

In-Context Symbolic Regression: Leveraging Language Models for Function Discovery

Matteo Merler and Nicola Dainese and Katsiaryna Haitsiukevich

Department of Computer Science

Aalto University

matteo.merler@aalto.fi

Abstract

Symbolic Regression (SR) is a task which aims to extract the mathematical expression underlying a set of empirical observations. Transformer-based methods trained on SR datasets detain the current state-of-the-art in this task, while the application of Large Language Models (LLMs) to SR remains unexplored. This work investigates the integration of pre-trained LLMs into the SR pipeline, utilizing an approach that iteratively refines a functional form based on the prediction error it achieves on the observation set, until it reaches convergence. Our method leverages LLMs to propose an initial set of possible functions based on the observations, exploiting their strong pre-training prior. These functions are then iteratively refined by the model itself and by an external optimizer for their coefficients. The process is repeated until the results are satisfactory. We then analyze Vision-Language Models in this context, exploring the inclusion of plots as visual inputs to aid the optimization process. Our findings reveal that LLMs are able to successfully recover good symbolic equations that fit the given data, outperforming SR baselines based on Genetic Programming, with the addition of images in the input showing promising results for the most complex benchmarks.

1 Introduction

Classical Machine Learning regression methods can be divided into two broad categories: statistical methods, which simply learn an implicit statistical model of the relationship between the observations typically represented with a black-box model, and rule-based methods, which instead attempt to extract an explainable set of rules that explicitly model the transformation between the inputs and outputs (Lample and Charton, 2019). Symbolic Regression (SR) is a particular subset of the latter category which searches the set of all possible explicit mathematical expressions to find an equation

that best fits the given set of observations, instead of simply predicting the value for unobserved data-points. This has the clear advantage of explainability, as well as a potential for better generalization if the trend holds outside of the observed data.

The most common approach for SR algorithms is Genetic Programming (Willis et al., 1997) (GP), combining fundamental blocks for mathematical expressions (e.g. basic operators, trigonometric functions, etc.) into more complex formulas using strategies borrowed by evolutionary biology, such as mutations and fitness. Out of the 14 methods tested in a popular Symbolic Regression Benchmark (Cava et al., 2021), 10 of them were based on GP.

The recent success of Transformer based Deep Learning models, first introduced by Vaswani et al. (2017), has revolutionized many fields, most notably Natural Language Processing (OpenAI, 2023; Brown et al., 2020; Touvron et al., 2023; Anil et al., 2023) and Computer Vision (Dosovitskiy et al., 2021) amongst many. Large Language Models (LLMs) in particular have been proven to possess unprecedented reasoning and generalization abilities, as well as the newly emergent capacity for In-Context Learning (ICL) (Dong et al., 2023), which refers to a prompting technique that includes demonstrations of a specific task in natural language into the input given to the LLM, which enables it to address new, potentially unseen, tasks. The emergence of ICL and few-shot prompting was first observed by Brown et al. (2020) and has since become a staple prompting technique. With the help of ICL, these models can be leveraged for a wide range of different tasks, suggesting a potential use case for Symbolic Regression. Many Transformer-based methods have been proposed, but to the best of our knowledge, the direct use of pre-trained LLMs for this task has not been explored yet.

In this paper, we examine the integration of

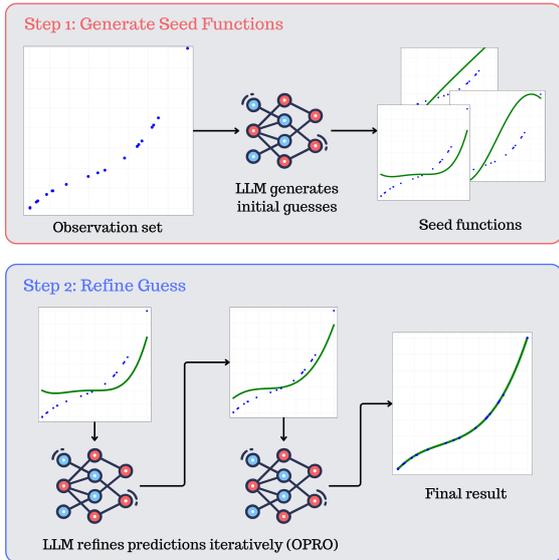


Figure 1: **High level overview of the approach.** Given an initial set of observations, we prompt the LLM or VLM to generate multiple initial guesses (seeds) of the true function that generated the observations. We then iteratively refine our guesses with the Optimization by Prompting (OPRO) method from Yang et al. (2023). The model only produces the functional form of a function, while the coefficients are fitted using non-linear least squares optimization.

LLMs into a Symbolic Regression pipeline, with the aim to use them as the main generator for new functions. Inspired by the Optimization by Prompting (OPRO) approach presented by Yang et al. (2023), we implement a similar system for Symbolic Regression, presented in Section 4. The LLM receives a number of previously tested functions and their scores on the dataset, and is then tasked to generate a new function that could be a better fit for the given observations. This approach is repeated until the error is low enough and we consider the algorithm to have converged to a solution. This approach has been proven to work on classical optimization problems, such as Linear Regression and on the challenging Travelling Salesman Problem.

We then investigate extending this approach to take advantage of Vision-Language Models (VLMs). This provides a richer context for this optimization task: plots are a natural way for humans to communicate about mathematical concepts (for example, a scatter plot is more informative compared to a string of values), arguably more so than text, as they can convey additional information in an intuitive way, such as the shape or general trend

of the data. We discuss the impact of this extension in Section 5, as well as discuss its limitations and potential in Section 6.

2 Related Work

Symbolic Regression. GP has traditionally formed the backbone for SR methods. A typical implementation can be seen in `gplearn`¹, which follows the classical definition of a GP approach by Koza and Poli (2005). From an initial population, an iterative tournament is played where functions with the lowest fitness are eliminated, before the best performing alternative can "reproduce" with some random mutation. More sophisticated approaches can make use of Pareto optimization (Smits and Kotanchek, 2005; Schmidt and Lipson, 2011), which take advantage of the Pareto front for more efficient convergence. A more recent GP method is GP-GOMEA by Virgolin et al. (2021), which takes advantage of a statistical model updated over time.

More recently, Deep Learning methods have also been proposed for symbolic regression. Petersen et al. (2021) use a Recurrent Neural Network with a risk-seeking policy to perform a hierarchical search over the space of possible smaller symbolic formulas and Cranmer et al. (2020) use Graph Neural Networks to extract physical equations. Many models based on the Transformer architecture have also been presented ever since their introduction and rise in popularity: Lample and Charton (2019) train a seq2seq model to directly solve integrals and differential equations and then followed up with Kamienny et al. (2022) to train a Transformer model specifically for Symbolic Regression. In parallel, Valipour et al. (2021) and Biggio et al. (2021) also trained different generative Transformer models for this task and later Li et al. (2023d) trained a model using joint supervised learning with a novel objective function, while Shojaee et al. (2023) present a method integrating Transformers with Monte Carlo Tree Search to guide the generation of the symbolic formulas. An example of image-based Symbolic Regression (extended to two-dimensional functions) was presented by Li et al. (2022a), in which the task is formulated as a reconstruction problem, with an autoregressive Vision Transformer model encoding and decoding functions, shown as heatmaps.

¹<https://gplearn.readthedocs.io>

Reasoning with Large Language Models. As LLMs form the backbone of the method presented in this work, we rely entirely on their reasoning capabilities, such as ICL (Brown et al., 2020), to explore the solution space. More specifically, we are interested in the mathematical understanding of LLMs, for which Liu et al. (2023d) provide a recent survey. Amongst many, Li et al. (2022c) introduces LEMMA, a fine-tuned model based on CodeLlama (Rozière et al., 2023) for theorem proving, and Luo et al. (2023) instead fine-tune Llama2 using a set of mathematical reasoning instructions generated with a novel Reinforcement Learning via Evol-Instruct Feedback approach. Mirchandani et al. (2023) show that LLMs are able to recognize patterns from in-context examples and can extrapolate them to complete related tasks in the input. Similarly, Gruver et al. (2023) find that LLMs can extrapolate zero-shot the pattern from a timeseries (although they do not extract any functional representation). Furthermore, Fu et al. (2023) present a study in which they find that Transformer models can learn higher-order optimization methods (similar to Newton’s method).

Vision-Language Models. VLMs have gained traction after Radford et al. (2021) introduced CLIP, aligning text and image representations using a contrastive objective. Various foundation models have been proposed, such as FLAVA (Singh et al., 2022), LLaVa (Liu et al., 2023b,a, 2024), Flamingo (Alayrac et al., 2022), OTTER (Li et al., 2023b,a), Fuyu (Bavishi et al., 2023) and more recently OpenAI’s GPT4’s vision extension. A thorough survey of VLM techniques and tasks was performed recently by Zhang et al. (2023). Typically, a VLM can be built on top of a pre-trained LLM, which is then paired with an image embedding network that can transfer the image into the same token space used by the model, attempting to keep semantical similarity. This approach is employed for example by BLIP (Li et al., 2022b) and its successor BLIP2 (Li et al., 2023c). Moreover, these models typically can only consume images as input, but are unable to generate them as an answer, but the general framework can be enhanced with methods for text-to-image generation, such as DALL-E (Ramesh et al., 2021, 2022) and GILL (Koh et al., 2023).

