}(\mathcal{D}_i) \rightarrow I_i$. As explained later, the image representation has several key benefits over raw numerical data.

2. **README model architecture.** [Section 3.2] The processed data will then pass through the README model, which consists of a pair of pre-trained image \mathcal{I} and text \mathcal{T} encoders as well as a text decoder \mathcal{W} that has been fine-tuned by a set of labeled training data $\{\mathcal{D}_i, F_i\}_i$ during the training phase. The README training process combines the feature extraction capabilities of the image encoder with the mathematical knowledge embedded in a pre-trained math text encoder, to obtain an informative embedding space \mathbb{S} that the image encoder maps datasets to (i.e., $I(\mathcal{D}_i) \rightarrow s_i$, where $s_i \in \mathbb{S}$), for inference as described below.
3. **Inference optimization.** [Section 3.3] For a given dataset \mathcal{D} during inference, we will use the trained README model to generate an initial candidate set, followed by our README optimization process to search for the best point $s^* \in \mathbb{S}$ that can be decoded to obtain the best symbolic expression $\mathcal{W}(s^*) \rightarrow G^*$. The base optimizer in the process is the Grey Wolf Optimizer (GWO), though for ultra-rapid scenarios ($\leq 10s$) we employ our novel Grey Wolf Optimizer with Bayesian Optimization (GWBO) method.

3.1. Data processing: Working with images

Unlike existing SR methods, README works by first converting raw numerical data to images. This is inspired by humans’ capabilities to more rapidly infer patterns and guess candidate symbolic expression skeletons (i.e., expression forms without specific numerical constants) for data by visualizing it, rather than just going through raw numerical data. For example, the oscillatory curves of a 1D sinusoidal function are immediately recognizable when visualized. For humans, visualizing and interpreting the plots quickly become very challenging for higher dimensional datasets though.

However, the growing capabilities of Multi-modal Large Language Models (MLLMs) suggests that their image encoders may have powerful feature extraction capabilities developed from large-scale training on diverse image datasets that may also be useful in capturing relevant patterns from data plots. If so, it may still be viable to use images to summarize relevant information from high-dimensional data, and use them for SR. Such images may even not appear human-interpretable, but could possibly be effectively used by fine-tuned image encoders and customized decoders. Hence, in README, for both model training and inference, we map every dataset \mathcal{D}_i to a corresponding image plot I_i through a standardized data processing step.

To demonstrate this, we propose starting with the most basic plotting approach: in a single graph, we generate and overlay line plots for each dimension of the data. Figure 1 shows an example plot. While the plots may not seem directly interpretable by humans and may also not uniquely represent a single symbolic expression (i.e., several expressions may correspond to the same plot), the general shape and features of the aggregated line plots, including information such as the axis magnitudes, provide sufficient details to significantly reduce the candidate search space and inform the optimization process for SR. This is similar to how humans can guess the expression skeleton but not necessarily the exact expression with its constants. In the README framework, the image is used only to narrow the search space for a more efficient optimization process to perform SR, as we will elaborate in section Section 3.2 and Section 3.3.

Furthermore, this approach also helps to standardize the input format (i.e. 1 image) across datasets that can consist of a wide range of dimensions and number of datapoints, unlike for numerical data where variations there will be variations in format and size. Such variation pose major challenges to other works (Meidani et al., 2023), leading to the lack of generalization of SR foundation models to datasets with larger sizes or dimensionality compared to the training dataset.

3.2. README model architecture and training process

A key innovation in our README model architecture component lies in our combination of the feature recognition capabilities of pre-trained image encoders with the mathematical knowledge contained in text encoders to generate an informative embedding space for SR. Note that the naive approach of directly using MLLMs for SR do not perform well (see Appendix F), hence past works (Kamienny et al., 2022; Meidani et al., 2023), have largely resorted to training transformers from scratch. Our approach allow us to obtain significantly better performance in SR with less training data. The README model architecture is adapted from the basic CLIP MLLM architecture (Radford et al., 2021b):

1. **Image encoder.** The image encoder $\mathcal{I} : \mathbb{I} \rightarrow \mathbb{S}_I$ maps each image plot $I_i = \mathcal{P}(\mathcal{D}_i) \in \mathbb{I}$ to its embedding vector representation $s_i \in \mathbb{S}$. Any general-purpose pre-trained encoder can be used, such as open-sourced ViT models (Dosovitskiy et al., 2021) which are trained on diverse image data. These models have powerful image feature recognition capabilities (e.g., earlier model layers), and although their original embedding space \mathbb{S}_I would not have the right structure for our SR tasks (see Section 4.2), they could be efficiently fine-tuned on our type of images from the data generation step.
2. **Text encoder.** The text encoder $\mathcal{T} : \mathbb{G} \rightarrow \mathbb{S}_T$ maps

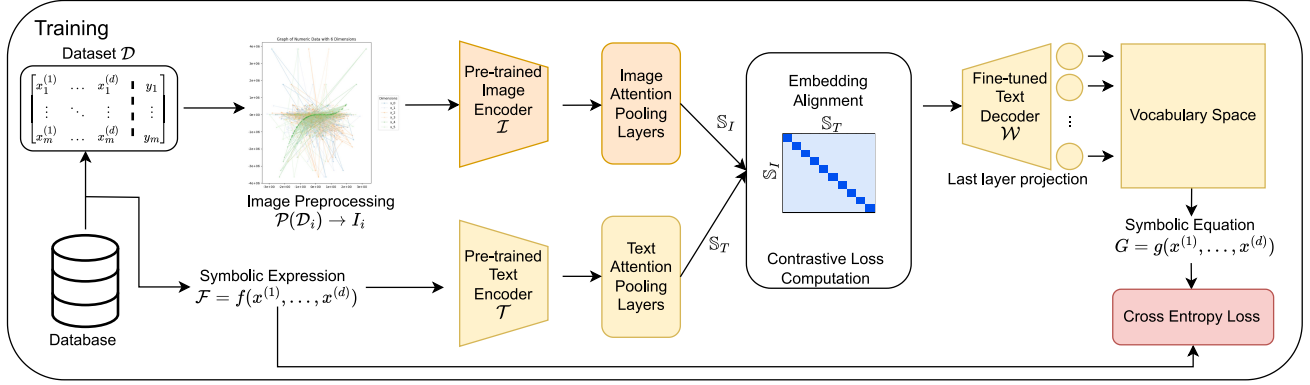


Figure 1. An overview of the README architecture

symbolic equations to its continuous embedding space \mathbb{S}_T . Crucially, we propose to use a text encoder pre-trained on math such as MathBERT (Shen et al., 2021) as it would contain inductive biases regarding symbolic expressions and have a relevant, well-structured embedding space \mathbb{S}_T for SR that can guide the training of \mathbb{S}_I for SR.

- Aligned embedding space.** The embedding spaces \mathbb{S}_I and \mathbb{S}_T are then aligned through joint contrastive learning, similar to the approach in CLIP (Radford et al., 2021a). We are primarily aiming for \mathbb{S}_I to inherit the relevant SR structure from \mathbb{S}_T and the training data, while preserving the image encoder’s feature extraction capabilities. In Section 4.2, we provide illustrations that this happens in our experiments. Additional training and architecture details are in Appendix A.
- Text decoder.** The final component is a text decoder $\mathcal{W} : \mathbb{S}_I \rightarrow \mathbb{G}$ that maps points in the aligned embedding space \mathbb{S} to symbolic expressions. To improve decoding performance for symbolic regression, we adopt an expression decoder (Kamienny et al., 2022), which overlays a multi-layer Transformer atop the numeric encoder to translate encodings into symbolic expressions. Similar to prior work (Meidani et al., 2023; Radford et al., 2021a), we first train the decoder with both encoders frozen, and then fine-tune all components jointly.

Leveraging the inductive biases of pretrained encoders enables us to reduce data requirements while improving performance in symbolic regression. Our model, **README**, trained on only ~ 1 million synthetic numeric-symbolic pairs, outperforms SNIP (Meidani et al., 2023) which is a recent multi-modal pretraining approach trained on 60 million examples. This demonstrates the strong data efficiency and generalization capabilities of our modular encoder choice.

Model	Pretraining Data	Mean R^2_{Test}
README	~ 1 million pairs	0.984 ± 0.004
SNIP	~ 60 million pairs	0.883 ± 0.091

Table 1. Comparison between README and SNIP models pre-trained on different volumes of synthetic data. Mean R^2 Test Score shown is for real-world Physics-Informed dataset. Similar results for other datasets are provided in Appendix B.4.1, with further discussion on evaluation and metrics in Section 4.3.

3.3. Inference optimization

In README, inference consists of two processes, as illustrated in Figure 2.

Inference decoding process. We first convert the numerical dataset \mathcal{D} into a plot image (Section 3.1), before mapping it through the image encoder to its embedding space representation $\tilde{s} = \mathcal{I}(\mathcal{P}(\mathcal{D}))$. In some cases, direct decoding using our text decoder $\mathcal{W}(\tilde{s})$ would already achieve a sufficiently good symbolic expression \tilde{G} for the dataset. However, README is designed to have the decoder just find the right symbolic expression skeleton, before doing ‘constants optimization’ via BFGS (Fletcher, 1987) similar to past works (Kamienny et al., 2022) where numerical constants in the expression may possibly be refined based on some metric (e.g., R^2) evaluated over \mathcal{D} . The decoding process is summarized in Algorithm 1.

