

Improving Genetic Programming for Symbolic Regression with Equality Graphs

Fabrício Olivetti de França
Federal University of ABC
Santo Andre, São Paulo, Brazil
folivetti@ufabc.edu.br

Gabriel Kronberger
University of Applied Sciences Upper Austria
Hagenberg, Upper Austria, Austria
gabriel.kronberger@fh-hagenberg.at

Abstract

The search for symbolic regression models with genetic programming (GP) has a tendency of revisiting expressions in their original or equivalent forms. Repeatedly evaluating equivalent expressions is inefficient, as it does not immediately lead to better solutions. However, evolutionary algorithms require diversity and should allow the accumulation of inactive building blocks that can play an important role at a later point. The equality graph is a data structure capable of compactly storing expressions and their equivalent forms allowing an efficient verification of whether an expression has been visited in any of their stored equivalent forms. We exploit the e-graph to adapt the subtree operators to reduce the chances of revisiting expressions. Our adaptation, called *egg*, stores every visited expression in the e-graph, allowing us to filter out from the available selection of subtrees all the combinations that would create already visited expressions. Results show that, for small expressions, this approach improves the performance of a simple GP algorithm to compete with PySR and Operon without increasing computational cost. As a highlight, *egg* was capable of reliably delivering short and at the same time accurate models for a selected set of benchmarks from SRBench and a set of real-world datasets.

CCS Concepts

• **Computing methodologies** → **Symbolic and algebraic algorithms**; • **Mathematics of computing** → *Genetic programming*.

Keywords

Symbolic regression, Genetic programming, Equality saturation, Equality graphs

ACM Reference Format:

Fabrício Olivetti de França and Gabriel Kronberger. 2024. Improving Genetic Programming for Symbolic Regression with Equality Graphs. In . ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 Introduction

Symbolic regression (SR) searches for a mathematical function that approximates a set of data points. It can be used to find nonlinear

regression models [17, 18] or for scientific discovery [6, 8, 18, 23, 31, 35].

The current state-of-the-art SR algorithms [10, 22] use genetic programming (GP) as the main search engine. One of such improvements, used by several state-of-the-art implementations, is the search for a parametric function that replaces the constants with adjustable parameters that are fitted using nonlinear optimization techniques.

Not only is the search for a symbolic model NP-hard [36] but when searching for a parametric model, it usually requires the solution to a multimodal optimization problem, which by itself is NP-hard [27] and can hinder the search for the optimal solution.

To make matters worse, the usual way of encoding mathematical expressions as symbolic expression trees, allows GP to visit different but equivalent expressions [21] that evaluate to the same values. These equivalent expressions may be unnecessarily large and contain redundant parameters, which reduces the probability of finding their optimal values [9, 19]. Even for the simple expression p_1x_1 we can produce an infinite number of equivalent expressions considering that p are fitting parameters, for example $((p_1x_1) + (p_2x_1))$, x_1/p_1 , $x_1^2/(p_1x_1)$, etc. are all different parameterizations of the same expression.

GP cannot differentiate between equivalent forms of a given expression, and simplification heuristics are often insufficient, as seen in [9]. Some authors [16, 26] argue that redundancy is necessary to allow the algorithm to navigate through the search space, as these equivalent expressions are guaranteed to have the same accuracy, allowing the search to keep multiple genetically different variations of solution candidates in the hopes of finding a better solution. However, we do not know what would happen to GP search dynamics if we try to prevent keeping semantic duplicates in the GP population.

Equality saturation [37] uses a set of equivalence rules and produces all equivalent forms of a given expression. It alleviates the phase ordering problem in the optimization phase when compiling computer programs. Given a program represented as a symbolic expression tree, and a set of equivalence rules, it iteratively applies the rules and stores all equivalent programs in a compact data structure called *equality graph* (*e-graph*). The main idea is that upon saturation, the graph contains all equivalent forms of the original program and the optimal form can be extracted from the e-graph using a heuristic. This technique was previously used in the context of SR in [9, 19] to investigate the problem of overparameterization that can negatively affect the fitting of numerical parameters. The e-graph has another interesting feature that can be exploited by SR algorithms: it contains a database of all visited patterns and their

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions to permissions@acm.org.
Conference'17, July 2017, Washington, DC, USA

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

equivalent forms that can be easily matched against new candidate expressions.

In this work, we investigate the benefits of exploiting the e-graph to increase the chances of generating a novel expression when applying crossover and mutation. In short, the proposed crossover operator will limit the selection of the second parent to a subset of subtrees that can only generate an unvisited expression. For the subtree mutation, it will replace a random subtree from the original tree t with a random subtree e generating the new tree t' . If t' has already been visited, the root of e will be changed to a random choice of nonterminals that generates an unvisited expression.

These operators are tested inside a simple GP implementation based on tinyGP [33] and compared against tinyGP with parameter fitting and two state-of-the-art algorithms: Operon [5] and PySR [8]. The results show that this *simple*¹ modification can improve the performance of tinyGP to an extent that it becomes competitive (and in some aspects better) than the state-of-the-art. The use of an e-graph as a support structure for GP brings new light to symbolic regression and GP with the possibility of exploring the accumulated history of the search process and even combining the history of multiple searches.

This paper is organized such that in Section 2 we will summarize the related works in symbolic regression. Section 3 will explain the basic concepts of equality saturation and the e-graph data structure. In Section 4 we will detail the proposed modification to simple GP. Section 5 will show the experiment methods used in this paper followed by Section 6 where we report and discuss the results. Finally, Section 7 gives some final remarks and expectations for the future.

2 Related work

The redundancy of GP search space has been investigated by many authors with conflicting conclusions to whether this is beneficial or not for the search. For example, Ebner [12] argued that this redundancy enables the search to reach the optima through different trajectories, increasing the chances of achieving one of the equivalent expressions. On the other hand, Gustafson et al. [13] observed that when the recombination between two similar solutions was forbidden, there was an increase in offsprings that changed the original behavior of their parents, leading to increased performance.

Several works made a detailed study about the redundancy and neutrality in GP (i.e., when a change in the solution has no effect on its outcome). For example, Hu, Banzhaf, and Ochoa [1, 14, 15] investigated linear GP for Boolean SR problems with the help of search trajectory networks showing that some phenotypes are overrepresented in the search space. Regarding subtree crossover, McPhee et al. [25] showed that over 75% of crossovers produced no immediately useful semantic changes.