3 Background

Optimization by Prompting. The OPRO framework was first introduced by Yang et al. (2023), with the initial goal of prompt optimization. The paper builds on previous work by Mirchandani et al. (2023), which shows the potential use of LLMs to complete and understand general patterns. Crucially, the authors also present experiments showing the method working on classical optimization problems (Linear Regression and TSP), which leads them to claim that OPRO can be used for classical optimization tasks. The key idea of OPRO is the use of a so-called meta-prompt, a higher level prompt which describes the task to be performed, as well as previous examples together with scores (evaluated externally). The LLM, with its pattern understanding capabilities, should then be able to extrapolate the trend in the examples using in-context learning, and propose a better alternative. The model’s output is then evaluated, and added to the meta-prompt if the score was good enough. This approach can then be iterated until convergence (i.e. when the score is good enough).

4 Method

To leverage the OPRO approach to perform Symbolic Regression, we need to design a meta-prompt suitable for the task and fill it with an initial set of functions $f_{t=0}$ (either hand-written or model-generated), together with a measure of their fitness on the observations X (the observations X themselves are also included). The goal is then to iteratively refine the set of functions $f_{t=i}, i \in [1, \dots, n]$, until a sufficiently low error score is obtained or a maximum number of iterations is reached. Note that since the context window has a finite fixed size, we can only keep the k best performing previous attempts in the prompt at the same time, where k is a design choice (for this study, we use $k = 5$). The full meta-prompt used in the experiment can be found in Appendix A.

Seed Functions. At the first timestep, $f_{t=0}$ is empty as there are no previous guesses from the model, thus an initial population of seed functions is required. Instead of relying on a fixed set of seed functions, which could be restrictive in general, we ask the model to generate the initial conditions (the prompt that can be found in Appendix A). This typically generates a complex and diverse set of functions, from which the model can explore and

refine its future predictions with OPRO. In our implementation we repeat this initial process up to 10 times, as some or all of the generated functions can be invalid for a given dataset if some of the points are undefined.

Error Score. The immediate choice for the score is using a simple Mean Squared Error (MSE) for every point $x \in X$. This frames the process as an optimization problem where the lower score is better, similar to the ones presented in Yang et al. (2023). Another option to obtain more bounded scores is using the coefficient of determination R^2 described in Section 5.1.1. Crucially, the LLM does not need to have an understanding of how the score is computed, but simply needs to know a lower score is better (or higher in the case of R^2). Ideally, the model will extrapolate the pattern in the previous trajectory, finding which changes in the functions improved the score without knowledge of how it was computed. For our implementation, we choose the MSE metric.

Parameter Fitting. We use the LLM only to generate generic functional forms (and thus task it with exploration), while using SciPy’s (Virtanen et al., 2020) implementation of Non-linear Least Squares (Kelley, 1999). This results in the model exploring a large space of functions in a short amount of time, allowing it to solve complex benchmarks. In our implementation, we repeat this process five times with different random initial values for the coefficients, to refine the results and avoid local minima, similar to Li et al. (2023d).

For other details about the OPRO implementation, we follow the original work. Specifically, we also sample multiple functions for every iteration in an attempt to improve stability and we experiment with a decreasing temperature parameter to balance exploration/exploitation (with a higher initial temperature encouraging the exploration of the underlying functional space, and a lower temperature at the later stages forcing smaller tweaks to the trajectory). To avoid crowding the input prompt with too much information, we limit the amount of training points that are included in written form to a certain threshold, arbitrarily set to 40. We discuss this issue in more detail in Section 6.1.

4.1 Vision-Language Extension

The input data used in the meta-prompt consists of the observations X , represented via a set of coordi-

nates (e.g. a pair x, y for a univariate function), as well as a trajectory composed of a set of the best the previously generated functions with their scores. Previous work by Gruver et al. (2023) shows that LLMs can understand a series of numerical points and process it autoregressively, but our extension also requires the model to reason on X together with the functions in the context; identifying why a function is not fitting the observations properly and where it could be improved is far more challenging when compared to simply predicting the next point in a timeseries. A natural way to enhance this information would be to provide a plot for the observations, as well as plots for the previously generated functions. This provides an intuitive way to analyze the performance of previous generations and identify areas where improvement is needed.

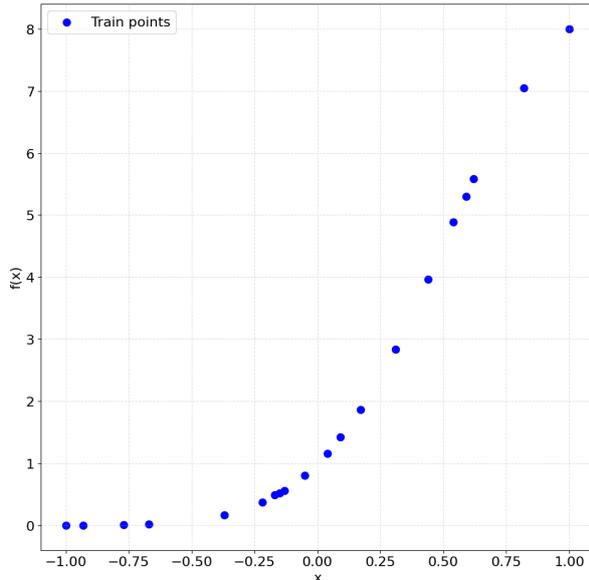
We thus extend the approach to use Vision-Language Models, by including either a plot of the observations (Figure 2a) when generating the seed functions or a plot of the best guess from the previous trajectory (Figure 2b) during the OPRO loop. We hypothesize an increase in accuracy and sample efficiency of the method, which we explore in Section 5.3. Crucially, the use of both vision and language as input comes with the restriction of dimensionality, as it is impossible to visualize inputs with more than two covariates. We discuss this limitation further in Section 6.1.

5 Experiments

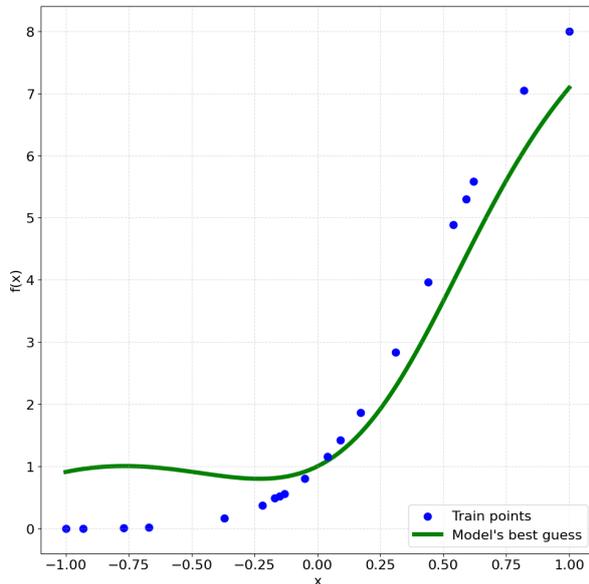
For our experimental setup, we employ two popular open source alternatives: Llama3 (Meta, 2024) as the text-only model and LLaVa-NeXT (Liu et al., 2024) as the vision-language model, which is based on Yi-34B (01.AI et al., 2024).

5.1 Benchmarks

Given the method’s inherent limitation in dimensionality due to the vision extension, we rely on classical SR benchmarks, containing functions in one and two dimensions, to evaluate our approach. We follow Li et al. (2023d) and choose four of the benchmarks tested in their work: Nguyen (Nguyen et al., 2011), Constant (a modified version of some of the Nguyen equations with different numerical values for the coefficients), Keijzer (Keijzer, 2003) and R (Krawiec and Pawlak, 2013). The formulas and ranges for both the training and testing points are taken from the original papers where available. The full list can be found in Appendix C.



(a) Scatter plot for the observations.



(b) Plot of the performance for a previous function.

Figure 2: Sample of plots used with the VLM.

5.1.1 Metrics

While we use MSE during the OPRO loop, we report the coefficient of determination R^2 (Glantz et al., 2017) to evaluate the quality of our method. This is a staple metric in SR and results in a more interpretable value: a function will get a positive score if it’s better compared to predicting the average value and will get a score of 1 for a perfect prediction. The coefficient is computed as follows:

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Benchmark	OPRO (Ours)	
	Llama3 R^2 (\uparrow)	gplearn R^2 (\uparrow)
Nguyen	0.9996	0.7524
Constant	0.9994	0.7441
Keijzer	0.9615	0.0917
R	0.9901	0.9697
Overall avg.	0.9877	0.6395

Table 1: Comparison between our proposed method and gplearn on several benchmarks.

where y_i is the ground truth value, \hat{y}_i is the predicted value and \bar{y} is the average of all y_i . We follow Li et al. (2023d) and Biggio et al. (2021) in removing the 5% worst predictions to ensure robustness against outliers.

5.2 Performance against GP methods

For a first exploratory study, we compare the performance of our method against the standard gplearn implementation, which can be seen as the simplest possible GP approach. As such, many hyperparameters can be tuned (population size, number of generations, etc.) which can drastically alter the accuracy and run-time of the method. Following the work by Li et al. (2023d), we keep the default values suggested by the library. The exact set of hyperparameters we use for both our method and gplearn can be found in Appendix B. We report the R^2 metric computed on a set of unseen testing points, which usually cover the same interval with a denser distribution, so that the resulting function can be evaluated on the whole point interval and we avoid reporting a metric that overfits on the training data. We repeat all experiments across five different random seeds, but we only randomize the set of input points (where they are random) only once, to ensure consistency across the different runs and methods. For this experiment we use a fixed number of OPRO iterations set to 50 and always perform a full run. The results can be seen in Table 1.