Inference optimization process. The optimization process involves searching in the image embedding space \mathbb{S}_I for the best point s^* that would be decoded to the best symbolic expression G^* . For most settings, we do this by first generating a candidate population within a region around s^* , and using the Grey Wolf Optimizer (GWO) (Mirjalili et al., 2014) to find s^* . Given the advantages from README’s data processing and model, applying a vanilla GWO optimizer would typically already give SOTA performance (see Section 4 for details). However, under very tight time constraints, e.g., $\leq 10s$ where none except one of our baseline

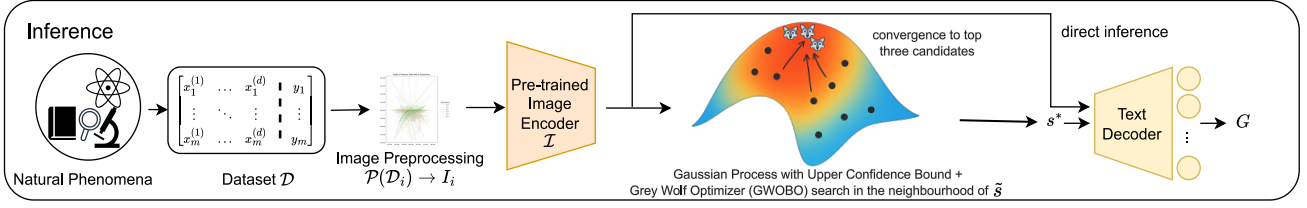


Figure 2. README Inference Optimization process

Algorithm 1 README inference decode algorithm

- 1: **Input:** Image encoder \mathcal{I} , Inference dataset \mathcal{D}
 - 2: **Output:** Symbolic expression \tilde{G} , Associated MSE loss on inference dataset \tilde{L}
 - 3: $I = \mathcal{P}(\mathcal{D})$ //Process dataset to image plot
 - 4: $\tilde{s} = \mathcal{I}(I)$ //Pass image through image encoder
 - 5: $G = \mathcal{W}(\tilde{s})$ //Decode to symbolic expression
 - 6: Run BFGS(G, \mathcal{D}) (Fletcher, 1987) to optimize for MSE loss L evaluated over \mathcal{D} to obtain \tilde{G} and final loss \tilde{L}
 - 7: **Return** \tilde{G} and \tilde{L}
-

methods manage to finish running, we need a faster optimization process. The bottleneck for GWO lies in the evaluation of entire population’s fitness score, which requires running the decoding process to get a symbolic expression and compute its R^2 .

Hence, we propose a hybrid GWO algorithm (GWOBO), that employs Bayesian Optimization (BO) as a supporting subroutine for GWO to (1) train and provide a Gaussian Process (GP) surrogate model to model the fitness value (R^2) given any s , and (2) pick the top three wolves (α, β, γ) that will influence the exploration of the rest of the population. Specifically, we iteratively run GWO and BO: after each iteration of GWO, we train the GP with past decoded points and run BO with an Upper Confidence Bound (UCB) acquisition function (see Appendix D for details) to pick the top three wolves for the next GWO iteration. The top three wolves will have their fitness score evaluated by the decoding process, while the rest of the population will have its fitness score estimated using the GP trained from the BO process. The algorithm is outlined in Algorithm 2. In high-resource settings, the decoder can also be parallelized across multiple GPUs, which allows a larger candidate pool to be explored and improves performance.

4. Experimental results

4.1. Experimental setup

Datasets. We evaluate README on three datasets with varying characteristics: the Strogatz dataset consisting of synthetic data of 2-state dynamic models (Strogatz, 2024; La Cava et al., 2016), the CP3-Bench astrophysical

Algorithm 2 README inference optimization process

- 1: **Input:** Image encoder \mathcal{I} , Inference dataset \mathcal{D} , Target loss \tilde{L} , Max iterations T
 - 2: **Output:** Best-fit symbolic expression and loss $r^* = (G^*, L^*)$
 - 3: Initialize GP
 - 4: **for** $t = 1, 2, \dots, T$ **do**
 - 5: Optimize via GWO to select top three wolves’ embedding $(s_t^\alpha, s_t^\beta, s_t^\gamma)$ based on UCB criterion
 - 6: Perform GP regression with $\{(s_t^a, R^2(s_t^a))\}_{a \in [\alpha, \beta, \gamma]}^{t-1}$
 - 7: Decode s_t^a with Algorithm 1 lines 5-6 to obtain $r_t^a = (\tilde{G}_t, \tilde{L}_t)$, $a \in [\alpha, \beta, \gamma]$
 - 8: Obtain corresponding R^2 score $R^2(s_t^a)$ for the expression \tilde{G}_t
 - 9: Find $r^* \in \{r_t\}_{t=1}^T$ with the lowest \tilde{L}^* .
 - 10: Exit if $L^* < \tilde{L}$
 - 11: **end for**
 - 12: **Return** r^*
-

dataset (Thing & Koksang, 2024) of synthetic data based on cosmological equations with added noise and varying precision, and problems from physics-informed experimental design (Hemachandra et al., 2025) (PIED) consisting of collated real-world, noisy experimental data. These datasets cumulatively provide 61 real-world regression problems on which README is benchmarked (see Appendix B for details).

Models. We primarily used ViT-Base as the image encoder and MathBERT as the text encoder for our experiments, but also observed strong results with encoders from other model families. This highlights the flexibility and performance of our modular framework (see Appendix A.2). Numeric data is converted into visual plots and processed using the Vision Transformer model google/vit-base-patch16-224-in21k (Dosovitskiy et al., 2021), which is pre-trained on large-scale image datasets for strong pattern recognition. Symbolic equations are encoded using tbs17/MathBERT (Shen et al., 2021), a model trained on mathematical texts to better capture the structure and semantics of symbolic expressions. This combination allows the model to benefit from both

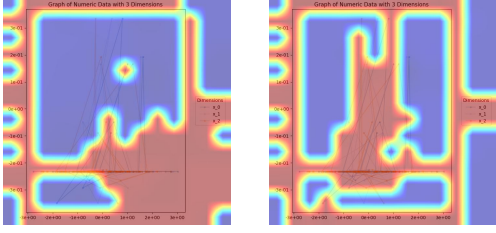


Figure 3. Image encoder attention rollout before (left) and after (right) training. Red indicates higher attention.

ViT’s visual representation capabilities and MathBERT’s inductive biases for symbolic reasoning.

Benchmarks. README is benchmarked against 9 algorithms, including Operon (Burlacu et al., 2020), Interaction-Transformation Evolutionary Algorithm (ITEA) (Aldeia & de França, 2022), Genetic Programming Gene-pool Optimal Mixing Evolutionary Algorithm for Genetic Programming (GPGOMEA) (Virgolin & Bosman, 2022), and Fast Function Extraction (FFX) (McConaghy, 2011) from SR-Bench (La Cava et al., 2016), and Meidani et al.’s Symbolic-Numeric Integrated Pretraining (SNIP). The selected algorithms are efficient and effective, and represent a diverse range of SR approaches.

Evaluation. We evaluate algorithms over three metrics, as described in Section 2. First, we evaluate the **accuracy** of the generated expression \hat{G} and its associated function $g(x)$ by computing its R^2 over test data points. Secondly, expression **complexity** or **parsimony** is evaluated as the number of nodes in its associated expression tree. A smaller expression complexity, or a more parsimonious expression, is better. Lastly, wall-clock **time** is used as a measure of the speed of the SR algorithm. This is measured as the time taken to train the algorithm. A smaller wall-clock time is better. All experiments are repeated for five random seeds.

4.2. Image encoder and embedding space structure

Attention visualization. We first analyze how the image encoder processes the input image plots, by visualizing the attention via attention rollout in Figure 3. Attention rollout involves recursively multiplying attention weights across all layers (Abnar & Zuidema, 2020), where we selected the minimum attention weight over all heads at each layer. We observed that the trained image encoder correctly focuses on important areas of the image plot, such as the graph, axes and legend, compared to the pre-trained ViT which misses portions of the graph and focuses on irrelevant blank spaces.

Families of equations. The t-SNE visualization of MathBERT’s embeddings in Figure 4a shows a well-structured embedding space, where different families of equations are

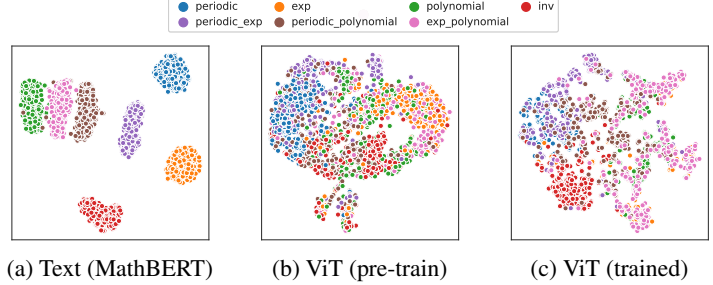


Figure 4. t-SNE plots of text and image embeddings.

separately clustered. Combinations of different families also result in embeddings being close to the original family cluster, where the ‘periodic_polynomial’ and ‘exp_polynomial’ clusters are close to the ‘polynomial’ cluster.

Transferring pre-trained math structure. The well-structured embedding space of MathBERT also transfers to the image encoder after training. While initial pre-trained ViT’s embeddings do not have meaningful structures (Figure 4b), a clearer separation between equation families emerges after training (Figure 4c). For example, the image embeddings of ‘inv’ equations become distinctly clustered, and ‘periodic’ and ‘periodic_exp’ equations are grouped together. Interestingly, ‘periodic_polynomial’ embeddings seem to be interpolated between ‘periodic’ and ‘polynomial’ embeddings, indicating the image encoder has recognized relationships between equation families.

4.3. Performance evaluation

We compare README against the baseline algorithms on the task of regressing 61 regression problems. For each problem, 75% of the data is designated as the inference set while the remaining 25% is designated as the test set. Each trial involves 200 randomly selected inference data points provided to the methods for them to generate symbolic expressions based on the data, with these expressions evaluated on randomly selected test set. Five trials are conducted for each problem. We consider two experimental settings. First, the rapid setting, where algorithms have to be trained within a 30-second cut-off. Second, the ultra rapid setting, where algorithms have to be trained within a 10-second cut-off.

4.3.1. RESULTS FOR RAPID SETTING

To evaluate both accuracy and parsimony of the symbolic equations produced by each method, we plot Pareto plots for each dataset where the y-axis represents accuracy measured by mean R^2_{Test} (larger is better) and the x-axis represents parsimony measured by mean equation length (smaller is better). The x-axis is plotted in descending order so that along both axes, points furthest from the origin are

Algorithm	Strogatz		CP3-Bench		Physics Informed	
	Mean R^2_{Test} (\uparrow better)	Equation Length (\downarrow better)	Mean R^2_{Test}	Equation Length	Mean R^2_{Test}	Equation Length
README	0.880 \pm 0.018	20.99 \pm 1.19	0.933 \pm 0.018	23.34 \pm 0.11	0.984 \pm 0.004	23.76 \pm 2.35
FFXRegressor	0.822 \pm 0.063	198.43 \pm 37.82	0.930 \pm 0.006	172.45 \pm 15.37	0.930 \pm 0.007	198.41 \pm 15.47
GPGOMEAREgressor (*)	0.849 \pm 0.054	29.27 \pm 3.46	0.910 \pm 0.004	29.51 \pm 0.83	0.930 \pm 0.003	29.85 \pm 0.38
AFPRRegressor	0.742 \pm 0.043	38.00 \pm 6.08	0.887 \pm 0.003	39.63 \pm 3.76	0.912 \pm 0.007	42.67 \pm 2.84
SNIP	0.780 \pm 0.046	22.39 \pm 1.98	0.909 \pm 0.011	25.07 \pm 1.10	0.895 \pm 0.045	28.04 \pm 1.44
DSRRegressor (*)	0.693 \pm 0.067	19.16 \pm 3.49	0.803 \pm 0.005	17.07 \pm 1.61	0.891 \pm 0.041	30.00 \pm 0.00
EPLEXRegressor (*)	0.591 \pm 0.076	45.09 \pm 3.69	0.886 \pm 0.021	51.96 \pm 1.45	0.869 \pm 0.026	47.40 \pm 3.53
EHRegressor	0.644 \pm 0.042	19.13 \pm 1.57	0.862 \pm 0.006	21.31 \pm 1.14	0.818 \pm 0.033	25.53 \pm 1.75
OperonRegressor	0.849 \pm 0.053	60.46 \pm 2.25	0.882 \pm 0.013	77.19 \pm 1.12	0.806 \pm 0.082	80.98 \pm 1.36
ITEAREgressor	0.736 \pm 0.052	12.83 \pm 0.27	0.918 \pm 0.002	16.08 \pm 0.56	0.777 \pm 0.017	11.67 \pm 0.40

Table 2. Comparison of symbolic regression algorithms across three datasets: Strogatz, CP3-Bench, and Physics Informed. Best performers are bolded. Algorithms marked with an asterisk (*) did not complete within the 30-second time budget (see Appendix B.4.2).

the best.