Kronberger et al. [21] studied the inefficiency of a simple GP comparing with the enumerated search space [2] and using equality saturation to count the percentage of unique expressions generated during the GP search. They found that from the total of visited expressions during the search, only around 40% were unique. This

not only wastes computational resource but it also shows that, at some point, GP fails to explore different regions of the search space.

Many authors observed improvements in the obtained solutions when applying any form of simplification during the search [7, 28, 30, 32] and, as a side effect, it stimulates the diversity of the population [3, 4].

Equality saturation has been used in the context of symbolic regression as a support tool to study the behavior of the search. Many state-of-the-art SR algorithms have a bias towards creating expressions with redundant numerical parameters [9, 19]. This redundancy can increase the chance of failing to correctly optimize such parameters, leading to sub-optimal solutions. In [21] this technique was used to detect the equivalent expressions visited during the GP search. So far, the equality saturation technique was not used during the GP search to improve the quality of the solutions.

3 Equality saturation and e-graphs

Equality saturation [34] was proposed as a solution to the phase ordering problem in compiler optimization. This problem occurs when optimizing a program by applying a set of rewrite rules sequentially while dropping the information about the previous versions of the program. If the optimization follows a non-optimal sequence, it will lead to a sub-optimal program.

Equality saturation solves this issue by applying all of the optimization rules in parallel while keeping the intermediate transformations in a compact form using the data structure called *e-graph*.

Fig. 1a illustrates an example of an e-graph. Each solid box represents an e-node that contains a symbol of the expression. The dashed boxes group a set of e-nodes together and it is called an e-class, each e-class is assigned an e-class id (number in the bottom right of an e-class box). The main property of an e-class is that, no matter which e-node is chosen during the traversal, it will lead to an equivalent expression with all other e-nodes.

We can see that in the middle box (e-class id 4), if we follow through \times it will generate the expression $2x$ and if we follow through $+$ it will generate $x + x$. The abstract description of the algorithm is very simple, though a concrete and optimal implementation requires the use of advanced techniques and data structures. The main idea is: i) match all the equivalence rules in the current state of the e-graph, ii) apply the rules creating new e-classes, iii) merge the equivalent e-classes, iv) repeat until saturation (i.e., no changes occur).

The current state-of-the-art implementation [37]², called *egg*³, encodes the e-graph as a map of an e-node to an e-class, a map of the e-class ids to their structure, and a set of e-classes that still needs analysis to maintain consistency. The structure of an e-class contains the information about the e-nodes it contains, a list of the parent e-nodes, and additional information also referred to as *semantic analysis*. This implementation also maintains a database of patterns that allows the algorithm to efficiently match patterns inside the e-graph structure.

Usually, the e-graph is used to represent a single program or expression and their equivalent forms. But, the structure can keep any number of expressions as long as we keep a list of the e-classes

¹We acknowledge that the modification is simple provided we have a working implementation of e-graph.

²<https://docs.rs/egg/latest/egg/>

³egg stands for e-graphs good

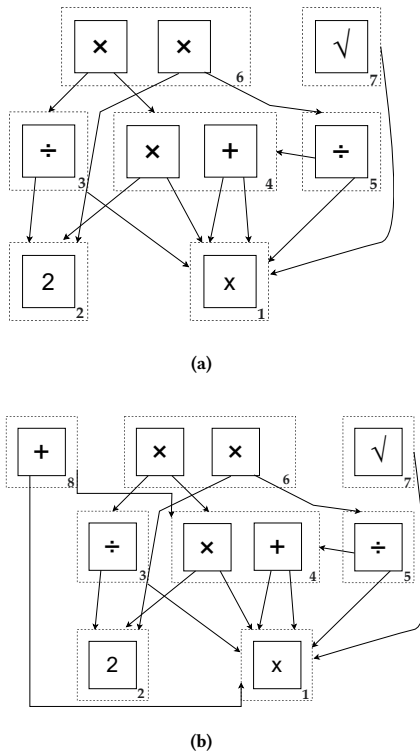


Figure 1: (a) Illustrative example of an e-graph and (b) the same e-graph after inserting the expression $x + 2x$.

ids representing the root of each expression. For example, if we insert the expressions $(2/x)(x + x)$, $2x$, \sqrt{x} into the e-graph, we would end up with Fig. 1a, minus the equivalent relations. As a result, we would keep the list $[6, 4, 7]$ representing the ids of the expressions we inserted. New expressions are added bottom-up. Starting from a terminal, the algorithm checks whether it already exists in the e-graph and returns its e-class id, if it does not exist, it creates a new id. When adding an internal node, the algorithm first converts it to an e-node by replacing its children by their e-class ids and then it checks whether it already exists in the e-graph, returning the corresponding e-class id or a new one. In our example from Fig. 1, if we try to add the expression $x + 2x$, it would first retrieve the e-class ids 1, 2 for the terminals x and 2, then it would return the e-class id 4 corresponding to the e-node 2×1 , where underlined numbers correspond to e-class ids. Finally, it would create a new e-class with id 8 and the e-node $\underline{1} + \underline{4}$. This mechanism allows us to compactly store a set of expressions and readily assert whether an expression already exists in the structure.

In our implementation, we made a series of improvements to this data structure as we will describe in the next section.

4 eggp: e-graph GP

The proposed algorithm, *eggp* (e-graph genetic programming), follows the same structure as the traditional GP. The initial population is created using ramped half-and-half respecting a maximum

size and maximum depth parameter [17] and, for a number of generations, it will choose two parents using tournament selection, apply the subtree crossover with probability pc followed by the subtree mutation with probability pm , when the offsprings replace the current population following a dominance criteria.

The key differences of eggp are:

- (1) new solutions are inserted into the e-graph followed by one step of equality saturation to find and store some of the equivalent expressions of the new offspring.
- (2) the current population is replaced by the set of individuals formed by: the Pareto front, the next front after excluding the first Pareto-front, and a selection of the last offspring at random until it reaches the desired population size.
- (3) the subtree crossover and mutation are modified to try to generate an unvisited expression.