The OPRO approach using Llama3 is very robust, achieving very good scores across all benchmarks, while the simpler gplearn approach already fails to fit some of the more complex patterns found in the more difficult functions, like Keijzer. It’s worth noting that the version of gplearn we used is restricted to a limited vocabulary of functions that

are defined in the real domain and is thus lacking for example the square root and logarithm. LLMs on the other hand are free to use these functions, as they have some prior bias of where they are allowed. An LLM will still try to use a square root for negative values, but is often smart enough to ensure that the function it produces will be defined across the full domain that it has seen.

In general, most SR methods tend to be limited to a predefined vocabulary of operators and tokens, while LLMs can virtually explore any possible function and combination. An example of this is with the $x_1^{x_2}$ function: in Li et al. (2023d), the authors mention that they had not seen this token in the training data and were thus limited to finding appropriate approximations, while our approach can recover the exact expression.

Another observation is that the best function (or one close to it) is often found very early into the 50 iterations, and the rest of the time is spent trying minor changes, sometimes failing to improve the function all-together. Ideally a threshold can be set using R^2 on the training points: if this exceeds a certain predefined value (e.g. 0.9999) we can end the iteration process early. We didn’t use this in the experiments to explore how many iterations are necessary to converge on a good result, as well as ensuring all runs have the same amount of iterations, but we would suggest integrating this for practical applications.

5.3 Comparison of Text-Only and Vision-Language Models

To evaluate the effectiveness of the additional image input, we compare our method with a variant using the LLaVa-NeXT vision-language model. While the LLM underlying LLaVa-NeXT, Yi, is larger than Llama3 (34B parameters versus 8B), the latter actually outperforms Yi on several relevant benchmarks, such as GSM8K for mathematical reasoning, and thus the comparison should be fair.

We repeat the experiments across five different random seeds. Both models are performing a full set of 50 OPRO iterations.

The vision approach matches or slightly underperforms the text-only approach on the simpler Nguyen and Constant benchmarks, and shows more significant improvements on the more complex Keijzer benchmark, suggesting that plots can help where the underlying expression is harder to produce and there is a risk of underfitting. The re-

Benchmark	LLaVa-NeXT R^2 (\uparrow)	Llama3 R^2 (\uparrow)
Nguyen	0.9985	0.9996
Constant	0.9984	0.9994
Keijzer	0.9924	0.9615
R	0.9568	0.9901
Overall avg.	0.9865	0.9877

Table 2: Comparison between LLaVa-NeXT and Llama3 using OPRO.

Benchmark	50 iterations R^2 (\uparrow)	Seed functions only R^2 (\uparrow)
Nguyen	0.9985	0.9875
Constant	0.9984	0.9979
Keijzer	0.9924	0.8270
R	0.9568	0.9954
Overall avg.	0.9865	0.9520

Table 3: Comparison between the full OPRO approach and the seed function generation step.

sults validate the usefulness of the additional visual input, allowing a less capable model to closely match the most powerful open source text-only LLM currently available, while also confirming the assumption that advancements in the underlying LLM backbone will further improve the approach.

The lower score obtained by LLaVa-NeXT in the R benchmark is mostly due to the random sampling of the training points for the R3 function, which can be seen in Figures 11 and 12. The wide area with no training points makes it hard to extrapolate the exact pattern followed by the underlying function. This is an inherent weakness of most Symbolic Regression methods, and can be solved by sampling more training points.

5.4 Ablation Study

In this section we investigate the importance of the iterative refinement of the functions with the OPRO loop. As mentioned in Section 4, at the beginning of a run we ask the model to generate a set of diverse seed functions that can serve as a starting point for optimization (the prompt can be found in Appendix A). This could be sufficient to generate good enough guesses, bypassing the OPRO loop entirely. We investigate the goodness of the seed functions generated for all experiments, and report the results below in Table 3.

The results suggest that the initial seed functions generation step plays a key role in our approach, immediately producing guesses that are only slightly below the results found by the full approach for the simpler benchmarks, like Nguyen and Constant. It’s also interesting to observe a better fit in R: as discussed previously, R3 has a large area not covered by training points, and this could perhaps indicate that the OPRO loop is actually overfitting to the observations if continued for too long.

However, on the more complex Keijzer benchmark, the seed functions underfit the observations, as demonstrated in Figure 13 in Appendix D. This suggests that for more complex benchmarks and functions the initial generation step is essential to provide an educated guess, but the iterative refining process is still necessary to achieve a good quality solution. The types of complexity that cause the initial guess to underfit can be either a complex shape (for example in keijzer3 and keijzer4, the first two functions in the top row of Figure 13) or a lack of datapoints, especially for two-dimensional functions (like the full bottom row of Figure 13).

6 Discussion

The empirical results show that LLMs are at least capable or even adept at Symbolic Regression, outperforming a simple GP approach. This exposes yet another task that large foundation models can be leveraged for thanks to the use of specialized prompting techniques such as OPRO and shows promise for future work into the integration of LLMs and mathematical reasoning. The impact of visual information should also not be understated, as the additional plots provided as input help the model to turn an underfitted guess into a good expression for more complex benchmarks.

The method’s performance relies on the initial generation step for the seed functions, just as any optimization algorithm depends on its initial conditions. Self-generating the initial set can potentially lead to an ill-conditioned problem if the generated functions have a poor fitness score and the OPRO loop can get stuck refining a local minimum, unable to explore further. However, we find that in general self-generated initial conditions work reasonably well. Prior knowledge can also be incorporated by using a manually crafted set of functions if some information about the observations is known; we leave to future work further exploration on how to incorporate such knowledge in our method.

As the approach is model-agnostic, advances in both LLMs and VLMs will also improve the quality and speed of this method, with larger and more capable models such as GPT4 and its vision extension being prime candidates for future research into scaling. As this is an initial exploratory study into the capabilities of these models, we consider this to be a success, and leave these potential improvements and more in-depth experiments for future work. As it stands, this approach likely is not practical enough to be used over some more advanced and established Transformer-based methods, but we have hope that the potential showed in this study, combined with scaling and innovations in the field will lead to the method proposed in this paper matching or surpassing state-of-the-art performance on more complex benchmarks.

A further benefit of the proposed method is the lack of necessity for an expensive training phase specific to Symbolic Regression: leveraging pre-trained LLMs produces a method that can be quickly adapted to incorporate more powerful foundation models while being able to take advantage of the full vocabulary and reasoning skills offered by natural language, as opposed to models pre-trained directly for Symbolic Regression, which are limited to the range of tokens they were trained on. An interesting direction for future work could be relying even more on the unique reasoning abilities offered by LLMs by employing explicit Chain of Thought-like techniques, allowing the model to output even more well-informed guesses at every step.

6.1 Limitations

Although promising, the approach presented in this work still suffers from some key limitations that hold back its potential as a full Symbolic Regression method.

Dimensionality. If the chosen model relies on images, this fundamentally limits it on working on inputs that have at most two variables. A potential solution could be to include projections into each dimension as the input, but this approach can quickly grow out of control as the number of variables increases, and even then the additional information is still of questionable utility. Using a text-only LLM for higher dimensional inputs is possible, but this comes with its own challenges. As the number of variables grows, so does the length of the input points in the prompt, which will natu-

rally confuse the model and obfuscate the structure in the datapoints even further. We already observe a similar trend if too many input points are included in the prompt. Specifically fine-tuning an LLM on this kind of examples might show some improvement, but scaling this approach for higher dimensional problems remains a challenge.

Context window. The maximum number of tokens that an LLM can process as input is referred to as its context window. This is a limiting factor for our approach: the prompt is already quite long and many factors can add to it significantly. For example, as discussed in the previous paragraph, higher dimensional inputs will increase the length of the observations string, as well as leading to longer functions on average.

The version of LLaVa used for the experiments in this paper is built on top of a backbone trained on sequence lengths of 4k tokens, which can be extended at 16k tokens during inference, while Llama3 has a context window of 8k tokens. Nevertheless, we observe that if too many points are included in the prompt in written form, the model often starts just repeating new invented points as the output, due to the sheer number of inputs provided for some functions. This happens more often with three-dimensional inputs, as there are more numbers (and thus more tokens) for each training example. To avoid this, we find that limiting the amount of training points shown in the prompt to an arbitrary number (which we set at 40) is a good solution. The full range of training points can still be shown in the input images, to give an idea of the overall shape, and can be used in its entirety to optimize the coefficients. This is an advantage that the text-only Llama3 loses, as it is now only relying on the information in the prompt. However, we argue that this is likely to be solved as research in the field advances and context-window size increases. There are already commercially available LLMs able to process more than 100k tokens, such as GPT-4 turbo (OpenAI, 2023) and Claude 3 (Anthropic, 2024), which may solve this limitation entirely, although some previous research suggests that LLMs are not able to take advantage of the full context (Liu et al., 2023c).

7 Conclusion

We show that LLMs paired with the OPRO approach are able to successfully perform Symbolic Regression tasks on benchmarks with up to two

independent variables. Additionally, we show that VLMs can outperform LLMs at more complex variants of this task thanks to the additional information provided by plots, which give intuitive meaning and structure to the input. The proposed method outperforms a simple GP approach and produces more streamlined functions, proving the validity of this method.