We indicate Pareto frontiers in different colors. Points within the same frontier are non-dominated with respect to each other, meaning no method in the frontier outperforms another on both accuracy and parsimony. A higher Pareto frontier contains at least one model that Pareto-dominates a model on a lower frontier.

As shown in Figure 5, README and ITEA are consistently in the top Pareto frontier. However, ITEA is in the top frontier because it has a strong bias for small expression size – its accuracy values tend to be among the lowest compared to the rest of the methods, especially for the Physics Informed dataset. In contrast, README identifies parsimonious and accurate equations rapidly (within 30 seconds) and consistently lies on the first Pareto frontier. It achieves the highest mean R^2_{Test} among all algorithms, demonstrating strong overall accuracy.

Table 2 shows the detailed results on accuracy and parsimony for all methods across the three datasets. Among the methods, GPGOMEA, DSR, and EPLEX Regressor exceeded the 30-second limit — GPGOMEA for the Strogatz dataset, and DSR and EPLEX across all three, sometimes taking up to 4 times the time budget to produce a result. However, despite taking longer time, these method still underperform README that is run within the 30-second time budget. Running configurations and detailed results for all algorithms are provided in Appendix B.3.1 and Appendix B.4.2, respectively.

4.3.2. EVALUATION RESULTS FOR ULTRA RAPID SETTING

Next, we analyze the ultra-rapid setting that requires methods to complete inference within 10 seconds. This setting is motivated by applications requiring low-latency predictions of physical movement, such as physics checking in synthetic video generation or real-time decision support, where inference must be done within a matter of seconds. We evaluate methods on the Physics Informed dataset, which is from

real-world experiments and hence serves as a test of method robustness in practical, noisy environment.

To achieve ultra-rapid inference, we use our **GWBO** optimization process when running README as described in Section 3.3, which speeds up the inference process while still achieving good results. Among the baselines, only README with GWBO and FFXRegressor are able to consistently return symbolic expressions under the ultra-rapid setting time limit of 10 seconds. In this setting, README continues to outperform FFXRegressor in both parsimony and accuracy. README achieved an R^2_{Test} of 0.958 with average equation length of approximately **18 terms**, while FFXRegressor only achieved R^2_{Test} of 0.930 with significantly higher (i.e., worse) average Equation Length of approximately **219 terms**.

In addition, we tested the methods’ sensitivity to noise by introducing additional Gaussian noise to data observations \tilde{y} and evaluating its performance. Figure 6 shows how the accuracy performance (R^2_{Test}) of both methods changes over increasing noise level (i.e., the Gaussian noise standard deviation is varied from 0.1 to 0.5 times the root mean square of the observation values). Note that README demonstrates greater robustness to noise with less performance degradation as noise is increased, e.g., README’s accuracy only decreased from 0.958 to 0.857 when noise of 0.1 noise level is added, but FFXRegressor’s accuracy had a much larger drop from 0.930 to 0.350.

4.4. Demonstration of Equation Prediction for Noisy Real Experimental Data

To demonstrate the practicality of our framework for an ultra-rapid setting, we analyze two real-world videos where object motion must obey physical laws. We adopt a 10-second runtime constraint to reflect real-time applications such as physics validation in synthetic video generation and decision-making systems that require relatively low latency and physically accurate predictions. In these scenarios, fast feedback is essential, and we benchmark symbolic

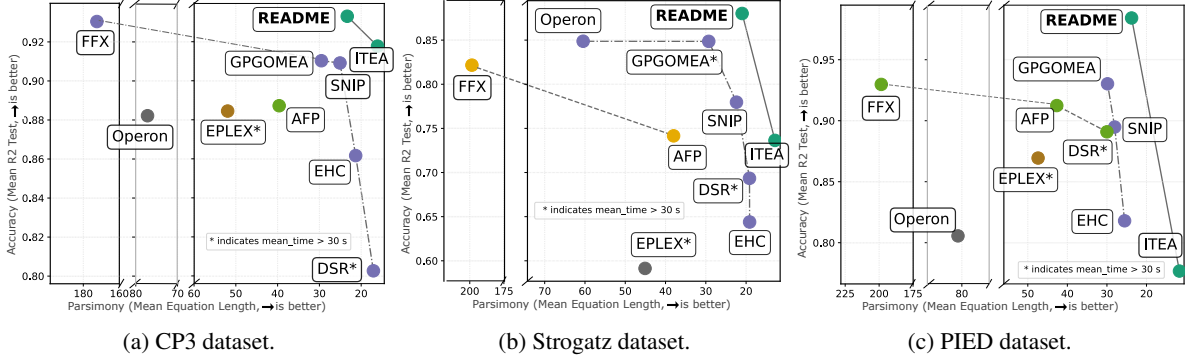


Figure 5. Pareto plots for all algorithms.

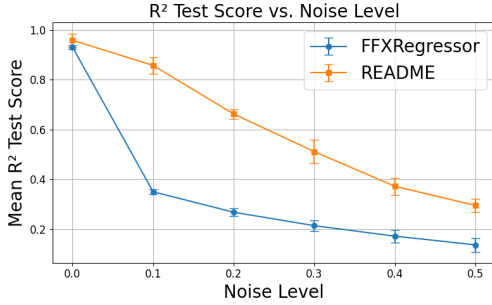


Figure 6. Mean R^2_{Test} vs. target noise for FFXRegressor (orange) and README (blue) in the ultra-rapid setting (<10 seconds).

regression (SR) algorithms based on their ability to generate interpretable equations within this strict time limit.

We developed two pipelines to extract object coordinates over time and apply symbolic regression to model their trajectories. For the pendulum video, we used Tracker software with basic computer vision techniques such as template matching to estimate and track the pendulum bob across 311 frames. For the ping pong video, we used a YOLOv8n model to detect the ball in each frame, extracted the center of its bounding box, and obtained 85 position points.

Table 3 summarizes the performance of each symbolic regression algorithm under the 10-second constraint on the *Pendulum Swinging* and *Ping Pong Ball Bouncing* datasets. Among all evaluated methods, README consistently achieves the best R^2 scores while maintaining compact Equation Lengths, making it the most effective approach under strict time constraints.

To assess performance beyond the 10-second limit, we conducted an additional evaluation with a 5-minute time budget using Grey Wolf Optimization across five random seeds. README achieved average R^2 scores above 0.99 on the pendulum dataset and above 0.96 on the ping pong dataset. The slightly lower score for the ping pong video is due to

Algorithm	Pendulum Swinging		Ping Pong Bouncing	
	R^2	Equation Length	R^2	Equation Length
README	0.686 ± 0.263	15.80 ± 1.79	0.862 ± 0.204	24.90 ± 6.92
GPGOMEAREgressor	0.018 ± 0.018	31.00 ± 0.00	—	—
FFXRegressor	0.012 ± 0.016	81.10 ± 59.31	0.000 ± 0.000	103.00 ± 0.00
ITEAREgressor	0.000 ± 0.000	9.75 ± 1.77	0.235 ± 0.000	10.00 ± 0.87
OperonRegressor	—	—	0.004 ± 0.008	79.60 ± 5.13

Table 3. Performance of symbolic regression algorithms on two real-world videos under a 10-second time constraint. A dash (—) indicates no result was returned within the limit.

missing frames around the bounce point, resulting in an incomplete trajectory. For both experiments, we trained on the first 75% of the time-sequenced data and tested on the final 25%, demonstrating README’s strong extrapolation capability on real-world data and showcasing its potential for applications in motion prediction.

5. Conclusion

We introduced README, a framework for rapid equation discovery that uses (1) an informative, compressed image representation of numerical data; (2) an efficiently-trained transformer-based model built on top of pre-trained image and text encoders ($\sim 60\times$ less training data compared to past works); and (3) a novel combination of BO and GWO for a rapid, effective optimization process to achieve state-of-the-art and robust SR results for challenging settings with realistic, noisy settings and tight time constraints (<10 seconds).

This work represents a first step toward quick and reliable symbolic regression that can be used as a module within real-world tasks. Potential applications include computer vision and robotics, where real-time, interpretable physics validation and decision-making are essential. This approach also holds promise for domains such as video analytics and synthetic video generation, where low-latency fast symbolic regression is crucial. Potential future work can expand on the type of equations that can be extracted from datasets, such as differential equations.

Acknowledgments

This research/project is supported by the National Research Foundation, Singapore under its AI Singapore Programme (AISG Award No: AISG2-PhD/2023-01-039J) and is part of the programme DesCartes which is supported by the National Research Foundation, Prime Minister's Office, Singapore under its Campus for Research Excellence and Technological Enterprise (CREATE) programme.