Apart from that, we do not make use of other mutation operators commonly used in the literature [5, 8] or advanced techniques to stimulate diversity such as island models. Notice that a single step of equality saturation will not guarantee the insertion of all equivalent expressions in the e-graph but, as we apply more iterations, the e-graph can grow exponentially large. This issue is amplified by the fact that we are storing multiple expressions. As we will see in Sec. 6, the single step seems to be sufficient to improve the results and the benefits of increasing the number of steps is a subject for a future research. Regarding the Pareto-front extracted from the entire history, this is equivalent to the traditional NSGA-II algorithm [11] as the dominance relation is transitive. Keeping two *ranks* of dominance and filling up the remainder of the population with new expressions is meant to stimulate the combination of new expressions (exploration) while keeping the best fronts (exploitation). As a side-note, in this implementation we keep a database of generated expressions sorted by fitness and size (both objectives used in this work), so the Pareto-front can be retrieved in $O(n)$ where n is the number of extracted individuals. A new expression can be inserted into this structure in $O(\log(m))$, where m is the number of evaluated expressions so far. Finally, if we add the expression $x + x + x$, equality saturation will generate the equivalent form θx (constants are replaced by parameters), this expression is stored in the database of expressions with size 3.

4.1 Crossover and mutation

Since the algorithm keeps the whole history of visited expression and its equivalent forms, we can exploit this information to increase the probability of generating unvisited expressions in the perturbation operators. For the crossover (Fig. 2a), the main idea is that it first chooses a crossover point from the first parent and then it randomly chooses a subtree of the second parent which generates an unvisited expression upon combination.

For the mutation operator (Fig. 2b), it chooses a point of the tree to generate a new subtree and, after generating the new subtree, it checks whether it generates an unvisited expression. With a negative answer, it randomly changes the root node of this subtree from a choice of all the possible non-terminals that can ensure novelty. In the event of impossibility of generating a new expression, a random recombination is returned.

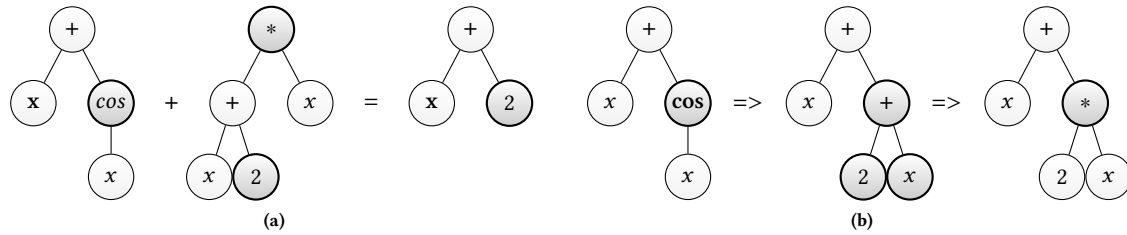


Figure 2: Examples using the e-graph in Fig. 1b of (a) recombination between two expressions: after choosing the recombination point marked in bold in the first tree, the second tree has only two points which will generate new expressions (marked in bold in the second expression), after picking one of these points, we generate the new solution illustrated in the tree to the right; (b) mutation: after choosing the mutation point, a new subtree is generated. If the new expression is already contained in the e-graph, the root of the subtree is changed by a random non-terminal that creates an unvisited expression.

Algorithm 1 `traverse-tree`, $\lambda e.f$ represents a λ -function that receives the argument e and evaluates the expression f .

Require: current node n , node index i , list of λ -parents p , subtree to be inserted s , e-graph g

```

1: if  $i = 0$  then
2:   return not-in-e-graph( $n, p, 2$ )
3: else
4:   if should go left or  $\text{arity}(n) = 1$  then
5:     traverse-tree(left( $n$ ),  $i - 1$ , ( $\lambda e.n(e, \text{right}(n)) : p, s$ ))
6:   else
7:      $nl \leftarrow \text{size}(\text{left}(n))$ 
8:     traverse-tree(right( $n$ ),  $i - nl - 1$ , ( $\lambda e.n(\text{left}(n), e) : p, s$ ))

```

To verify whether the combination of the original tree n with the subtree s replacing the i -th node will generate an unvisited expression, we traverse the tree in pre-order traversal until we reach the node we want to replace (i -th node). While traversing the tree, we build a list of λ -functions equivalent to a function $f(e, n)$ that, given a tree e , it creates a new tree with the operator of n as the root, replacing one child with e and keeping the other child of n (if it is a binary operator). So, if the next step in the traversal is to the left, this function will build a tree rooted at n with the left child being replaced by e and the right child as $\text{right}(n)$, this is illustrated in Alg. 1. Upon reaching the desired node, we traverse the tree upwards by checking whether the current node n exists in the e-graph and, if the assertion is false, it will apply the first function of the list to n , generating the parent e-node, and calling this same function recursively (Alg. 2). Given that none of the ancestors exist in the e-graph, the algorithm confirms that this is an unvisited expression.

The full implementation of `egg` has nine hyperparameters: number of generations, population size, maximum expression size, likelihood function (MSE, Gaussian, Poisson, Bernoulli, ROXY [24]), number of iterations and retries for the parameter optimization, probabilities of crossover and mutation, and the list of non-terminals.

The algorithm is implemented in Haskell using the `srtree`⁴ library for symbolic regression and with an adapted implementation of

⁴<https://github.com/folivetti/srtree>

Algorithm 2 `not-in-e-graph`

Require: current e-node n , list of λ -parents p , e-graph g .

```

1:  $n\text{-is-new} \leftarrow \text{notIn}(n, g)$ 
2: if empty( $p$ ) or  $n\text{-is-new}$  then
3:   return  $n\text{-is-new}$ 
4: else
5:   ( $\text{parent}, p'$ )  $\leftarrow \text{uncons}(p)$ 
6:   not-in-e-graph( $\text{parent}(n), p', g$ )

```

equality saturation based on `hegg`⁵. The binaries are available at <https://github.com/folivetti/srtree/releases/tag/v2.0.1.0>.

5 Experiments

To measure the benefit of stimulating novelty using the history of visited expressions and their equivalent form stored in the e-graph, we chose three baseline algorithms: a version of `tinyGP` [33] implemented using the same backend library, `Operon` [5] and `PySR` [8].

`Operon` is a carefully crafted implementation of GP for symbolic regression with runtime performance in mind and a good set of default hyperparameters. It incorporates multiple mutation operators that allow a finer perturbation of a solution. Besides, it envelops every variable node with a scaling parameter adjusted using nonlinear optimization. `PySR` also supports the same mutation operators and nonlinear optimization of the parameters, it stands out with the use of an island model capable of keeping the diversity of the population to stimulate the exploration of the search space. It also applies a simplification heuristic on a selection of the expressions. Both `PySR` and `Operon` uses multi-objective optimization optimizing accuracy and size as default. The research questions we want to address with the experiments are:

- (1) What is the impact of increasing the probability of generating novelty, when compared to a minimalistic implementation such as `tinyGP`?
- (2) How close does `egg` get to the state-of-the-art without resorting to more advanced concepts such as specialized mutation operators, enforcing the placement of numerical parameters, and island model?