7.1 Future Work

As this is merely a first exploration into this topic, much work remains to be done. The next steps for this project are experimenting with model scale, to test the impact of more powerful LLMs and VLMs, with the assumption being that the approach should greatly benefit from it. We also leave testing the method against current state-of-the-art methods for future work. For this, an exploration into extending the text-only method with Llama3 for higher-dimensional inputs would be necessary to solve a reasonable subset of tasks from the more contemporary benchmark SRBench (Cava et al., 2021).

As also discussed in Section 6.1, fine-tuning can also be an option to improve the model’s performance. We see two general options for this: either fine-tune a VLM to improve its visual mathematical reasoning skills on plots, or fine-tune a text-only LLM to improve the overall performance without relying on images. Both avenues show potential and could drastically improve on the results we show in this work.

References

- 01.AI, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. 2024. [Open foundation models by 01.AI](#). *Preprint*, arXiv:2403.04652.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, et al. 2022. Flamingo: a visual language model for few-shot learning. In *Advances in Neural Information Processing Systems*.
- Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng

- Chen, et al. 2023. [Palm 2 technical report](#). *Preprint*, arXiv:2305.10403.
- Anthropic. 2024. [Introducing the next generation of Claude](#).
- Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Sağnak Taşlılar. 2023. Introducing our multimodal models.
- Luca Biggio, Tommaso Bendinelli, Alexander Neitz, Aurelien Lucchi, and Giambattista Parascandolo. 2021. [Neural symbolic regression that scales](#). *Preprint*, arXiv:2106.06427.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- William La Cava, Patryk Orzechowski, Bogdan Burlacu, Fabrício Olivetti de França, Marco Virgolin, Ying Jin, Michael Kommenda, and Jason H. Moore. 2021. [Contemporary symbolic regression methods and their relative performance](#). *Preprint*, arXiv:2107.14351.
- Miles Cranmer, Alvaro Sanchez Gonzalez, Peter Battaglia, Rui Xu, Kyle Cranmer, David Spergel, and Shirley Ho. 2020. Discovering symbolic models from deep learning with inductive biases. In *Advances in Neural Information Processing Systems*, volume 33, pages 17429–17442. Curran Associates, Inc.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and Zhifang Sui. 2023. [A survey on in-context learning](#). *Preprint*, arXiv:2301.00234.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. [An image is worth 16x16 words: Transformers for image recognition at scale](#). *Preprint*, arXiv:2010.11929.
- Deqing Fu, Tian-Qi Chen, Robin Jia, and Vatsal Sharan. 2023. [Transformers learn higher-order optimization methods for in-context learning: A study with linear models](#). *Preprint*, arXiv:2310.17086.
- Stanton A. Glantz, Bryan K. Slinker, and Torsten B. Neilsands. 2017. *Dedication*. McGraw-Hill Education, New York, NY.
- Nate Gruver, Marc Anton Finzi, Shikai Qiu, and Andrew Gordon Wilson. 2023. Large language models are zero-shot time series forecasters. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Pierre-Alexandre Kamienny, Stéphane d’Ascoli, Guillaume Lample, and Francois Charton. 2022. End-to-end symbolic regression with transformers. In *Advances in Neural Information Processing Systems*.
- Maarten Keijzer. 2003. [Improving Symbolic Regression with Interval Arithmetic and Linear Scaling](#). In *Genetic Programming*, Lecture Notes in Computer Science, pages 70–82, Berlin, Heidelberg. Springer.
- C. T. Kelley. 1999. *Iterative Methods for Optimization*, pages 22–25. Society for Industrial and Applied Mathematics.
- Jing Yu Koh, Daniel Fried, and Ruslan Salakhutdinov. 2023. Generating images with multimodal language models. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- John R. Koza and Riccardo Poli. 2005. [Genetic programming](#). In Edmund K. Burke and Graham Kendall, editors, *Search Methodologies: Introductory Tutorials in Optimization and Decision Support Techniques*, pages 127–164. Springer US, Boston, MA.
- Krzysztof Krawiec and Tomasz Pawlak. 2013. [Approximating geometric crossover by semantic backpropagation](#). In *Proceedings of the 15th Annual Conference on Genetic and Evolutionary Computation, GECCO ’13*, page 941–948, New York, NY, USA. Association for Computing Machinery.
- Guillaume Lample and François Charton. 2019. [Deep learning for symbolic mathematics](#). *Preprint*, arXiv:1912.01412.
- Bo Li, Peiyuan Zhang, Jingkang Yang, Yuanhan Zhang, Fanyi Pu, and Ziwei Liu. 2023a. [Otterhd: A high-resolution multi-modality model](#). *Preprint*, arXiv:2311.04219.
- Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. 2023b. [Otter: A multi-modal model with in-context instruction tuning](#). *Preprint*, arXiv:2305.03726.
- Jiachen Li, Ye Yuan, and Hong-Bin Shen. 2022a. [Symbolic expression transformer: A computer vision approach for symbolic regression](#). *Preprint*, arXiv:2205.11798.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023c. [Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models](#). *Preprint*, arXiv:2301.12597.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022b. [Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation](#). *Preprint*, arXiv:2201.12086.
- Wenqiang Li, Weijun Li, Linjun Sun, Min Wu, Lina Yu, Jingyi Liu, Yanjie Li, and Songsong Tian. 2023d. [Transformer-based model for symbolic regression via joint supervised learning](#). In *The Eleventh International Conference on Learning Representations*.

- Zhenyong Li, Gabriel Poesia, Omar Costilla-Reyes, Noah Goodman, and Armando Solar-Lezama. 2022c. [Lemma: Bootstrapping high-level mathematical reasoning with learned symbolic abstractions](#). *Preprint*, arXiv:2211.08671.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. Improved baselines with visual instruction tuning. In *NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following*.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024. [Llava-next: Improved reasoning, ocr, and world knowledge](#).
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023c. [Lost in the middle: How language models use long contexts](#). *Preprint*, arXiv:2307.03172.
- Wentao Liu, Hanglei Hu, Jie Zhou, Yuyang Ding, Junsong Li, Jiayi Zeng, Mengliang He, Qin Chen, Bo Jiang, Aimin Zhou, and Liang He. 2023d. [Mathematical language models: A survey](#). *Preprint*, arXiv:2312.07622.
- Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. 2023. [Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct](#). *Preprint*, arXiv:2308.09583.
- Meta. 2024. [Meta Llama 3](#).
- Suvir Mirchandani, Fei Xia, Pete Florence, Brian Ichter, Danny Driess, Montserrat Gonzalez Arenas, Kanishk Rao, Dorsa Sadigh, and Andy Zeng. 2023. [Large language models as general pattern machines](#). *Preprint*, arXiv:2307.04721.
- Quang Uy Nguyen, Nguyen Hoai, Michael O’Neill, Robert McKay, and Edgar Galván-López. 2011. [Semantically-based crossover in genetic programming: Application to real-valued symbolic regression](#). *Genetic Programming Theory and Evolvable Machines*, 12:91–119.
- OpenAI. 2023. [Gpt-4 technical report](#). *Preprint*, arXiv:2303.08774.
- Brenden K Petersen, Mikel Landajuela Larra, Terrell N. Mundhenk, Claudio Prata Santiago, Soo Kyung Kim, and Joanne Taery Kim. 2021. Deep symbolic regression: Recovering mathematical expressions from data via risk-seeking policy gradients. In *International Conference on Learning Representations*.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. [Hierarchical text-conditional image generation with clip latents](#). *Preprint*, arXiv:2204.06125.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. [Zero-shot text-to-image generation](#). *Preprint*, arXiv:2102.12092.
- Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2023. [Code llama: Open foundation models for code](#). *Preprint*, arXiv:2308.12950.
- Michael Schmidt and Hod Lipson. 2011. [Age-fitness pareto optimization](#). In Rick Riolo, Trent McConaghy, and Ekaterina Vladislavleva, editors, *Genetic Programming Theory and Practice VIII*, pages 129–146. Springer New York, New York, NY.
- Parshin Shojaee, Kazem Meidani, Amir Barati Farimani, and Chandan K. Reddy. 2023. Transformer-based planning for symbolic regression. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. 2022. Flava: A foundational language and vision alignment model. In *CVPR*, pages 15617–15629.
- Guido F. Smits and Mark Kotanchek. 2005. [Pareto-front exploitation in symbolic regression](#). In Una-May O’Reilly, Tina Yu, Rick Riolo, and Bill Worzel, editors, *Genetic Programming Theory and Practice II*, pages 283–299. Springer US, Boston, MA.
- Hugo Touvron, Louis Martin, Kevin Stone, et al. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *Preprint*, arXiv:2307.09288.
- Mojtaba Valipour, Bowen You, Maysum Panju, and Ali Ghodsi. 2021. [Symbolicgpt: A generative transformer model for symbolic regression](#). *Preprint*, arXiv:2106.14131.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz

- Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- M. Virgolin, T. Alderliesten, C. Witteveen, and P. A. N. Bosman. 2021. [Improving model-based genetic programming for symbolic regression of small expressions](#). *Evolutionary Computation*, 29(2):211–237.
- Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, SciPy 1.0 Contributors, et al. 2020. [SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python](#). *Nature Methods*, 17:261–272.
- Mark Willis, Hugo Hiden, P. Marenbach, Ben McKay, and Gary Montague. 1997. [Genetic programming: An introduction and survey of applications](#). pages 314 – 319.
- Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun Chen. 2023. [Large language models as optimizers](#). *Preprint*, arXiv:2309.03409.
- Jingyi Zhang, Jiaxing Huang, Sheng Jin, and Shijian Lu. 2023. [Vision-language models for vision tasks: A survey](#). *Preprint*, arXiv:2304.00685.