References

- Learning symbolic forward models for robotic motion planning and control (full article)*, volume ECAL 2011: The 11th European Conference on Artificial Life of *Artificial Life Conference Proceedings*, 08 2011. doi: 10.7551/978-0-262-29714-1-ch086. URL <https://doi.org/10.7551/978-0-262-29714-1-ch086>.
- Abnar, S. and Zuidema, W. Quantifying attention flow in transformers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4190–4197, 2020.
- AL NAJAR, M., ALMAR, R., BERGSMA, E., DELVIT, J.-M., and Wilson, D. Improving a shoreline forecasting model with symbolic regression. In *ICLR 2023 Workshop on Tackling Climate Change with Machine Learning*, 2023. URL <https://www.climatechange.ai/papers/iclr2023/21>.
- Aldeia, G. S. I. and de França, F. O. Interaction-transformation evolutionary algorithm with coefficients optimization. In *Proceedings of the Genetic and Evolutionary Computation Conference Companion, GECCO '22*, pp. 2274–2281, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450392686. doi: 10.1145/3520304.3533987. URL <https://doi.org/10.1145/3520304.3533987>.
- Anđelić, N., Lorencin, I., Mrzljak, V., and Car, Z. On the application of symbolic regression in the energy sector: Estimation of combined cycle power plant electrical power output using genetic programming algorithm. *Engineering Applications of Artificial Intelligence*, 133:108213, 2024. ISSN 0952-1976. doi: <https://doi.org/10.1016/j.engappai.2024.108213>. URL <https://www.sciencedirect.com/science/article/pii/S0952197624003713>.
- Biggio, L., Bendinelli, T., Neitz, A., Lucchi, A., and Parascandolo, G. Neural symbolic regression that scales. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 936–945. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/biggio21a.html>.
- Brown, D. Tracker video analysis and modeling tool. <https://opensourcephysics.github.io/tracker-website/>, 2024. Accessed: 2025-05-23.
- Brunton, S. L., Proctor, J. L., and Kutz, J. N. Discovering governing equations from data by sparse

- identification of nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*, 113(15):3932–3937, 2016. doi: 10.1073/pnas.1517384113. URL <https://www.pnas.org/doi/abs/10.1073/pnas.1517384113>.
- Brunton, S. L., Nathan Kutz, J., Manohar, K., Aravkin, A. Y., Morgansen, K., Klemisch, J., Goebel, N., Buttrick, J., Poskin, J., Blom-Schieber, A. W., Hogan, T., and McDonald, D. Data-driven aerospace engineering: Reframing the industry with machine learning. *AIAA Journal*, pp. 1–26, July 2021. doi: 10.2514/1.j060131. URL <http://dx.doi.org/10.2514/1.J060131>.
- Burlacu, B., Kronberger, G., and Kommenda, M. Operon c++: An efficient genetic programming framework for symbolic regression. In *Proceedings of the 2020 Genetic and Evolutionary Computation Conference Companion*, GECCO ’20, pp. 1562–1570, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450371278. doi: 10.1145/3377929.3398099. URL <https://doi.org/10.1145/3377929.3398099>.
- Champion, K., Lusch, B., Kutz, J. N., and Brunton, S. L. Data-driven discovery of coordinates and governing equations. *Proceedings of the National Academy of Sciences*, 116(45):22445–22451, 2019. doi: 10.1073/pnas.1906995116. URL <https://www.pnas.org/doi/abs/10.1073/pnas.1906995116>.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houshy, N. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- Fitzsimmons, J. and Moscato, P. Symbolic regression modeling of drug responses. In *2018 First International Conference on Artificial Intelligence for Industries (AI4I)*, pp. 52–59, 2018. doi: 10.1109/AI4I.2018.8665684.
- Fletcher, R. *Practical methods of optimization; (2nd ed.)*. Wiley-Interscience, USA, 1987. ISBN 0471915475.
- Garnett, R. *Bayesian Optimization*. Cambridge Univ. Press, 2022.
- Goldberg, D. E. *Genetic Algorithms in Search, Optimization and Machine Learning*. Addison-Wesley Longman Publishing Co., Inc., USA, 1st edition, 1989. ISBN 0201157675.
- Grundner, A., Beucler, T., Gentine, P., and Eyring, V. Data-driven equation discovery of a cloud cover parameterization, February 2024. URL <http://dx.doi.org/10.1029/2023MS003763>.
- Guimerà, R., Reichardt, I., Aguilar-Mogas, A., Massucci, F. A., Miranda, M., Pallarès, J., and Sales-Pardo, M. A bayesian machine scientist to aid in the solution of challenging scientific problems. *Science Advances*, 6(5):eaav6971, 2020. doi: 10.1126/sciadv.aav6971. URL <https://www.science.org/doi/abs/10.1126/sciadv.aav6971>.
- Hemachandra, A., Lau, G. K. R., Ng, S.-K., and Low, B. K. H. PIED: Physics-informed experimental design for inverse problems. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=w7P92BESb2>.
- Johnson, W. B. and Lindenstrauss, J. Extensions of lipschitz mappings into a hilbert space. *Contemporary Mathematics*, 26:189–206, 1984.
- Kamienny, P.-A., d’Ascoli, S., Lample, G., and Charton, F. End-to-end symbolic regression with transformers, April 2022.
- Koza, J. Genetically breeding populations of computer programs to solve problems in artificial intelligence. In *[1990] Proceedings of the 2nd International IEEE Conference on Tools for Artificial Intelligence*, pp. 819–827, 1990. doi: 10.1109/TAI.1990.130444.
- Koza, J. R. Hierarchical genetic algorithms operating on populations of computer programs. In *International Joint Conference on Artificial Intelligence*, 1989. URL <https://api.semanticscholar.org/CorpusID:17882725>.
- Koza, J. R. Genetic programming as a means for programming computers by natural selection. *Statistics and Computing*, 4(2):87–112, Jun 1994. ISSN 1573-1375. doi: 10.1007/BF00175355. URL <https://doi.org/10.1007/BF00175355>.
- La Cava, W., Spector, L., Danai, K., and Lackner, M. Evolving differential equations with developmental linear genetic programming and epigenetic hill climbing. In *Proceedings of the Companion Publication of the 2014 Annual Conference on Genetic and Evolutionary Computation*, GECCO Comp ’14, pp. 141–142, New York, NY, USA, 2014. Association for Computing Machinery. ISBN 9781450328814. doi: 10.1145/2598394.2598491. URL <https://doi.org/10.1145/2598394.2598491>.

- La Cava, W., Danai, K., and Spector, L. Inference of compact nonlinear dynamic models by epigenetic local search. *Engineering Applications of Artificial Intelligence*, 55:292–306, 2016. ISSN 0952-1976. doi: <https://doi.org/10.1016/j.engappai.2016.07.004>. URL <https://www.sciencedirect.com/science/article/pii/S0952197616301294>.
- La Cava, W., Helmuth, T., Spector, L., and Moore, J. H. A probabilistic and multi-objective analysis of lexicase selection and ϵ -lexicase selection. *Evol. Comput.*, 27(3): 377–402, September 2019. ISSN 1063-6560. doi: 10.1162/evco_a_00224. URL https://doi.org/10.1162/evco_a_00224.
- Liu, Y., Zhang, Z., and Schaeffer, H. Prose: Predicting operators and symbolic expressions using multimodal transformers, 2023. URL <https://arxiv.org/abs/2309.16816>.
- Makke, N. and Chawla, S. Interpretable scientific discovery with symbolic regression: a review. *Artificial Intelligence Review*, 57(1):2, Jan 2024. ISSN 1573-7462. doi: 10.1007/s10462-023-10622-0. URL <https://doi.org/10.1007/s10462-023-10622-0>.
- McConaghy, T. *FFX: Fast, Scalable, Deterministic Symbolic Regression Technology*, pp. 235–260. Springer New York, New York, NY, 2011. ISBN 978-1-4614-1770-5. doi: 10.1007/978-1-4614-1770-5_13. URL https://doi.org/10.1007/978-1-4614-1770-5_13.
- Meidani, K., Shojaei, P., Reddy, C. K., and Farimani, A. B. SNIP: Bridging Mathematical Symbolic and Numeric Realms with Unified Pre-training. In *The Twelfth International Conference on Learning Representations*, October 2023.
- Mirjalili, S., Mirjalili, S. M., and Lewis, A. Grey wolf optimizer. *Advances in engineering software*, 69:46–61, 2014.
- Oord, A. v. d., Li, Y., and Vinyals, O. Representation learning with contrastive predictive coding. In *Advances in Neural Information Processing Systems*, pp. 10203–10214, 2018.
- Petersen, B. K., Larma, M. L., Mundhenk, T. N., Santiago, C. P., Kim, S. K., and Kim, J. T. Deep symbolic regression: Recovering mathematical expressions from data via risk-seeking policy gradients. In *International Conference on Learning Representations*, October 2020.
- Petersen, B. K., Larma, M. L., Mundhenk, T. N., Santiago, C. P., Kim, S. K., and Kim, J. T. Deep symbolic regression: Recovering mathematical expressions from data via risk-seeking policy gradients. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=m5Qsh0kBQG>.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., and Sutskever, I. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021a.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., and Sutskever, I. Learning Transferable Visual Models From Natural Language Supervision. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 8748–8763. PMLR, July 2021b.
- Rudin, C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206–215, May 2019. ISSN 2522-5839. doi: 10.1038/s42256-019-0048-x. URL <https://doi.org/10.1038/s42256-019-0048-x>.
- Schmidt, M. D. and Lipson, H. Age-fitness pareto optimization. In *Proceedings of the 12th Annual Conference on Genetic and Evolutionary Computation, GECCO '10*, pp. 543–544, New York, NY, USA, 2010. Association for Computing Machinery. ISBN 9781450300728. doi: 10.1145/1830483.1830584. URL <https://doi.org/10.1145/1830483.1830584>.
- Shen, J. T., Yamashita, M., Prihar, E., Heffernan, N., Wu, X., Graff, B., and Lee, D. Mathbert: A pre-trained language model for general nlp tasks in mathematics education. *arXiv preprint arXiv:2106.07340*, 2021.
- Shojaei, P., Meidani, K., Gupta, S., Farimani, A. B., and Reddy, C. K. LLM-SR: Scientific Equation Discovery via Programming with Large Language Models, June 2024.
- Smits, G. F. and Kotanchek, M. *Pareto-Front Exploitation in Symbolic Regression*, pp. 283–299. Springer US, Boston, MA, 2005.
- Strogatz, S. H. *NONLINEAR DYNAMICS AND CHAOS, THIRD EDITION: With applications to physics, biology, chemistry,... And engineering, third edition, student's solution*. 2024.
- Thing, M. E. and Koksang, S. M. cp3-bench: A tool for benchmarking symbolic regression algorithms tested with cosmology, 2024. URL <https://arxiv.org/abs/2406.15531>.
- Valipour, M., You, B., Panju, M., and Ghodsi, A. SymbolicGPT: A Generative Transformer Model for Symbolic Regression. <https://arxiv.org/abs/2106.14131v1>, June 2021.