⁵<https://github.com/alt-romes/hegg>

Table 1: Symbolic regression algorithms hyperparameters. Operators enveloped with $|\cdot|$ apply the absolute value to the first argument. The population size for PySR is $1/10$ of the reported values in this table to allow the use of ten islands.

Parameter	Value
number of evaluations	100 000
prob. mutation	0.3
prob. crossover	0.9
non-terminal set	$+, -, *, \div, \log(\cdot), \exp, \sqrt{ \cdot }, x ^y$
pop. size / gens. / tourn. size	{500/200/5, 200/500/3}
max depth	10
objectives	[MSE, size]
optimization steps	2×50 (100 for Operon)

We should stress that we have kept eggp with only the subtree crossover and mutation to be directly comparable with tinyGP, thus measuring the benefits of this modification.

We have fixed all the common hyperparameters to a standard value and performed a hyperparameter tuning with two different settings for each algorithm varying only whether to have a larger population with fewer iterations or a smaller population with more iterations. The only differences in settings are: for PySR we are using 10 populations in its island model, so the size of each population is $1/10$ of the population size for the other algorithms, earlier experiments with the first set of benchmarks revealed that PySR performs significantly worse when using a single population; for Operon, we perform a maximum of 100 optimization iterations instead of 50 iterations with 2 different starting points since it does not support multiple restarts for the parameter optimization; for eggp, $1/3$ of the training data is separated and used as a validation set to calculate the fitness, while the parameters are fitted using the remaining $2/3$. Notice that tinyGP is the only single-objective approach and will be left out of the hypervolume comparison.

For the first set of experiments, we picked four datasets from the original SRBench [10] and split each one into training and test sets with a ratio of 0.7/0.3. We executed each configuration 30 times and stored the test errors of each run. Once we have the optimal hyperparameters, we evaluate these algorithms in the benchmarks of the current version of SRBench⁶ for the reduced SRBench track. This set is supposed to be representative as it contains datasets with different characteristics. For this experiment, we applied a 3-fold cross-validation repeating the experiment 10 times, creating a total of 30 runs. Finally, we picked some real-world datasets from the literature corresponding to data from different fields, these data already have pre-determined training and test sets, as such we will run 30 repetitions of each experiment. For every experiment, we store the final Pareto front and will report the performance plot, the AUC and average rank among all datasets and statistical test of the ranks, the average and standard deviation of the running time, the hypervolume of the Pareto front. All algorithms were restricted to a single core to ensure equal conditions.

Table 2: Datasets, number of points and variables, and corresponding max. size. Every training set of the SRBench group was capped at 1 000 data points chosen at random. For 192_vineyard we ensured that the rows with $x_0 = x_1 = 0$ were contained in the training set to avoid misbehaving models.

Name	Points	Features	max. size
hyperparameter tuning			
523_analcatdata_neavote	100	2	50(50)
527_analcatdata_election2000	67	14	50(50)
635_fri_c0_250_10	250	10	50(50)
1029_LEV	1 000	4	50(50)
SRBench			
192_vineyard	52	2	50(33)
210_cloud	108	5	50(33)
522_pm10	500	7	50(33)
557_analcatdata_apnea1	475	3	50(33)
579_fri_c0_250_5	250	5	50(33)
606_fri_c2_1000_10	1 000	10	50(33)
650_fri_c0_500_50	500	50	50(33)
678_visualizing_environmental	111	3	50(33)
1028_SWD	1 000	10	50(33)
1089_USCrime	47	13	50(33)
1193_BNG_lowbwt	31 104	9	50(33)
1199_BNG_echoMonths	17 496	9	50(33)
Real world			
Chemical_1_tower	4 999	25	30(20)
Chemical_2_competition	1 066	57	30(20)
Friction_stat_one-hot	2 016	16	30(20)
Friction_dyn_one-hot	2 016	17	30(20)
Flow_stress_phip0.1	7 800	2	20(13)
Nasa_battery_1_10min	636	6	20(13)
Nasa_battery_2_20min	1 638	5	20(13)
Nikuradse_1	362	2	20(13)
Nikuradse_2	362	1	20(13)

Table 2 shows the datasets of each set of experiment with their corresponding number of data points, features and the chosen maximum size parameter.⁷

6 Results and Discussion

6.1 Hyperparameter tuning

Applying the Wilcoxon rank test to the obtained results using $\alpha = 0.05$, we observed that only tinyGP expressed a significant difference in their results (p -value 0.002) favoring 200 generations and population size 500. There was no significant difference observed in any algorithm, thus, we chose the configuration which found the best result for most datasets. With this criteria, PySR and Operon was set to 200 generations and eggp with 500 generations.

⁶<https://github.com/cavalab/srbench/discussions/174#discussioncomment-10285133>

⁷The maximum size for Operon is set to 0.67 of the maximum size because, internally, Operon will not count the scale coefficients of a terminal towards the model size. This factor enforces Operon to search on a similar search space as the other algorithms.

6.2 New SRBench datasets

In Fig. 3 we can see the performance plots for each SRBench dataset of the best solution according to the highest R^2 on the training set. The x-axis of these plots represents the R^2 measured on the test set, and the y-axis shows the percentage of runs that the algorithm found an R^2 equal or larger than x . The ideal algorithm would cover the whole area from (0, 0) to (1, 1). As we can see from these plots, eggp always covers an area similar or better than the competing algorithms. Considering the area under the curve values (AUC), eggp covered the most area in 7 out of the 12 datasets, Operon in 5, PySR in 4 and tinyGP in 0. It is also evident from these plots how eggp displays a significant improvement over tinyGP.

In Table 3 we can see the ranks when considering the median R^2 of the test set and the AUC. We can see that, using this criteria, eggp is consistently ranked second, while Operon and PySR take turns in the first and third place. The statistical test reveals that we can reject the null hypothesis when comparing to PySR with the alternative of having greater median rank. On the other hand, for the AUC values, we can see that eggp is greater than the other algorithms on average while rejecting the null hypotheses for each comparison. Unlike the median of the R^2 , the AUC is the average R^2 weighted by the probability of obtaining that value or greater, acting as a reliability measure of obtaining that value or greater.