APPENDIX

A Prompts

The prompt used to generate the seed functions is reported in Figure 3 while the prompt used during the OPRO loop is reported in Figure 4. It’s worth noting that when using a text-only model we remove all mentions of images.

B Hyperparameters

We report the hyperparameters used when sampling from LLMs (Table 4) and for gplearn (Table 5).

Parameter	Value
temperature	1.0
top_p	0.9
top_k	60
num_beams	1
max_new_tokens	512

Table 4: Sampling parameters for the LLMs.

Parameter	Value
population_size	1000
generations	20
tournament_size	20
const_range	(-5, 5)

Table 5: Hyperparameters for gplearn.

C Benchmark functions

The list of functions and point ranges for all the benchmarks can be found in Table 6. The range for training and testing points was taken from the original source where available. Nguyen and Constant don’t include a range for the testing points, so we used the same range as the training points but with more sample points. $\mathcal{R}[\min, \max, \text{num}]$ indicates points sampled between the interval from a uniform distribution, while $\mathcal{U}[\min, \max, \text{num}]$ indicates point sampled at equal distance from each other, forming a meshgrid.

D Sample results

We present a sample of one solution for each function in the benchmarks found by our method, to qualitatively investigate the generated expressions. The true function is seen in red, while the model’s guess is seen in green. The samples can be seen in Figures 5, 6, 7, 8, 9, 10, 11 and 12.

Some of the failures of the models are apparent: in areas where there is a low density of training points the model sometimes makes guesses that ignore the overall trend, as seen more dramatically in R3 (Figures 11 and 12). However, generally the method seems to hold up well, with the large majority of the guesses being appropriate for the benchmarked function.

We also report a sample of the results in the Keijzer benchmark for the ablation study performed in Section 5.4, using only the seed functions generation step, in Figure 13. The results for the other three benchmarks are very similar to the ones reported for the full method.

I want you to act as a mathematical function generator. Given a set of points below, you are to come up with 5 potential functions that would fit the points. Don't worry too much about accuracy: your task is to generate a set of functions that are as diverse as possible, so that they can serve as starting points for further optimization. To generate the functions, you will start from a set of basic operators and expressions, and combine them into something more complex.

Your options are:

- An independent variable symbol: x .
- A coefficient symbol: c (there is no need to write a number - write this generic coefficient instead).
- Basic operators: $+$, $-$, $*$, $/$, $\sqrt{\quad}$, \exp , \log , abs
- Trigonometric expressions: \sin , \cos , \tan , \sinh , \cosh , \tanh

Make sure there are no numbers in the functions, use the coefficient token ' c ' instead. Analyze the points carefully: if there are any negative points in the input, $\sqrt{\quad}$ and \log can not be used unless the input is combined with abs .

The functions should all begin with the indicators " $f_1(x) =$ ", " $f_2(x) =$ "... Your task is to combine an arbitrary number of these basic blocks to create a complex expression. Don't be afraid to be creative and experiment! The functions should be as complex as possible, combining many different operations. Variety is key!

Points: {points}

Functions:

Figure 3: Prompt used to generate the seed functions.

I want you to act as a mathematical function generator.

You are given an image of a graph with a set of point plotted as a scatter plot in blue. The coordinates of the points are: points

Below are some previous functions and the error they make on the points above. The errors are arranged in order of their fit values, with the highest values coming first, and lower is better. These functions are also shown in the image in the form of a green line and the formula for the plotted function can be seen in the title.

Your task is to give me a list of five new potential functions that are different from all the ones reported below, and have a lower error value than all of the functions below. Only output the new functions and nothing else.

Remember that the functions you generate should always have at most {num_variables} variables {variables_list}.

The functions should have parametric form, using ' c ' in place of any constant or coefficient. The coefficients will be optimized to fit the data. Make absolutely sure that the functions you generate are completely different from the ones already given to you.

The functions should all begin with the indicators " $f_1(x) =$ ", " $f_2(x) =$ "...

Remember that you can combine the simple building blocks (operations, constants, variables) in any way you want to generate more complex functions. Don't be afraid to experiment!

{previous_trajectory}

Figure 4: Prompt used during the OPRO loop.

Experiment	Function	Train Points	Test Points
nguyen1	$x^3 + x^2 + x$	$\mathcal{R}[-1, 1, 20]$	$\mathcal{U}[-1, 1, 200]$
nguyen2	$x^4 + x^3 + x^2 + x$	$\mathcal{R}[-1, 1, 20]$	$\mathcal{U}[-1, 1, 200]$
nguyen3	$x^5 + x^4 + x^3 + x^2 + x$	$\mathcal{R}[-1, 1, 20]$	$\mathcal{U}[-1, 1, 200]$
nguyen4	$x^6 + x^5 + x^4 + x^3 + x^2 + x$	$\mathcal{R}[-1, 1, 20]$	$\mathcal{U}[-1, 1, 200]$
nguyen5	$\sin(x^2) \cdot \cos(x) - 1$	$\mathcal{R}[-1, 1, 20]$	$\mathcal{U}[-1, 1, 200]$
nguyen6	$\sin(x) + \sin(x + x^2)$	$\mathcal{R}[-1, 1, 20]$	$\mathcal{U}[-1, 1, 200]$
nguyen7	$\log(x + 1) + \log(x^2 + 1)$	$\mathcal{R}[0, 2, 20]$	$\mathcal{U}[0, 2, 200]$
nguyen8	\sqrt{x}	$\mathcal{R}[0, 4, 20]$	$\mathcal{U}[0, 4, 200]$
nguyen9	$\sin(x_1) + \sin(x_2^2)$	$\mathcal{R}[-1, -1], [1, 1], 100]$	$\mathcal{U}[-1, -1], [1, 1], 500]$
nguyen10	$2 \cdot \sin(x_1) \cdot \cos(x_2)$	$\mathcal{R}[-1, -1], [1, 1], 100]$	$\mathcal{U}[-1, -1], [1, 1], 500]$
nguyen11	$x_1^{x_2}$	$\mathcal{R}[0, 0], [1, 1], 100]$	$\mathcal{U}[0, 0], [1, 1], 500]$
nguyen12	$x_1^4 - x_1^3 + \frac{1}{2} \cdot x_2^2 - x_2$	$\mathcal{R}[-1, -1], [1, 1], 100]$	$\mathcal{U}[-1, -1], [1, 1], 500]$
constant1	$3.39x^3 + 2.12x^2 + 1.78x$	$\mathcal{R}[-1, 1, 20]$	$\mathcal{U}[-1, 1, 200]$
constant2	$\sin(x^2) \cdot \cos(x) - 0.75$	$\mathcal{R}[-1, 1, 20]$	$\mathcal{U}[-1, 1, 200]$
constant3	$\sin(1.5x_1) \cdot \cos(0.5x_2)$	$\mathcal{R}[-1, -1], [1, 1], 100]$	$\mathcal{U}[-1, -1], [1, 1], 500]$
constant4	$2.7x_1^{x_2}$	$\mathcal{R}[0, 0], [1, 1], 100]$	$\mathcal{U}[0, 0], [1, 1], 500]$
constant5	$\sqrt{1.23x}$	$\mathcal{R}[0, 4, 20]$	$\mathcal{U}[0, 4, 200]$
constant6	$x^{0.426}$	$\mathcal{R}[0, 4, 20]$	$\mathcal{U}[0, 4, 200]$
constant7	$2 \sin(1.3x_1) + \cos(x_2)$	$\mathcal{R}[-1, -1], [1, 1], 100]$	$\mathcal{U}[-1, -1], [1, 1], 500]$
constant8	$\ln(x + 1.4) + \ln(x^2 + 1.3)$	$\mathcal{R}[0, 2, 20]$	$\mathcal{U}[0, 2, 200]$
keijzer3	$0.3x \cdot \sin(2\pi x)$	$\mathcal{U}[-1, 1, 100]$	$\mathcal{U}[-1, 1, 10000]$
keijzer4	$x^3 \cdot \exp(-x) \cdot \cos(x) \sin(x) \cdot (\sin(x)^2 \cdot \cos(x) - 1)$	$\mathcal{U}[0, 10, 200]$	$\mathcal{U}[0.05, 10.05, 200]$
keijzer6	$(x \cdot (x + 1))/2$	$\mathcal{U}[-1, 1, 50]$	$\mathcal{U}[-1, 1, 100]$
keijzer7	$\ln(x)$	$\mathcal{U}[1, 100, 100]$	$\mathcal{U}[1, 100, 1000]$
keijzer8	\sqrt{x}	$\mathcal{U}[0, 100, 100]$	$\mathcal{U}[0, 100, 1000]$
keijzer9	$\ln(x + \sqrt{x^2 + 1})$	$\mathcal{U}[0, 100, 100]$	$\mathcal{U}[0, 100, 1000]$
keijzer10	$x_1^{x_2}$	$\mathcal{R}[0, 1, 100]$	$\mathcal{U}[0, 1, 1000]$
keijzer11	$x_1 \cdot x_2 + \sin((x_1 - 1) \cdot (x_2 - 1))$	$\mathcal{R}[-3, 3, 20]$	$\mathcal{U}[-3, 3, 1000]$
keijzer12	$x_1^4 - x_1^3 + \frac{(x_2^2)}{2} - x_2$	$\mathcal{R}[-3, 3, 20]$	$\mathcal{U}[-3, 3, 1000]$
keijzer13	$6 \cdot \sin(x_1) \cdot \cos(x_2)$	$\mathcal{R}[-3, 3, 20]$	$\mathcal{U}[-3, 3, 1000]$
keijzer14	$8/(2 + x_1^2 + x_2^2)$	$\mathcal{R}[-3, 3, 20]$	$\mathcal{U}[-3, 3, 1000]$
keijzer15	$\frac{x_1^3}{5} + \frac{x_2^3}{2} - x_2 - x_1$	$\mathcal{R}[-3, 3, 20]$	$\mathcal{U}[-3, 3, 1000]$
R1	$(x + 1)^3/(x^2 - x + 1)$	$\mathcal{R}[-1, 1, 20]$	$\mathcal{U}[-1, 1, 20]$
R2	$(x^5 - 3 \cdot x^3 + 1)/(x^2 + 1)$	$\mathcal{R}[-1, 1, 20]$	$\mathcal{U}[-1, 1, 20]$
R3	$(x^6 + x^5)/(x^4 + x^3 + x^2 + x + 1)$	$\mathcal{R}[-1, 1, 20]$	$\mathcal{U}[-1, 1, 20]$