-
- Virgolin, M. and Bosman, P. A. N. Coefficient mutation in the gene-pool optimal mixing evolutionary algorithm for symbolic regression. In *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, GECCO '22, pp. 2289–2297, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450392686. doi: 10.1145/3520304.3534036. URL <https://doi.org/10.1145/3520304.3534036>.
- Virgolin, M. and Pissis, S. P. Symbolic regression is NP-hard. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL <https://openreview.net/forum?id=LTiaPxqe2e>.
- Wahlquist, Y., Sundell, J., and Soltesz, K. Learning pharmacometric covariate model structures with symbolic regression networks. *Journal of Pharmacokinetics and Pharmacodynamics*, 51(2):155–167, Apr 2024. ISSN 1573-8744. doi: 10.1007/s10928-023-09887-3. URL <https://doi.org/10.1007/s10928-023-09887-3>.
- Wang, C., Zhang, Y., Wen, C., Yang, M., Lookman, T., Su, Y., and Zhang, T.-Y. Symbolic regression in materials science via dimension-synchronous-computation. *Journal of Materials Science & Technology*, 122:77–83, 2022. ISSN 1005-0302. doi: <https://doi.org/10.1016/j.jmst.2021.12.052>. URL <https://www.sciencedirect.com/science/article/pii/S1005030222002055>.
- Wang, Y., Wagner, N., and Rondinelli, J. M. Symbolic regression in materials science. *MRS Communications*, 9(3):793–805, Sep 2019. ISSN 2159-6867. doi: 10.1557/mrc.2019.85. URL <https://doi.org/10.1557/mrc.2019.85>.
- Zhang, Z. and Chen, Z. Modeling and control of robotic manipulators based on symbolic regression. *IEEE Transactions on Neural Networks and Learning Systems*, 34(5): 2440–2450, 2023. doi: 10.1109/TNNLS.2021.3106648.

A. Additional training and architecture details

A.1. Model Training Details

A.1.1. TRAINING DATA

To train our model, we generated synthetic pairs of numeric and symbolic data using the publicly available codebase in (Kamienny et al., 2022), following the data generation settings used by SNIP (Meidani et al., 2023). This includes operator downsampling and restricting expressions to at most 10 input dimensions. The only difference is that we generated a total of ~1 million (image, equation) pairs for training, whereas SNIP used ~60 million pairs to pretrain their numeric and symbolic encoders.

Numeric Data Visualization For each equation, input data x with dimensionality $n \leq 10$ was generated, comprising 200 data points represented as $200 \times n$ matrices. Each input dimension was paired with targets y , represented as a 200×1 vector. Each input dimension was plotted individually against the target y using Matplotlib, assigning different colors for clarity. Each graph includes the dimensionality information in its title. Figure 7 shows a sample graph from our dataset, illustrating how patterns are captured across different dimensions.

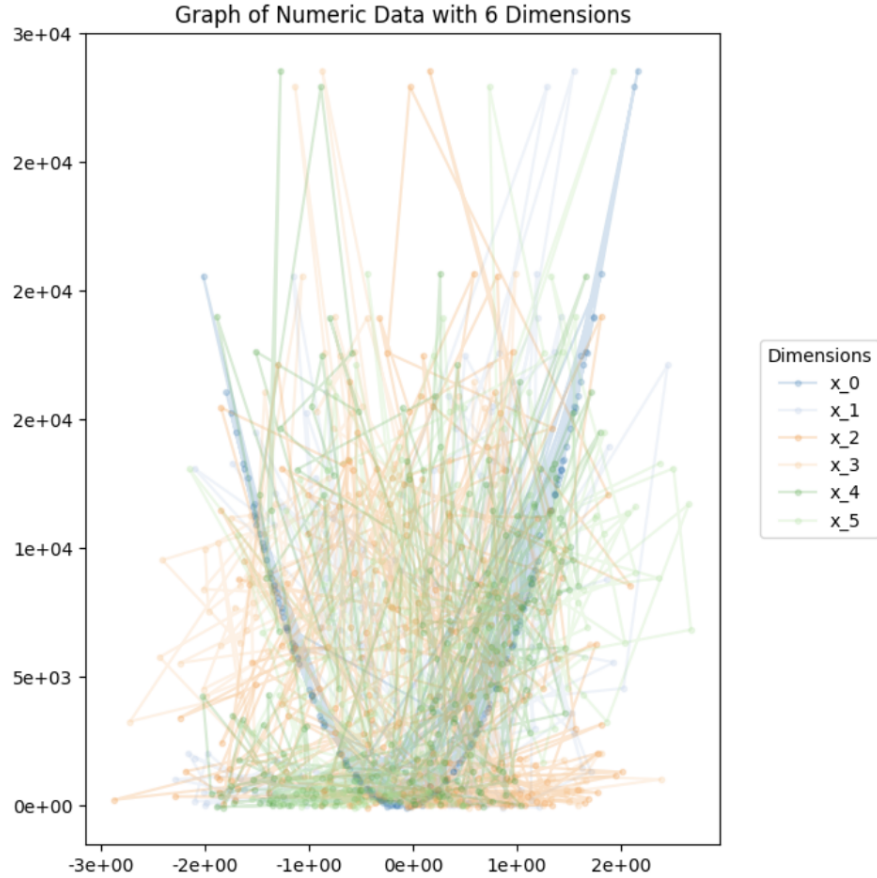


Figure 7. Sample visualization of numeric data with multiple dimensions plotted against the target. Patterns across dimensions are captured effectively.

Equation Representation Expression trees were converted into their equivalent infix notation, providing readable symbolic equations. Each visual plot and corresponding symbolic equation formed a training pair.

A.1.2. IMAGE ENCODER

We employed a pre-trained Vision Transformer model `google/vit-base-patch16-224-in21k` (Dosovitskiy et al., 2021). This model, trained extensively on diverse image data, excels in pattern recognition.

This approach offers two primary advantages:

1. **Pattern Recognition:** The model swiftly identifies patterns within numeric data visualizations, leveraging robust feature extraction capabilities from its pre-training.
2. **Regularization Effect:** By plotting each input dimension against the target in the same visualization, the model is naturally regularized, treating all dimensions uniformly. This method helps prevent overfitting to any specific dimension.

We also experimented with the larger `google/vit-huge-patch14-224-in21k` model and observed improved results, as expected due to the increased model capacity. However, for the purpose of balancing performance and computational efficiency in our experiments, we chose to use the smaller base Vision Transformer for performance evaluation.

A.1.3. TEXT ENCODER

The symbolic equations were encoded using `tbs17/MathBERT` (Shen et al., 2021), a model specifically pre-trained on mathematical text, enabling effective interpretation and encoding of symbolic equations. Equations were provided directly in infix notation, aligning well with MathBERT’s pre-training on mathematical textbooks and notations. Utilizing plain text inputs avoids constraints typical of tree-based representations, offering greater flexibility and leveraging the inductive biases inherent to MathBERT.

A.1.4. ALIGNED EMBEDDING SPACE

The README model aligns numeric and symbolic representations in a shared latent space. Inspired by the joint training approach used in CLIP (Radford et al., 2021a), README optimizes a symmetric cross-entropy loss over similarity scores. A contrastive loss based on the InfoNCE objective (Oord et al., 2018) effectively aligns embeddings of matching numeric-symbolic pairs while pushing apart non-matching pairs.

The loss function is defined as:

$$\mathcal{L} = - \sum_{(v,s) \in B} \left[\log \text{NCE}(Z_S, Z_V) + \log \text{NCE}(Z_V, Z_S) \right] \quad (1)$$

where B represents a batch of (symbolic, numeric) data pairs, and $\text{NCE}(Z_S, Z_V)$ and $\text{NCE}(Z_V, Z_S)$ are the symbolic-to-numeric and numeric-to-symbolic contrastive losses, respectively. The symbolic-to-numeric contrastive loss is computed as:

$$\text{NCE}(Z_S, Z_V) = \frac{\exp(Z_S \cdot Z_V^+ / \tau)}{\sum_{Z \in \{Z_V^+, Z_V^-\}} \exp(Z_S \cdot Z / \tau)}$$

where τ is a temperature parameter, Z_V^+ represents positive numeric embeddings that correspond to the symbolic embedding Z_S , and Z_V^- are negative embeddings from other batch data. This symmetric contrastive loss encourages alignment of numeric and symbolic pairs while separating unrelated pairs.

A.1.5. TEXT DECODER

For decoding symbolic equations, we also adopted the decoder architecture detailed in (Kamienny et al., 2022), consisting of 16 transformer decoder layers. This architecture effectively leverages attention mechanisms to autoregressively generate equations from the aligned embedding representations, benefiting from its deep, layered structure which facilitates complex symbolic regression tasks.

Following prior work (Meidani et al., 2023), training is conducted in two stages. First, the decoder is trained with the image and text encoders frozen, allowing it to learn how to decode from the latent space. Next, the encoders and decoder are

fine-tuned together, so the representations become better suited for decoding symbolic equations. The decoder is supervised using cross-entropy loss over the target symbolic sequence, encouraging accurate reconstruction of symbolic expressions from the shared representation.

A.2. Ablation Studies

For the ultra-rapid setting, where the models are expected output an expression within 10 seconds, we introduce a key contribution: a novel hybrid algorithm that combines the Grey Wolf Optimizer with Bayesian Optimization (GWOBO) to efficiently identify high-quality symbolic expressions under tight runtime constraints.

As shown in Section 4.3.2, our model is also more robust to noise compared to FFXRegressor across varying noise levels. Note that in the section we showed GWOBO results with a candidate set size of 70.

We show that increasing the number of wolf candidates leads to substantial gains in performance while maintaining the same setup described in Section 4.3.2. Specifically, we continue to select the top 3 candidates based on Upper Confidence Bound (UCB) scores, computed using a Gaussian Process with an RBF kernel as detailed in D. Only these top 3 are decoded using Algorithm 1, while the remaining candidates are evaluated using surrogate scores from the GP.

These combined real and surrogate scores are then used to update the population via GWO. As shown in Table 4, GWOBO consistently outperforms pure GWO across all candidate configurations.

Candidates	GWO		GWOBO	
	R^2	Equation Length	R^2	Equation Length
10	0.865 ± 0.058	18.34 ± 1.35	0.880 ± 0.037	16.54 ± 2.36
30	0.903 ± 0.046	19.27 ± 0.76	0.934 ± 0.044	17.05 ± 1.43
50	0.924 ± 0.037	18.24 ± 1.43	0.937 ± 0.026	17.26 ± 1.33
70	0.946 ± 0.029	17.29 ± 1.50	0.958 ± 0.027	16.67 ± 1.57

Table 4. Comparison of GWO and GWOBO at different candidate counts under ultra rapid setting (10 seconds). Each entry reports mean \pm standard deviation. Best values are bolded.