6.3 Real-world datasets

In Fig. 4 we observe a similar behavior with eggp covering an area close to or better than the best competing algorithm. Considering AUC, it found 4 out of 9 best results, while PySR only 2, Operon 6 and tinyGP 1. Again, considering the ranks on the median of the R^2 (Table 4), it consistently achieved second place, but in this set, Operon was ranked first for most of the datasets. We should notice, though, that eggp results were always close to the best competing algorithm, while even Operon misbehaved in two datasets (flow and niku-2). When observing the statistical test results, we can conclude that there are no significant differences between the eggp and Operon, but the hypotheses of eggp being equivalent to PySR and tinyGP can be rejected.

6.4 Hypervolume

Table 5 shows the average and standard deviation of the hypervolume of the final Pareto fronts. For SRBench, eggp covered a similar area to the other algorithms. PySR returned the largest hypervolume in most of the SRBench datasets. For the real-world datasets, eggp covered the most hypervolume overall.

6.5 Computational time

Regarding the runtime, since Operon is the fastest symbolic regression implementation, as noted in [5], we calculated the ratio between the average runtime of each algorithm to Operon. Fig. 5 shows the relative runtime per dataset. In this plot we can see that eggp and tinyGP were both between 5 to 10 times slower than Operon. PySR varied from 5 to 30 times the runtime of Operon depending on the dataset. The higher ratios were observed on high-dimensional or larger datasets. We should stress that all algorithms were constrained to run with a single thread, thus both Operon and PySR runtime could be smaller when exploiting multi-threading.

Table 3: Ranks of the median (1st block) and AUC (2nd block) of the test set R^2 for the SRBench. The p -values were calculated with a Wilcoxon signed-rank test using as alternative hypotheses ($\alpha = 0.05$) being greater ($>$) than eggp.

dataset	eggp	Operon	PySR	tinyGP
192_vineyard	1	4	3	2
210_cloud	1	4	2	3
522_pm10	3	4	1	2
557_analcatdata_apnea1	2	4	3	1
579_fri_c0_250_5	2	1	3	4
606_fri_c2_1000_10	1	2	3	4
650_fri_c0_500_50	2	1	3	4
678_visualizing_environmental	1	4	2	3
1028_SWD	4	1	2	3
1089_USCrime	1	3	4	2
1193_BNG_lowbwt	2	1	3	4
1199_BNG_echoMonths	4	3	1	2
mean	2.00	2.67	2.50	2.83
p -value $>$		0.08	0.02	0.06
192_vineyard	0.24	0.08	0.16	0.17
210_cloud	0.71	0.34	0.71	0.58
522_pm10	0.16	0.14	0.19	0.18
557_analcatdata_apnea1	0.75	0.40	0.58	0.68
579_fri_c0_250_5	0.91	0.95	0.83	0.79
606_fri_c2_1000_10	0.96	0.96	0.82	0.70
650_fri_c0_500_50	0.87	0.92	0.83	0.40
678_visualizing_environmental	0.28	0.08	0.27	0.20
1028_SWD	0.38	0.39	0.39	0.36
1089_USCrime	0.68	0.59	0.59	0.59
1193_BNG_lowbwt	0.56	0.57	0.56	0.55
1199_BNG_echoMonths	0.40	0.30	0.43	0.35
mean	0.58	0.48	0.53	0.46
p -value $>$		0.05	0.03	0.00

7 Conclusions

In this paper we explored the use of e-graphs and equality saturation as a mechanism to keep track of the history of the search engine of symbolic regression and exploiting its pattern matching capabilities to propose a variation to the traditional subtree crossover and mutation that increases the probability of generating novel expressions.

The e-graph data structure efficiently stores and queries for parts of expressions. Exploiting this capability, we modified the subtree operators to only sample subtrees that would form unvisited expressions. The expectation is that this simple modification would render a significant improvement in the search procedure.

We tested the proposed algorithm, called eggp, in 21 different benchmarks from the literature and compared with a simple GP with the original subtree operators, and two state-of-the-art algorithms, PySR and Operon. The results showed that the modified operators are capable of improving the performance of a simple GP to compete with the state-of-the-art. The main highlight of this