Table 6: Functions and point ranges for all benchmarks.

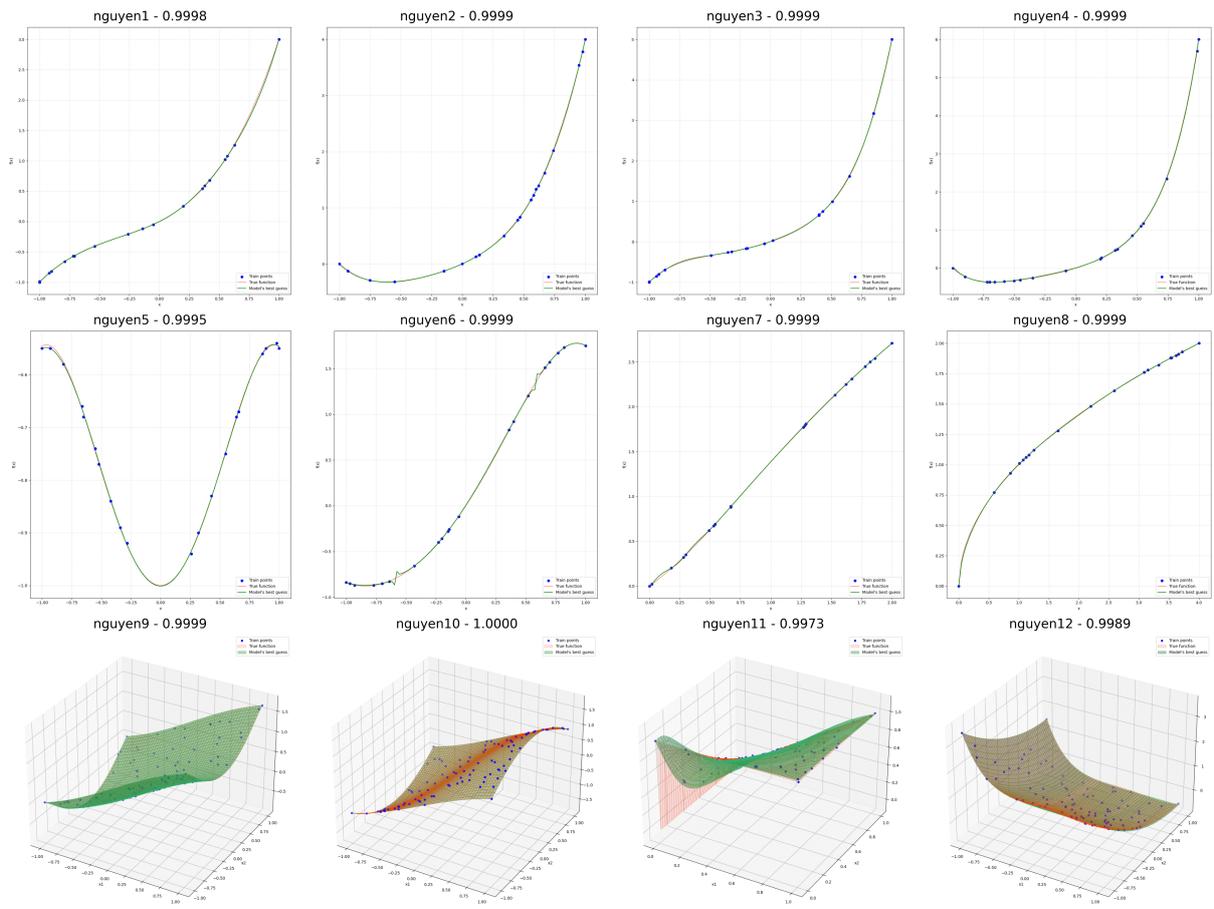


Figure 5: Results for the Nguyen benchmark with Llama3.

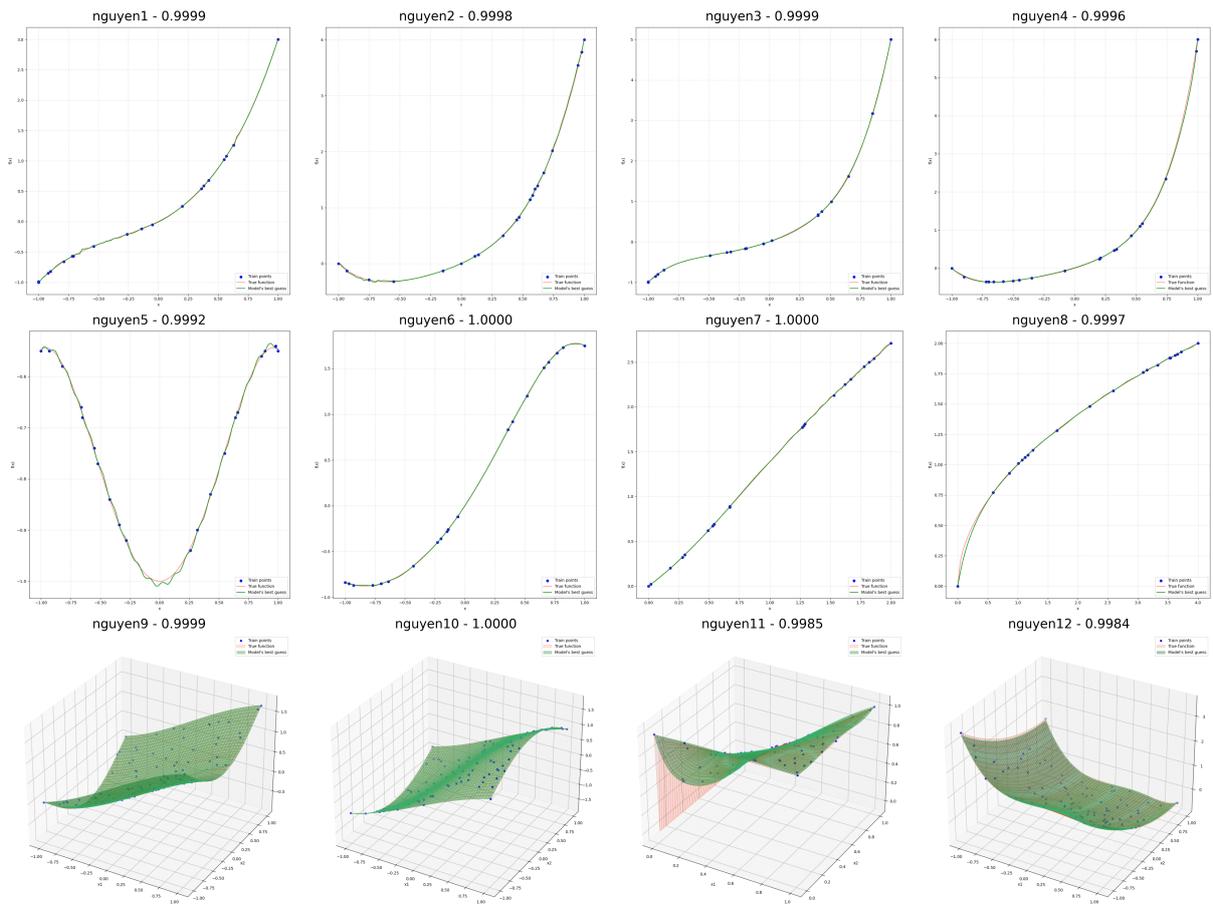


Figure 6: Results for the Nguyen benchmark with LLaVa-NeXT.