We also conducted an ablation study to compare different encoder architectures, using encoders trained on a smaller set of 512,000 pairs for quick evaluation. As shown in Table 5, ViT-base with MathBERT outperforms CLIP with T5. This is likely due to MathBERT’s strong inductive bias for mathematical structure based on its pretraining (Shen et al., 2021). While ViT is not specifically trained on math-related content (Dosovitskiy et al., 2021), it performed better than CLIP in this setting. We also evaluated ViT-huge and observed marginal gains, but selected ViT-base to ensure efficiency during timed experiments.

Algorithm	Mean R^2 Test Score	Mean Equation Length
vit-huge-with-mathbert	0.767 ± 0.016	14.69 ± 1.60
vit-base-with-mathbert	0.765 ± 0.012	18.56 ± 0.82
clip-with-t5	0.738 ± 0.018	15.13 ± 0.75

Table 5. Comparison of ViT-based and multimodal models on symbolic regression. Best R^2 and smallest Equation Length are bolded. Higher is better for R^2 , lower is better for Equation Length.

B. Additional Experimental Details

B.1. Dataset Details

To comprehensively evaluate our symbolic regression framework, we curated a diverse set of 61 problems drawn from publicly available data sources in multiple scientific domains. These datasets were selected to balance a range of characteristics, including equation complexity, noise levels, and real-world relevance. They span both simulated and experimentally grounded physical systems, allowing us to assess model performance in controlled and practical scenarios.

B.1.1. STROGATZ DATASET

This dataset comprises 14 canonical equations modeling nonlinear dynamical systems, originally drawn from the textbook *Nonlinear Dynamics and Chaos: With Applications to Physics, Biology, Chemistry, and Engineering* by Steven H. Strogatz (Strogatz, 2024; La Cava et al., 2016). These systems describe a range of physical and biological processes such as population dynamics, chemical oscillators, and mechanical systems, and are commonly used in the literature to benchmark symbolic regression algorithms due to their compact closed-form representations and interpretable dynamics.

B.1.2. CP3-BENCH

This dataset consists of 28 equations derived from a diverse collection of real-world and simulated problems across physics, engineering, and environmental science. These problems were introduced as part of the CP3-Bench benchmark (Thing & Koksang, 2024), which was designed to evaluate the capability of scientific equation learning models in recovering compact symbolic expressions from noisy data. The dataset includes systems such as gas solubility estimation, diffusion, and biochemical rate equations, and reflects varying levels of complexity and noise, making it a rigorous testbed for symbolic regression.

B.1.3. PHYSICS INFORMED DATASET

To evaluate the effectiveness of our model on physics-informed problems, we curated a dataset comprising 19 equations. These equations span simulated, experimentally validated, and real-world scenarios, as detailed below:

- **Groundwater Flow (11 equations):**

These equations model steady-state unconfined groundwater flow and are sourced from the paper “*Investigating Steady Unconfined Groundwater Flow using Physics Informed Neural Networks*” by Mohammad Afzal Shadab, Dingcheng Luo, Yiran Shen, Eric Hiatt, and Marc Andre Hesse.

GitHub: <https://github.com/dc-luo/seepagePINN/tree/main>

- **Chromatography (4 equations):**

These equations originate from the experimental validation study “*Can a Computer ‘Learn’ Nonlinear Chromatography?: Experimental Validation of Physics-Based Deep Neural Networks for the Simulation of Chromatographic Processes*” by Sai Gokul Subraveti, Zukui Li, Vinay Prasad, and Arvind Rajendran. The equations simulate nonlinear solute transport in chromatographic columns.

- **Fluid Dynamics (2 equations):**

We include the *nikuradse_1* and *nikuradse_2* equations, which describe the friction factor for turbulent flow in rough pipes based on experimental studies by Johann Nikuradse. These expressions capture the nonlinear relationship between the Darcy friction factor, Reynolds number, and relative roughness in turbulent pipe flow.

- **Pendulum Motion (2 equations):**

These equations are derived from video recordings of a swinging pendulum captured using the Tracker software. The motion was tracked using a fixed camera setup, and position-time data was extracted to recover the underlying physical relationship governing the pendulum’s oscillatory dynamics.

B.2. Experiment Settings

All experiments in section 4.3 were conducted on 2× AMD EPYC 7763 64-Core CPUs and 1× NVIDIA L40 GPUs (CUDA 12.1, Driver 545.23.06), and running Ubuntu 22.04.3 LTS. All software was implemented in Python 3.11.2 using PyTorch 2.0.0, Transformers 4.44.2, and BoTorch 0.11.3.

For each problem described in each dataset (see Appendix B.1.1), we applied a 75%–25% train-test split. To ensure consistent computational constraints across problems, if a problem contained more than 200 training points, the training set was randomly subsampled to retain only 200 points. In experiments involving noise, Gaussian noise with the specified standard deviation was added to the target values of the training set. All experiments were repeated across five random seeds: 23654, 15795, 860, 5390, and 16850. These seeds controlled both the train-test split and the noise sampling, ensuring that our model was evaluated on different subsets of the data in each problem and that its performance was robust to variation in data sampling.

For our experiments under the *rapid* and *ultra-rapid* settings, we initialized both GWO and GWOBO with 71 wolf candidates: 70 perturbations around the base latent encoding plus the original encoding.

B.3. Evaluation Details

B.3.1. OTHER BASELINES CONFIGURATIONS

The baselines used were McConaghy’s FFX Regressor (McConaghy, 2011), de Franca and Aldeia’s ITEA (Aldeia & de França, 2022), Virgolin et al.’s GP-GOMEA (Virgolin & Bosman, 2022), La Cava et al.’s EPLEX (La Cava et al., 2019), Schmidt and Lipson’s AFP regressor (Schmidt & Lipson, 2010), La Cava et al.’s EHC regressor (La Cava et al., 2014), Burlacu et al.’s Operon (Burlacu et al., 2020), and Meidani et al.’s SNIP (Meidani et al., 2023). The baselines were run using their default settings in SRBench. Hyperparameter tuning was skipped. This section clarifies these configurations.

McConaghy’s FFX Regressor (McConaghy, 2011) has no tunable parameters. de Franca and Aldeia’s ITEA (Aldeia & de França, 2022) has a default population size of 1000 and 5000 generations. The minimum and maximum exponents of the interactions are -1 and $+1$. The minimum and maximum number of terms in an expression is 2. The number of non-zero exponents in each term of the initial population is 1. The transformation functions supported are the identity function, the hyperbolic tangent function, the sine function, the cosine function, the logarithmic function, the exponential function, and the square root function. Virgolin et al.’s GP-GOMEA (Virgolin & Bosman, 2022) sets a budget of 500000 evaluations over a single run (as opposed to interleaved multiple runs). It uses the default functions for addition, subtraction, multiplication and division. The initialized maximum tree height is 4. La Cava et al.’s EPLEX (La Cava et al., 2019) uses a selection mechanism called `epsilon_lexicase`. It uses 500 generations, with 500 individuals in the genetic programming population. The genetic programming algorithm also forces survival of the best individual in the population. The maximum number of nodes at the initialization of the genetic program is 20, and increases to 64 for the rest of the program. Schmidt and Lipson’s AFP regressor (Schmidt & Lipson, 2010) uses the parametric hill climber algorithm. It uses 250 generations, with 1000 individuals in the genetic programming population. The genetic programming algorithm also forces survival of the best individual in the population. The maximum number of nodes at the initialization of the genetic program is 20, and increases to 64 for the rest of the program. La Cava et al.’s EHC regressor (La Cava et al., 2014) uses an epigenetic hill climbing algorithm. It uses 100 generations, with 1000 individuals in the genetic programming population. The genetic programming algorithm also forces the survival of the best individual in the population. The maximum number of nodes at the initialization of the genetic program is 20, and increases to 64 for the rest of the program. Burlacu et al.’s Operon (Burlacu et al., 2020) uses five local iterations with 10000 generations. It sets a maximum evaluation budget of 5×10^5 . The population size is 500. The default allowed symbols are addition, subtraction, multiplication, and division. These baselines were run from the SRBench github, maintained by Cava Lab, at <https://github.com/cavalab/srbench>. The number of trials and number of jobs are set to 1. The code is run locally as opposed to on LPC. The time limit is set depending on the time budget of the experiment. The `skip_tuning` hyperparameter is used to skip tuning. Meidani et al.’s SNIP (Meidani et al., 2023) uses the grey wolf optimizer with a population of 50. It uses a beam search size of 2 with a stopping criterion of $R^2 = 0.9998$. The maximum iteration budget is changed from 80 to a time budget depending on the experiment.

B.3.2. README CONFIGURATIONS

In the rapid setting, each algorithm was allocated a runtime of **30 seconds** to improve the R^2 training fit. For README, we used the Grey Wolf Optimizer (GWO) to explore the neighborhood around a base latent encoding. A total of 70 additional candidates were generated by adding scaled Gaussian noise to the base encoding, resulting in 71 latent vectors (1 original + 70 perturbed). In each iteration, all 71 candidates were decoded into symbolic expressions and evaluated on the test set (see Algorithm 1). Their scores were then used to perform a population update using the GWO algorithm. This full decoding strategy was feasible due to the relatively generous runtime budget.

In the ultra rapid setting, each algorithm was given a strict time limit of **10 seconds** to maximize the R^2 training fit. We enabled README to operate effectively within this constraint by introducing a novel GWOBO algorithm, detailed in Appendix D. GWOBO combines Grey Wolf Optimization with Bayesian optimization by scoring candidate embeddings using a Gaussian Process with an Upper Confidence Bound (UCB) acquisition function. Only the top 3 candidates are selected for decoding, significantly reducing computational cost compared to decoding all candidates with pure GWO. As shown in Appendix A.2, GWOBO achieves higher accuracy than the standard GWO within the tight time budget.

B.4. Detailed Evaluation of Accuracy, Parsimony, and Runtime

B.4.1. EXTENDED RESULTS FOR TABLE 1

Table 1 compares README and SNIP on the Physics-Informed dataset. Here, we provide similar comparisons on the Strogatz and CP3-Bench datasets. Despite using 60 times less pretraining data, README consistently outperforms SNIP across all three datasets, highlighting the effectiveness of our README framework.