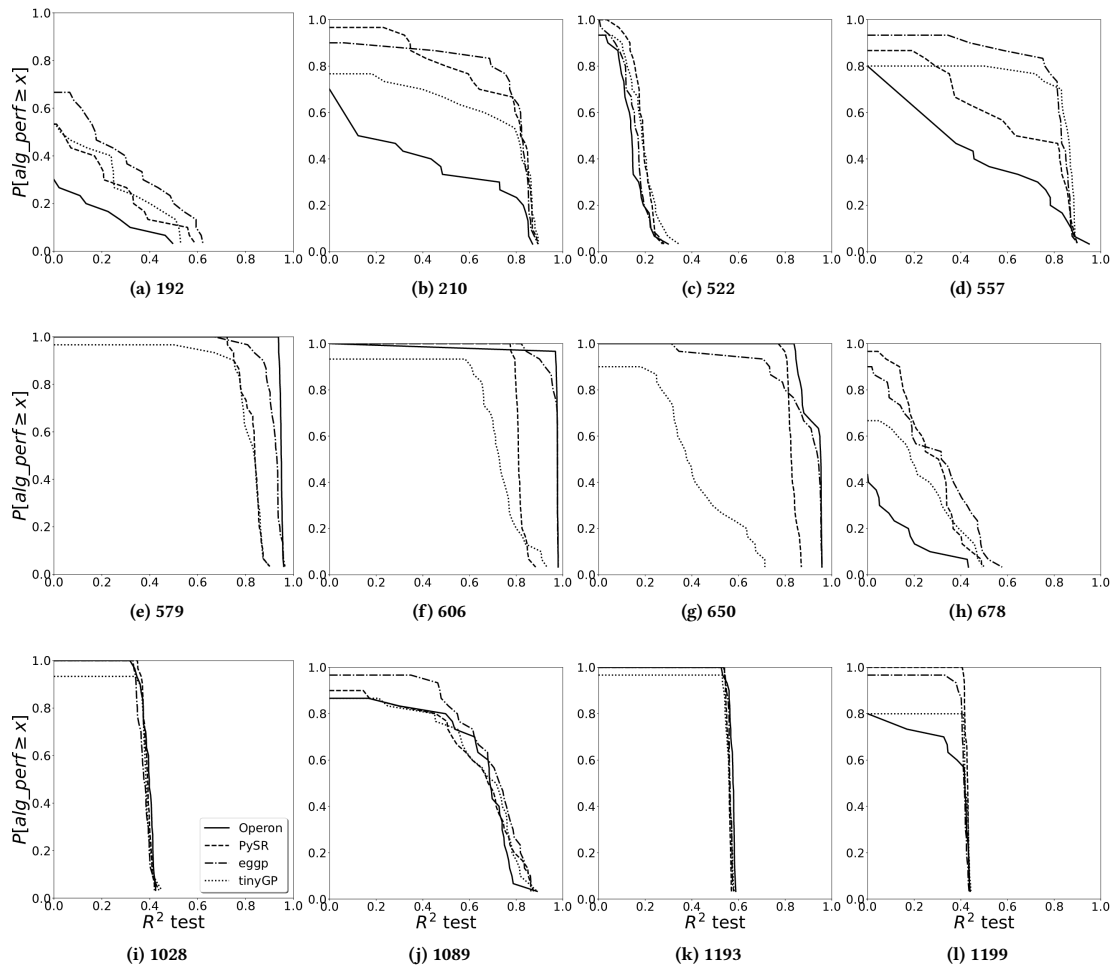


Figure 3: Performance plots for the SRBench datasets. This plot shows the probability of returning an R^2 equal or larger than x on a random run of each algorithm.

approach is that it consistently performs equal or better than the best competing approach.

The main limitations of these experiments lie in the use of default or reasonable values for the hyperparameters. We should notice that eggp contains eight parameters that should be adjusted accordingly in a practical scenario. Meanwhile, PySR contains about 30 hyperparameters that may affect the algorithm performance⁸ and Operon contains about 20 hyperparameters. With careful experimentation, both PySR and Operon could possibly achieve similar results. Regarding the runtime, eggp is consistently faster than PySR but significantly slower than Operon. When comparing with tinyGP, we can see that the use of e-graph and equality saturation does not increase the runtime significantly.

In conclusion, the expressiveness and capabilities of the e-graph data structure enabled us to make a simple modification to the original subtree operators while significantly improving the performance of GP for symbolic regression. This allowed us to obtain a

more robust algorithm delivering better performance more reliably than state-of-the-art implementations. As for the next steps, this same modifications can be applied to any other mutation operator used by the state-of-the-art algorithms. Not only that, but the e-graph opens up many new possibilities for improving the search as it allows us to query expressions with a combination of properties, which can translate to diversity-preserving population and easy to integrate prior knowledge [20, 29]. In addition, the storage of the search history allows us to analyze the learned building blocks and exploit this information to generate new solutions.⁹

Acknowledgments

F.O.F. is supported by Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq) grant 301596/2022-0.

G.K. is supported by the Austrian Federal Ministry for Climate Action, Environment, Energy, Mobility, Innovation and Technology,

⁸we are not considering hyperparameters that expect prior knowledge.

⁹The authors proudly made no use of LLMs for this work.

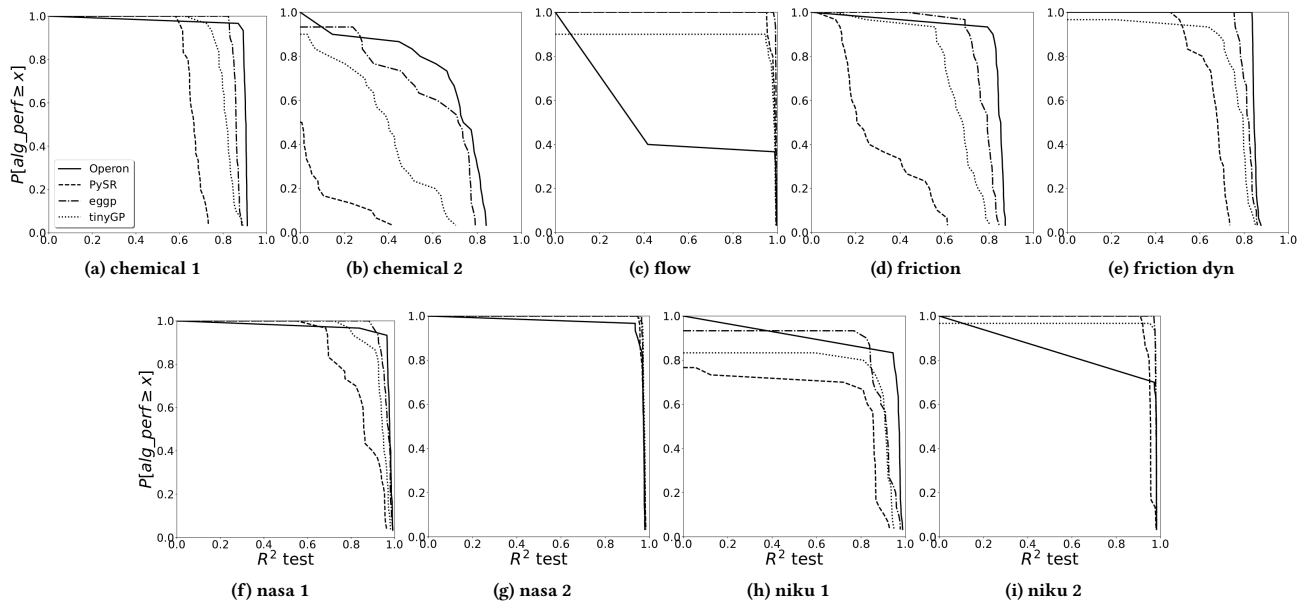


Figure 4: Performance plots for the real-world datasets. This plot shows the probability of returning an R^2 equal or larger than x on a random run of each algorithm.

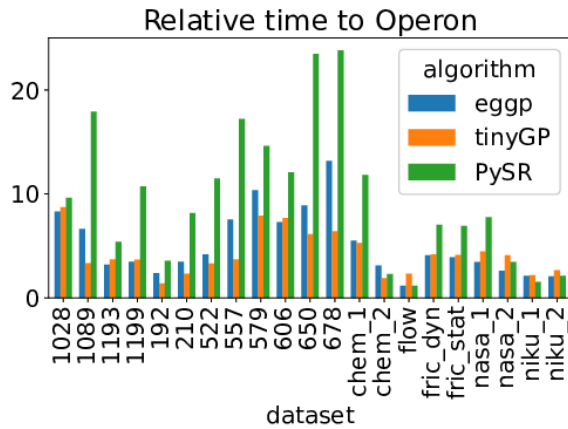


Figure 5: Relative avg. runtime of each algorithm using Operon as a baseline.

the Federal Ministry for Labour and Economy, and the regional government of Upper Austria within the COMET project ProMetHeus (904919) supported by the Austrian Research Promotion Agency (FFG).

No LLMs were used for this work.