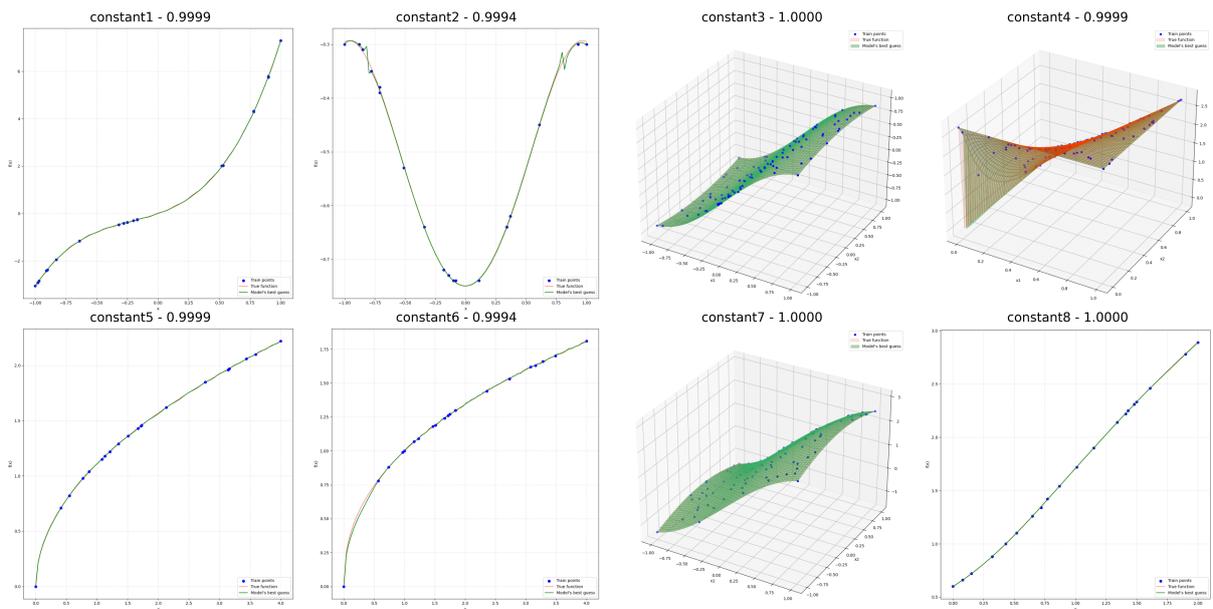


Figure 7: Results for the Constant benchmark with Llama3.

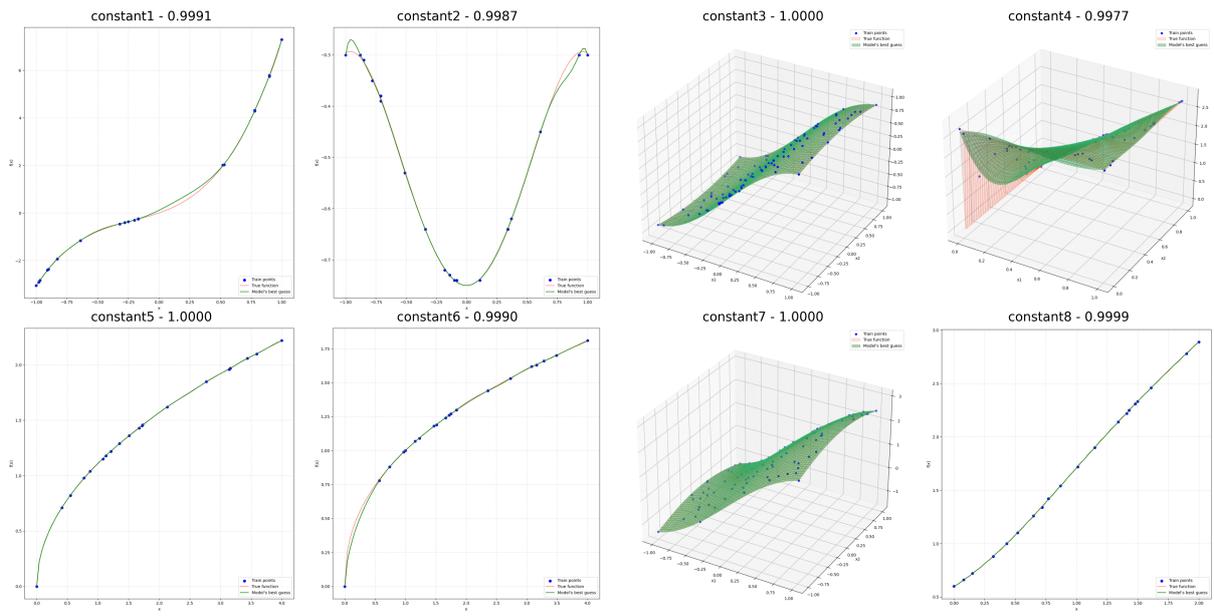


Figure 8: Results for the Constant benchmark with LLaVa-NeXT.

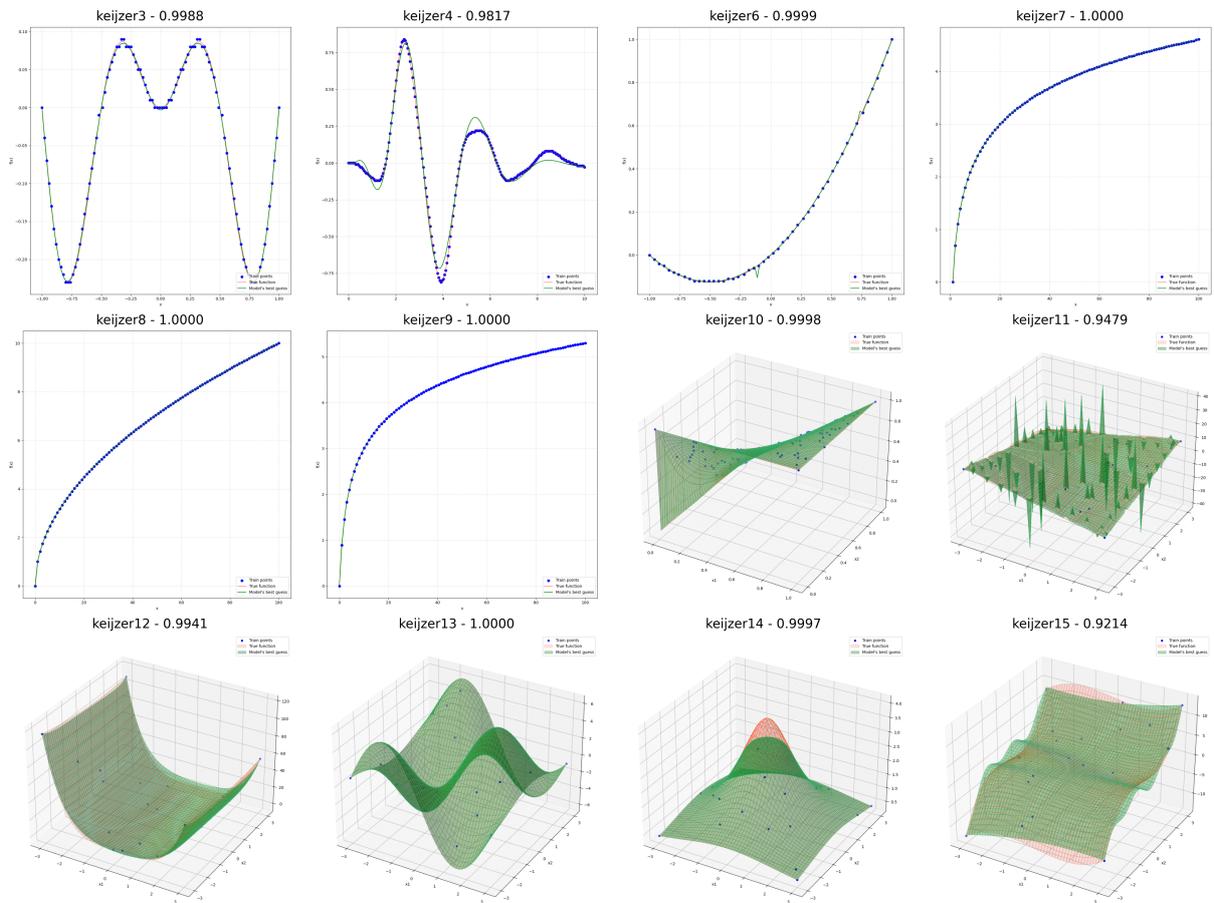


Figure 9: Results for the Keijzer benchmark with Llama3.