Model	Pretraining Data	Mean R^2_{Test}
README	~1 million pairs	0.880 ± 0.018
SNIP	~60 million pairs	0.791 ± 0.069

Table 6. Comparison on the Strogatz dataset.

Model	Pretraining Data	Mean R^2_{Test}
README	~1 million pairs	0.933 ± 0.018
SNIP	~60 million pairs	0.923 ± 0.017

Table 7. Comparison on the CP3-Bench dataset.

B.4.2. EXTENDED RESULTS FOR TABLE 2 AND FIGURE 5

We provide a detailed explanation of the evaluation metrics used—Mean R^2_{Test} for accuracy, Mean Equation Length for parsimony, and Mean Time for runtime—for all algorithms referenced in Table 2.

To compute the R^2_{Test} score, we follow the standard SRBench setup: each dataset is first split into 75% training and 25% testing. From the training split, we randomly subsample 200 points for training, and the fitted equation is evaluated on the held-out 25% to compute R^2_{Test} . The R^2 score measures the proportion of variance in the dependent variable that is predictable from the independent variables, with a score of 1.0 indicating perfect fit and lower values reflecting greater error.

Parsimony is measured by Equation Length, defined as the number of nodes in the expression tree. Each constant, variable, and operator (e.g., "3.5", "x_2", and "mul") counts as a single node. Shorter equations are considered more parsimonious and are generally easier to interpret and more robust to noise.

Time refers to the duration each algorithm takes to produce a final equation, excluding all data preprocessing and splitting steps. It reflects the computational efficiency of the symbolic regression process itself.

The term "Mean" in all metrics indicates that values are averaged across all problems and five random seeds for each dataset. The random seed affects both the subsampling of 200 training points and the train-test split. The error bars shown represent the standard deviation across these five seeds.

As shown in Figure 5, **README** consistently lies on the first Pareto frontier across all three datasets, achieving the highest accuracy while remaining reasonably parsimonious. ITEA is also Pareto-optimal as it achieves lower accuracy but with smaller equation length, making it non-dominated in the accuracy-parsimony space. Tables 8, 9, and 10 report the detailed results for these metrics.

Algorithm	Mean R^2 Test Score	Mean Equation Length	Mean Time (s)
README	0.880 ± 0.018	20.99 ± 1.19	24.95 ± 1.77
OperonRegressor	0.849 ± 0.053	60.46 ± 2.25	10.56 ± 0.11
GPGOMEAREgressor	0.849 ± 0.054	29.27 ± 3.46	39.63 ± 2.66
FFXRegressor	0.822 ± 0.063	198.43 ± 37.82	6.65 ± 9.41
SNIP	0.791 ± 0.069	22.97 ± 1.08	26.42 ± 0.44
AFPREgressor	0.742 ± 0.043	38.00 ± 6.08	20.02 ± 0.77
ITEAREgressor	0.736 ± 0.052	12.83 ± 0.27	15.67 ± 0.34
DSRRegressor	0.693 ± 0.067	19.16 ± 3.49	138.12 ± 3.58
EHCREgressor	0.644 ± 0.042	19.13 ± 1.57	12.44 ± 0.30
EPLEXRegressor	0.591 ± 0.076	45.09 ± 3.69	60.74 ± 1.18

Table 8. Performance comparison on the Strogatz Dataset. Bolded entries indicate the best R^2 score and the most parsimonious model.

Algorithm	Mean R^2 Test Score	Mean Equation Length	Mean Time (s)
README	0.933 ± 0.018	23.34 ± 0.11	23.04 ± 1.07
FFXRegressor	0.930 ± 0.006	172.45 ± 15.37	2.78 ± 0.16
ITEAREgressor	0.918 ± 0.002	16.08 ± 0.56	16.14 ± 0.29
GPGOMEAREgressor	0.910 ± 0.004	29.51 ± 0.83	17.35 ± 5.20
AFPREgressor	0.887 ± 0.003	39.63 ± 3.76	25.13 ± 0.63
EPLEXRegressor	0.885 ± 0.021	51.96 ± 1.45	68.70 ± 3.07
SNIP	0.883 ± 0.091	27.25 ± 0.92	29.19 ± 0.26
OperonRegressor	0.882 ± 0.013	77.19 ± 1.12	11.15 ± 0.06
EHCREgressor	0.862 ± 0.006	21.31 ± 1.14	13.73 ± 0.30
DSRRegressor	0.803 ± 0.005	17.07 ± 1.61	136.84 ± 0.59

Table 9. Performance comparison on the CP3-Bench Dataset. Bolded entries indicate the best R^2 score and the most parsimonious model.

Algorithm	Mean R^2 Test Score	Mean Equation Length	Mean Time (s)
README	0.984 ± 0.004	23.76 ± 2.35	28.91 ± 0.12
GPGOMEAREgressor	0.930 ± 0.003	29.85 ± 0.38	21.26 ± 3.54
SNIP	0.923 ± 0.017	24.52 ± 0.45	28.08 ± 0.24
FFXRegressor	0.930 ± 0.007	198.41 ± 15.47	3.09 ± 0.07
AFPREgressor	0.912 ± 0.007	42.67 ± 2.84	26.06 ± 1.48
DSRRegressor	0.891 ± 0.041	30.00 ± 0.00	135.32 ± 1.61
EPLEXRegressor	0.869 ± 0.026	47.40 ± 3.53	74.29 ± 0.32
EHCREgressor	0.818 ± 0.033	25.53 ± 1.75	13.94 ± 0.23
OperonRegressor	0.806 ± 0.082	80.98 ± 1.36	14.47 ± 0.15
ITEAREgressor	0.777 ± 0.017	11.67 ± 0.40	16.77 ± 0.46

Table 10. Performance comparison on the Physics Informed Dataset. Bolded entries indicate the best R^2 score and the most parsimonious model.

B.5. Latent Space Analysis

For Section 4.2, the equations are classified into the following family types.

- periodic: contains sin and/or cos, and is periodic. E.g. $\sin(x_1) + \cos(x_2)$, $\cos(x_1) \cos(x_2)$.
- exp: contains exp only. E.g. $\exp(x_1) + \exp(x_2)$.

-
- polynomial: contains $+$, \times , and/or the power function. E.g. $x_1^2 + x_2$.
 - inv: contains the inverse function. E.g. x_1^{-1} .
 - periodic_exp: contains \sin , \cos and/or \exp . E.g. $\sin(x_1) + \exp(x_2)$.
 - periodic_polynomial: contains \sin , \cos , $+$, \times , and/or power. E.g. $\sin(x_1^2) \times x_2$.
 - exp_polynomial: contains \exp , $+$, \times , and/or power. E.g. $\exp(x_1 + x_2) + x_2^2$.

Examples provided are for a 2-dimensional input, where $x = [x_1, x_2]^\top$.

In Figure 4, the t-SNE plots were generated with 512 equations and the t-SNE perplexity parameter was set as 30. We noted that ‘periodic_polynomial’ embeddings seem to be interpolated between ‘periodic’ and ‘polynomial’ embeddings for trained image embeddings, while this characteristic was not obvious for text embeddings from MathBERT. Thus, we hypothesize that some relationships between equation families may be more learnable when visualized in plots, as they may appear different symbolically but can be more easily captured by the image encoder through visual patterns.

B.6. Video Analytics Pipeline

We implemented two pipelines to extract the position of objects from video frames: one based on template matching and the other using YOLOv8 for object detection. While these methods are applicable to any video source, we demonstrate their effectiveness using two specific examples. Template matching using the Tracker Software is showcased on a video of a pendulum swinging, and bounding box identification with YOLOv8 is demonstrated on a video of a ping pong ball dropping. If the object of interest is not supported by YOLOv8’s predefined classes, template matching offers a flexible alternative for tracking custom objects.

B.6.1. TRACKER SOFTWARE

To extract motion data from real-world footage, we utilized the Tracker software (Brown, 2024), which is based on *template matching*. In our experiment, we used a publicly available video of a simple pendulum in motion to demonstrate the effectiveness of this pipeline.² Template matching works by selecting a region of interest and then scanning each subsequent frame to find the region that most closely resembles the original template. By identifying the best match in each frame, the software tracks the object’s position over time.

Using this method, we tracked the pendulum bob across frames and extracted a total of 311 data points representing its x and y coordinates over time. The resulting $(x, y, time)$ coordinates were exported to a CSV file and used as input to our README model for symbolic regression.

B.6.2. YOLOv8

For automated tracking in a separate experiment, we used the pre-trained YOLOv8n model (?) to detect and track the trajectory of a falling ping pong ball from a publicly available YouTube video.³ The video was trimmed to a 10-second segment from 348 to 358 seconds (5 minutes and 48 seconds to 5 minutes and 58 seconds). YOLOv8 identified the object of interest, classified as a `sports ball`, and recorded the center coordinates of its bounding box in each frame. In total, 85 data points were extracted.

We used the `yolov8n.pt` model with the detection class set to `sports ball` over the specified time window. The resulting $(x, y, time)$ coordinates were exported to a CSV file and used as input to our README model for symbolic regression.

C. Related works

Regression-based models Regression-based approaches, as their name suggests, use regression on a fixed basis to find an accurate representation of the input and output data of a system. Regression-based approaches tend to focus on using regularization to find a parsimonious basis (Brunton et al., 2016; Champion et al., 2019). However, they predefine the

²Video source: https://www.youtube.com/watch?v=02w91Sii_Hs

³Video source: <https://www.youtube.com/watch?v=pZ1Y10121Fs>

structure of the equation they aim to find, reducing the SR problem into one solving a system of linear equations (Makke & Chawla, 2024). This makes regression-based approaches very fast, but limits the generalizability of regression-based approaches. For example, McConaghy’s Fast Function Extraction (McConaghy, 2011) uses regularization to prune the search space of functions, and is a fast and deterministic algorithm for solving symbolic regression algorithms.

Genetic programming-based models These include seminal works by Koza (Koza, 1989; 1990; 1994), which represent each approximation of an unknown equation as a genetic program with a tree-like data structure, with traits (or nodes in the tree) representing functions or operations and variables representing real numbers. The fitness of each genetic program is its prediction error. Fitter genetic programs undergo a set of transition rules comprising selection, crossover, and mutation to find the optimal equation form iteratively. Genetic programming algorithms perform well in SR tasks as the transition rules allow for large variations in the population to adequately explore the search space.