References

- [1] Wolfgang Banzhaf, Ting Hu, and Gabriela Ochoa. 2024. How the Combinatorics of Neutral Spaces Leads Genetic Programming to Discover Simple Solutions. In *Genetic Programming Theory and Practice XX*. Springer, 65–86.
- [2] Deaglan J Bartlett, Harry Desmond, and Pedro G Ferreira. 2023. Exhaustive symbolic regression. *IEEE Transactions on Evolutionary Computation* (2023).
- [3] Bogdan Burlacu, Michael Affenzeller, Gabriel Kronberger, and Michael Kommenda. 2019. Online diversity control in symbolic regression via a fast hash-based tree similarity measure. In *2019 IEEE congress on evolutionary computation (CEC)*. IEEE, 2175–2182.
- [4] Bogdan Burlacu, Lukas Kammerer, Michael Affenzeller, and Gabriel Kronberger. 2020. Hash-based tree similarity and simplification in genetic programming for symbolic regression. In *Computer Aided Systems Theory–EUROCAST 2019: 17th International Conference, Las Palmas de Gran Canaria, Spain, February 17–22, 2019, Revised Selected Papers, Part I 17*. Springer, 361–369.
- [5] Bogdan Burlacu, Gabriel Kronberger, and Michael Kommenda. 2020. Operon C++ an efficient genetic programming framework for symbolic regression. In *Proceedings of the 2020 Genetic and Evolutionary Computation Conference Companion*. 1562–1570.
- [6] Haotian Cao, Garrett W Merz, Kyle Cranmer, and Gary Shiu. [n. d.]. Learning Conformal Field Theory with Symbolic Regression: Recovering the Symbolic Expressions for the Energy Spectrum. ([n. d.]).
- [7] Lulu Cao, Zimo Zheng, Chenwen Ding, Jinkai Cai, and Min Jiang. 2023. Genetic Programming Symbolic Regression with Simplification-Pruning Operator for Solving Differential Equations. In *International Conference on Neural Information Processing*. Springer, 287–298.
- [8] Miles Cranmer. 2023. Interpretable Machine Learning for Science with PySR and SymbolicRegression.jl. <https://doi.org/10.48550/ARXIV.2305.01582>
- [9] Fab rio Olivetti de Franca and Gabriel Kronberger. 2023. Reducing Overparameterization of Symbolic Regression Models with Equality Saturation. In *Proceedings of the Genetic and Evolutionary Computation Conference*. 1064–1072.
- [10] Fab rio O de Franca, Marco Virgolin, M Kommenda, MS Majumder, M Cranmer, G Espada, I Ingelse, A Fonseca, M Landajuela, B Petersen, et al. 2024. SRBench++: Principled benchmarking of symbolic regression with domain-expert interpretation. *IEEE transactions on evolutionary computation* (2024).
- [11] Kalyanmoy Deb, Samir Agrawal, Amrit Pratap, and Tanaka Meyarivan. 2000. A fast elitist non-dominated sorting genetic algorithm for multi-objective optimization: NSGA-II. In *Parallel Problem Solving from Nature PPSN VI: 6th International Conference Paris, France, September 18–20, 2000 Proceedings 6*. Springer, 849–858.
- [12] Marc Ebner, Mark Shackleton, and Rob Shipman. 2001. How neutral networks influence evolvability. *Complex*. 7, 2 (Nov. 2001), 19–33. <https://doi.org/10.1002/cplx.10021>
- [13] Steven Gustafson, Edmund K Burke, and Natalio Krasnogor. 2005. On improving genetic programming for symbolic regression. In *2005 IEEE Congress on Evolutionary Computation*, Vol. 1. IEEE, 912–919.
- [14] Ting Hu and Wolfgang Banzhaf. 2018. Neutrality, robustness, and evolvability in genetic programming. *Genetic Programming Theory and Practice XIV* (2018), 101–117.

Table 4: Ranks of the median (1st block) and AUC (2nd block) of the test set R^2 for the real-world. The p -values were calculated with a Wilcoxon signed-rank test using as alternative hypotheses ($\alpha = 0.05$) being greater ($>$) than eggp.

dataset	eggp	Operon	PySR	tinyGP
Chemical_1_tower	2	1	4	3
Chemical_2_competition	2	1	4	3
Flow_stress_phip0.1	1	4	2	3
Friction_dyn_one-hot	2	1	4	3
Friction_stat_one-hot	2	1	4	3
Nasa_battery_1_10min	2	1	4	3
Nasa_battery_2_20min	2	3	4	1
Nikuradse_1	2	1	4	3
Nikuradse_2	3	1	4	2
mean	2.00	1.56	3.78	2.67
p -value $>$		0.94	0.00	0.01
Chemical_1_tower	0.86	0.89	0.67	0.81
Chemical_2_competition	0.58	0.66	0.07	0.37
Flow_stress_phip0.1	0.99	0.51	0.98	0.88
Friction_dyn_one-hot	0.81	0.85	0.66	0.74
Friction_stat_one-hot	0.77	0.82	0.30	0.65
Nasa_battery_1_10min	0.96	0.96	0.85	0.93
Nasa_battery_2_20min	0.97	0.95	0.97	0.98
Nikuradse_1	0.85	0.89	0.62	0.75
Nikuradse_2	0.98	0.83	0.95	0.95
mean	0.87	0.82	0.67	0.79
p -value $>$		0.63	0.00	0.00

Table 5: Hypervolume of the negative R^2 and model size with a reference point of $[0.0, 50.0]$.