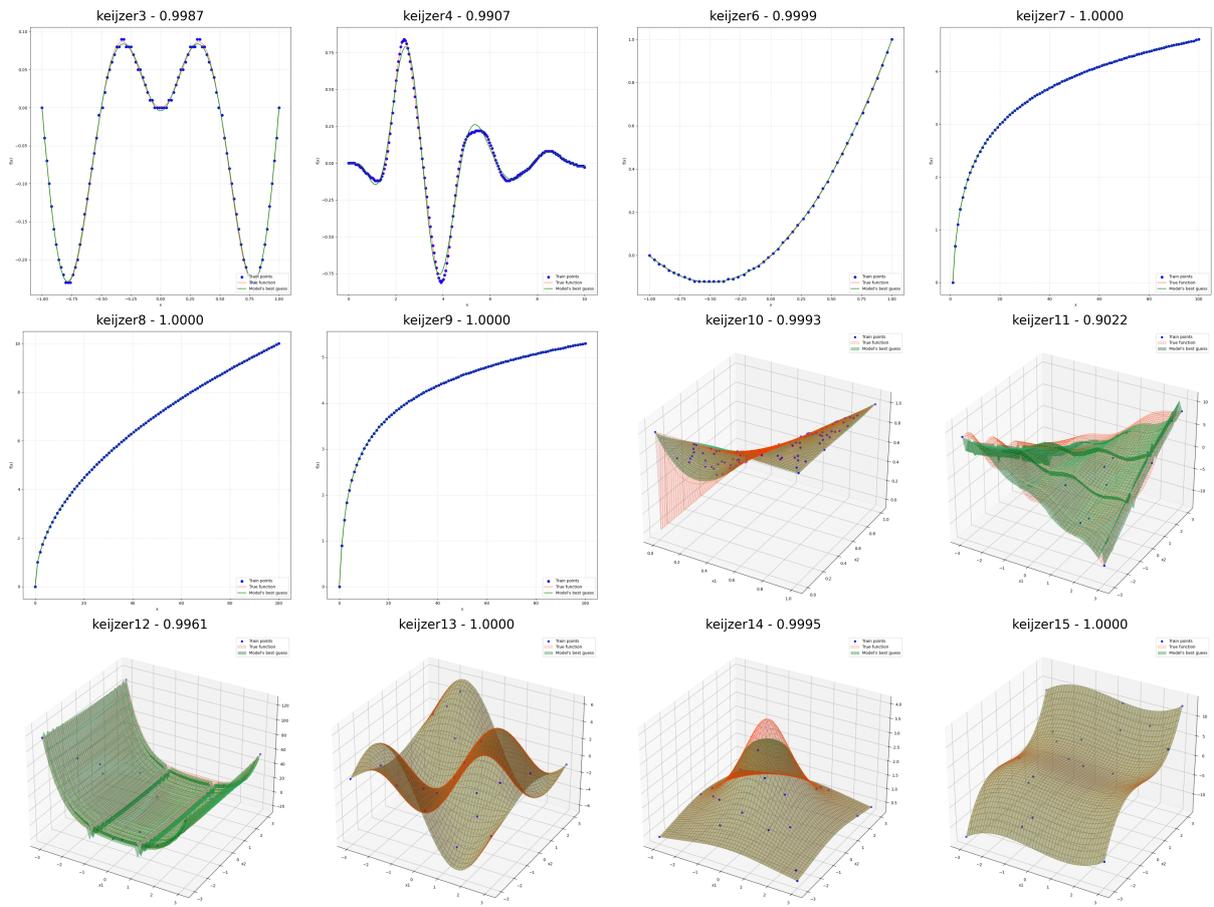


Figure 10: Results for the Keijzer benchmark with LLaVa-NeXT.

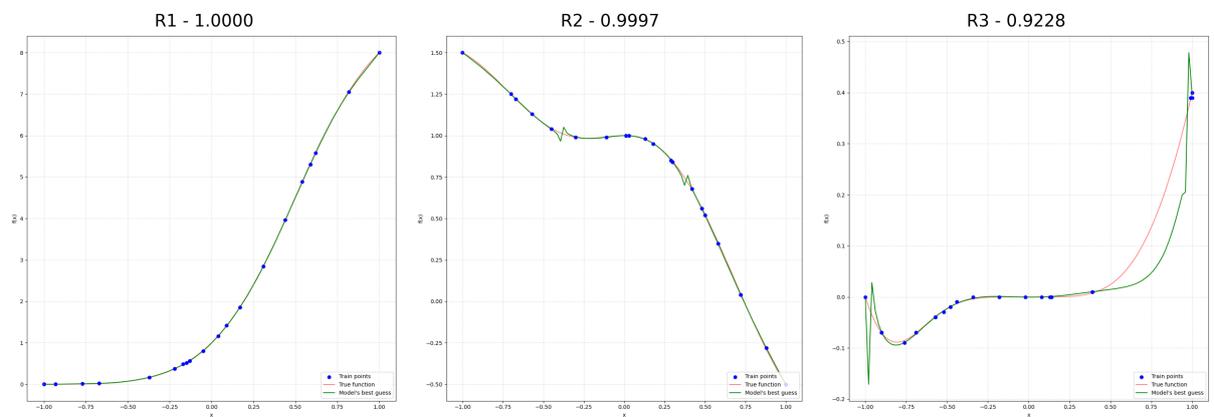


Figure 11: Results for the R benchmark with Llama3.

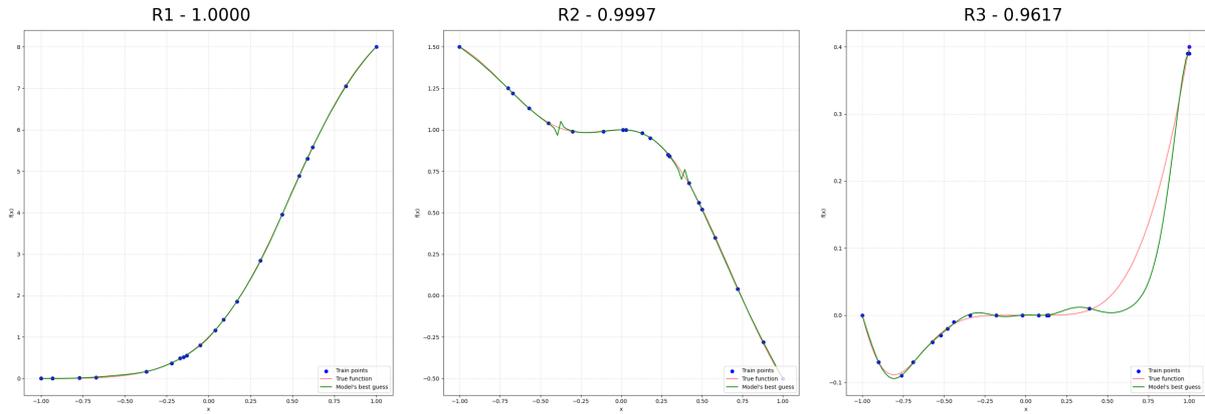


Figure 12: Results for the R benchmark with LLaVa-NeXT.

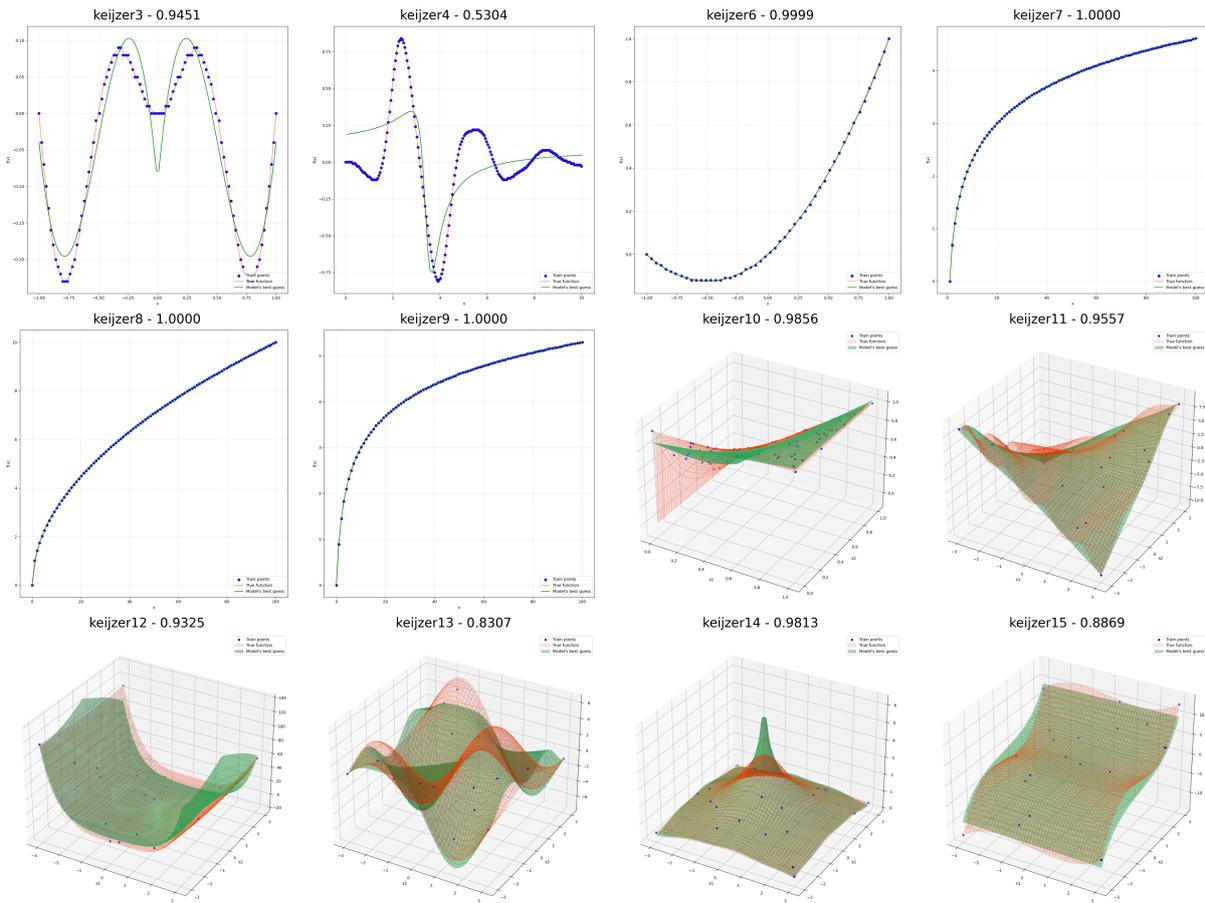


Figure 13: Results for the Keijzer benchmark using only the seed functions generation step.