Recent SR algorithms that use genetic programming to tackle common issues such as coefficient optimization include for example, de Franca and Aldeia’s Interaction-Transformation Evolutionary Algorithm (ITEA) (Aldeia & de França, 2022) which uses a search space which contains only mathematical expressions described as an affine combination of nonlinear transformations of different interactions between the original variables. ITEA then uses a mutation-based evolutionary algorithm to search for the optimal coefficients to express a linear relationship between the nonlinear transformations and the target variable. Likewise, Virgolin et al.’s GP-GOMEA (Virgolin & Bosman, 2022) searches for optimal values of coefficients by estimating interdependencies between model components and using this information to cross-over interdependent components en block, to preserve their concerted action improving mutation in genetic programming, and La Cava et al.’s EPLEXRegressor (La Cava et al., 2019) uses lexicase selection as a parent selection method that considers training cases individually, rather than in aggregate, to select elite parents for mutation.

Genetic programs may greedily mimic nuances of the unknown equation (Smits & Kotanchek, 2005), limiting generalisability. David Goldberg (Goldberg, 1989) therefore proposed to use Pareto optimization to balance the objectives of fit and parsimony in SR. At each iteration, the fittest genetic programmes lie on the non-dominated Pareto-frontier. Other works that use the Pareto frontier to evolve a population include Schmidt and Lipson’s age-fitness Pareto (AFP) optimization regressor (Schmidt & Lipson, 2010).

Lastly, some genetic algorithms explicitly minimize a target. For example, Burlacu et al.’s Operon (Burlacu et al., 2020), a genetic programming symbolic regression algorithm written in C++, minimizes speed, while La Cava et al.’s epigenetic hill climbing symbolic regression algorithm (EHC) (La Cava et al., 2014) minimize complexity of the equation and computational cost.

However, the transition rules of genetic programming algorithms mean that they are by design highly sensitive to hyperparameters and do not scale well to high-dimensional data (Petersen et al., 2021). This motivates the study of other types of symbolic regression algorithms.

Foundation model-based models. Deep learning algorithms are a recent advancement in the field of symbolic regression. An early example of a deep learning approach to symbolic regression is Petersen’s Deep Symbolic Regression (Petersen et al., 2021) which uses a recurrent neural network to emit a distribution over tractable mathematical expressions and employ a novel risk-seeking policy gradient to train the network to generate better-fitting expressions. Deep learning approaches have evolved following the progress in the field. Likewise, a few works (Guimerà et al., 2020) have also looked into adopting a Bayesian approach to symbolic regression, where a prior can be established based on a pool of past expressions which incorporates some domain knowledge, as well as naturally encode some balance between model complexity represented by the prior and data fit. Radford et al. (Radford et al., 2021a), in 2021, proposed multimodal architectures trained on symbolic expressions and numerical data to speed up genetic programming-based symbolic regression methods, while Kamienny et al. (Kamienny et al., 2022) introduced the use of transformers for symbolic regression to directly predict symbolic equations and Biggio et al. (Biggio et al., 2021) popularized the use of pre-trained transformers for symbolic regression. As the proposed architectures have grown bigger, the amount of data required to train these models has also grown. Pre-trained transformers start with a robust understanding of general symbolic patterns and syntax, and can be fine tuned to specific tasks such as for SR with less task-specific data. Since the model has already learned generic features of mathematical equations, the optimization process during fine-tuning focuses on symbolic regression-specific nuances. This significantly reduces training time and computational costs, and the pre-trained transformers converge faster during evaluation as they were trained on richer datasets. Meidani et al.’s SNIP (Meidani et al., 2023) proposed pre-training numeric and symbolic encoders jointly to produce a structured latent space that could be used for cross-domain tasks such as symbolic regression. In our work, we build on all these foundational works by (1) introducing an informative, compressed image representation

of numerical data that can be efficiently used in our framework for symbolic regression, (2) using pre-trained image and text encoders along with customized components to significantly reduce the training data and resources needed, (3) a novel method GWOBO that enables symbolic regression at the ultra-rapid setting ($< 10s$) that have not been explored before in past works, and (4) achieving significantly performance improvements over past methods.

D. GWOBO: Grey Wolf Optimizer with Bayesian Optimization

For our novel GWOBO algorithm, we begin by generating an initial population of latent encodings by perturbing a base encoding with scaled Gaussian noise, similar to the initialization used in the standard Grey Wolf Optimizer (GWO). Each candidate in this set is decoded using Algorithm 1 to obtain its symbolic expression and corresponding score.

To select the next embeddings to perform decoding on, we attempt to construct a Gaussian process (GP) surrogate to predict the score for an embedding. For each latent embedding in the original 512-dimensional space that we have already decoded, we compute its difference from the base encoding, then perform some Johnson-Lindenstrauss transformation (Johnson & Lindenstrauss, 1984) to project the difference down to a smaller 20-dimensional space. Given the transformed lower-dimensional vectors and their corresponding actual scores, we fit a GP with radial basis function (RBF) kernel with automatic relevance determination (ARD) and appropriate input and output normalization. The GP surrogate can then be used to compute the upper confidence bound (UCB) score for each latent embedding, which can be seen as an optimistic estimate of the actual score.

To reduce the high cost of decoding, in subsequent iterations only the **top 3 candidates** as ranked by UCB scores are decoded using Algorithm 1. Their actual scores are then used to update the GP model, while the remaining candidates are assigned surrogate scores predicted by the GP. These updated scores are then used by the GWO algorithm to perform a population update, guiding the search toward more promising regions. This allows reduction in the running time since the UCB scores can be computed much more rapidly compared to the true score while being reasonably accurate.

This loop continues until either the runtime exceeds the 10-second time budget or a sufficiently good score (R^2 train > 0.9998) is achieved. This hybrid approach, combining population-based candidate generation and optimization with GP-based surrogate modeling, enables strong performance while ensuring the entire inference process completes within the ultra-rapid 10-second runtime constraint.

We also conducted an ablation study on the number of initial candidates used to fit the GP, which is detailed in Appendix A.2.

E. Limitations and Broader Impacts

We introduced README, a framework for symbolic regression that leverages image representations of numerical and pre-trained multimodal foundation models for efficient learning. Compared to other foundation-model-based approaches, README requires significantly less training data and time, and can benefit from capability advancement of open-sourced text and image encoders. While a foundation model approach may allow for faster inference time and better performance, and is also capable of making full use of modern hardware such as GPUs, there are settings where this approach is less suitable, e.g., in Internet of Things (IoT) deployment settings where the hardware is constrained to lightweight devices/CPU. As most works on symbolic regression do not consider very high-dimensional datasets, we have similarly only considered up to 10-dimensional problems. We leave it to future works to examine the performance of README in higher-dimensional problems.

As README allows users to rapidly identify equations to describe data, it has the potential to support applications such as interactive/iterative scenarios such as adaptive scientific experimental and close to real-time decision making. README might also be used as components in AI/machine learning systems where interpretability in terms of symbolic equations would be useful. While we expect that the majority of such applications will lead to societal benefits, there may be malicious actors who might come up with applications that are to the detriment of society – general regulations and efforts to prevent such abuse of AI/machine learning tools are needed.

F. Directly using multi-modal large language models for symbolic regression

To demonstrate that our approach is effective and necessary for symbolic regression, we used GPT-4o out of the box for a naive comparison. On the physics-informed dataset with real-world measurements, our README model achieved a Mean

R^2 Test score of 0.958 under the ultra rapid 10 second setting with 71 candidates as detailed in A.2, while GPT-4o achieved only 0.015. This highlights that GPT-4o, used out of the box, is unable to perform symbolic regression meaningfully.

For GPT-4o, we followed the same 5-seed train-test splits described in Appendix B.2. We provided 200 subsampled data points and explicitly informed GPT-4o of the number of input dimensions, instructing it to return a valid NumPy expression that fits the data. Notably, GPT-4o was given an advantage by being explicitly told to produce equations with the correct number of input dimensions, a constraint that was not enforced for README. Table 12 shows the prompt and example outputs from GPT-4o on the Physics Informed Dataset (seed 23654), illustrating that although the expressions are syntactically valid, they fail to fit the data well.

Algorithm	Mean R^2 Test Score	Mean Equation Length
README (71 initial candidates)	0.958 \pm 0.027	16.67 \pm 1.57
GPT-4o	0.015 \pm 0.023	10.46 \pm 0.67

Table 11. Comparison between README and GPT-4o on Physics Informed dataset for ultra-rapid setting (10 seconds). Bold indicates the better score per metric.

Prompt	GPT-4o Discovered Equations for seed 23654
<p>You are given training data with input dimension = {num_dim}.</p> <p>The X array contains input points, and the Y array contains the corresponding target values.</p> <p>X = {X_to_fit}</p> <p>Y = {Y_to_fit}</p> <p>Provide a single NumPy compatible expression f(x) that takes an (n, {num_dim}) array x and returns an (n,) array of predictions.</p> <p>Reply with only the expression itself (e.g., np.sin(x[:,0]) + 0.5*x[:,1]), without any explanation or quotes.</p>	0.05 * np.sin(2 * np.pi * x[:, 0]) + 0.1
	0.5 * np.sin(x[:, 0]) + 0.1
	0.5 * np.cos(x[:, 0]) + 0.5
	0.5 * np.cos(x[:, 0]) + 0.5
	0.5 * np.tanh(x[:, 0]) + 0.3 * np.tanh(x[:, 1]) + 0.5
	0.5 * np.cos(x[:, 0]) + 0.1
	0.05 + 0.03 * np.tanh(2 * x[:, 0])
	np.where(x[:, 1] > 0, 0.001, np.exp(x[:, 0]) / (1 + np.exp(x[:, 0])))
	1.75 - 0.1 * np.tanh(x[:, 0])
	-0.1 * np.tanh(x[:, 0]) - 0.05
	0.5 * np.sin(x[:, 0]) - 0.5 * x[:, 0] - 0.5
	np.where(x[:, 1] > 0, 0.00005, np.exp(-x[:, 0]**2))
	0.05 + 0.02 * np.sin(2 * np.pi * x[:, 0])
	0.05 * np.sin(2 * np.pi * x[:, 0]) + 0.1 * np.cos(np.pi * x[:, 0]) + 0.1
	np.where(x[:, 1] > 0, 0.001, np.exp(-x[:, 0]**2))
	np.where(x[:, 1] > 0, 0.5 * (x[:, 0] + 1.5)**2, np.exp(-x[:, 0]**2))
	0.05 * np.sin(2 * np.pi * x[:, 0]) + 0.1 * np.cos(np.pi * x[:, 0]) + 0.1
	-100 * x[:, 0] + 700
	0.5 * np.cos(x[:, 0]) + 0.5
	np.exp(-x[:, 0]**2)

Table 12. Prompt and discovered expressions for 19 symbolic regression problems in the Physics Informed dataset.