dataset	egraphGP	operon	PySR
srbench			
192	30.62 \pm 5.53	35.09 \pm 1.17	37.54 \pm 1.19
210	41.05 \pm 2.15	37.80 \pm 0.34	44.74 \pm 0.63
522	9.36 \pm 2.57	13.26 \pm 1.29	12.47 \pm 1.05
557	30.94 \pm 7.54	32.97 \pm 0.86	30.02 \pm 10.69
579	37.13 \pm 1.88	32.37 \pm 0.58	36.49 \pm 1.15
606	38.85 \pm 1.36	33.77 \pm 0.80	35.29 \pm 0.54
650	36.08 \pm 5.94	32.49 \pm 0.91	36.50 \pm 0.91
678	16.22 \pm 6.03	24.05 \pm 2.12	24.18 \pm 2.83
1028	14.85 \pm 3.74	16.37 \pm 0.27	17.98 \pm 0.29
1089	34.92 \pm 3.51	37.72 \pm 0.37	39.15 \pm 1.99
1193	24.98 \pm 0.87	23.61 \pm 0.18	25.14 \pm 0.27
1199	18.97 \pm 0.88	17.88 \pm 0.53	19.46 \pm 0.58
real-world			
Chemical_1	36.66 \pm 1.68	36.87 \pm 0.84	27.99 \pm 1.57
Chemical_2	28.52 \pm 5.63	31.64 \pm 0.99	22.48 \pm 2.91
Flow_stress	46.20 \pm 0.11	25.31 \pm 16.18	45.81 \pm 0.43
Friction_dyn	34.97 \pm 6.66	31.29 \pm 2.01	33.49 \pm 1.61
Friction_stat	30.17 \pm 10.29	28.79 \pm 3.23	28.56 \pm 1.60
Nasa_1	45.24 \pm 0.65	41.37 \pm 1.44	44.92 \pm 0.42
Nasa_2	40.76 \pm 13.82	39.53 \pm 2.23	43.99 \pm 0.40
Nikuradse_1	44.53 \pm 0.64	43.33 \pm 0.41	45.22 \pm 0.17
Nikuradse_2	40.31 \pm 0.69	33.29 \pm 2.02	40.26 \pm 0.85

- [15] Ting Hu, Gabriela Ochoa, and Wolfgang Banzhaf. 2023. Phenotype search trajectory networks for linear genetic programming. In *European Conference on Genetic Programming (Part of EvoStar)*. Springer, 52–67.
- [16] Robert E Keller and Wolfgang Banzhaf. 1999. The evolution of genetic code in genetic programming. In *Proceedings of the Genetic and Evolutionary Computation Conference*, Vol. 2. Morgan Kaufmann Orlando, 1077–1082.
- [17] John R. Koza. 1992. *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. MIT Press, Cambridge, MA, USA.
- [18] Gabriel Kronberger, Bogdan Burlacu, Michael Kommenda, Stephan M. Winkler, and Michael Affenzeller. 2024. *Symbolic Regression*. Chapman & Hall / CRC Press.
- [19] Gabriel Kronberger and Fabricio Olivetti de França. 2024. Effects of reducing redundant parameters in parameter optimization for symbolic regression using genetic programming. *Journal of Symbolic Computation* (2024), 102413.
- [20] Gabriel Kronberger, Fabricio Olivetti de França, Bogdan Burlacu, Christian Haider, and Michael Kommenda. 2022. Shape-constrained symbolic regression—improving extrapolation with prior knowledge. *Evolutionary Computation* 30, 1 (2022), 75–98.
- [21] Gabriel Kronberger, Fabricio Olivetti de Franca, Harry Desmond, Deaglan J. Bartlett, and Lukas Kammerer. 2024. The Inefficiency of Genetic Programming for Symbolic Regression. In *Parallel Problem Solving from Nature – PPSN XVIII*, Michael Affenzeller, Stephan M. Winkler, Anna V. Kononova, Heike Trautmann, Tea Tušar, Penousal Machado, and Thomas Bäck (Eds.). Springer Nature Switzerland, Cham, 273–289.
- [22] William La Cava, Bogdan Burlacu, Marco Virgolin, Michael Kommenda, Patryk Orzechowski, Fabricio Olivetti de França, Ying Jin, and Jason H Moore. 2021. Contemporary symbolic regression methods and their relative performance. *Advances in neural information processing systems* 2021, DB1 (2021), 1.
- [23] Kara D Lamb and Jerry Y Harrington. [n. d.]. Discovering How Ice Crystals Grow Using Neural ODE’s and Symbolic Regression. ([n. d.]).
- [24] Federico Lelli, Stacy S. McGaugh, James M. Schombert, and Marcel S. Pawlowski. 2017. One Law to Rule Them All: The Radial Acceleration Relation of Galaxies. *Astrophysical Journal* 836, 2, Article 152 (Feb. 2017), 152 pages. <https://doi.org/10.3847/1538-4357/836/2/152> arXiv:1610.08981 [astro-ph.GA]
- [25] Nicholas Freitag McPhee, Brian Ohs, and Tyler Hutchison. 2008. Semantic building blocks in genetic programming. In *Genetic Programming: 11th European Conference, EuroGP 2008, Naples, Italy, March 26-28, 2008. Proceedings 11*. Springer, 134–145.
- [26] J.F. Miller and S.L. Smith. 2006. Redundancy and computational efficiency in Cartesian genetic programming. *IEEE Transactions on Evolutionary Computation* 10, 2 (2006), 167–174. <https://doi.org/10.1109/TEVC.2006.871253>
- [27] Katta G Murty and Santosh N Kabadi. 1985. *Some NP-complete problems in quadratic and nonlinear programming*. Technical Report.
- [28] David L Randall, Tyler S Townsend, Jacob D Hochhalter, and Geoffrey F Bomarito. 2022. Bingo: a customizable framework for symbolic regression with genetic programming. In *Proceedings of the Genetic and Evolutionary Computation Conference Companion*. 2282–2288.
- [29] Julia Reuter, Viktor Martinek, Roland Herzog, and Sanaz Mostaghim. 2024. Unit-Aware Genetic Programming for the Development of Empirical Equations. In *International Conference on Parallel Problem Solving from Nature*. Springer, 168–183.
- [30] Daniel Rivero, Enrique Fernandez-Blanco, and Alejandro Pazos. 2022. DoME: A deterministic technique for equation development and Symbolic Regression. *Expert Systems with Applications* 198 (2022), 116712.
- [31] Etienne Russeil, Fabricio Olivetti de França, Konstantin Malanchev, Bogdan Burlacu, Emille Ishida, Marion Leroux, Clément Michelin, Guillaume Moinard, and Emmanuel Gangler. 2024. Multiview Symbolic Regression. In *Proceedings of the Genetic and Evolutionary Computation Conference*. 961–970.
- [32] Guilherme Seidyo Imai Aldeia, Fabricio Olivetti de Franca, and William G La Cava. 2024. Inexact Simplification of Symbolic Regression Expressions with Locality-sensitive Hashing. *arXiv e-prints* (2024), arXiv-2404.
- [33] M. Sipper. 2019. Tiny Genetic Programming in Python. https://github.com/moshesipper/tiny_gp.
- [34] Ross Tate, Michael Stepp, Zachary Tatlock, and Sorin Lerner. 2009. Equality saturation: a new approach to optimization. In *Proceedings of the 36th annual ACM SIGPLAN-SIGACT symposium on Principles of programming languages*. 264–276.
- [35] Silviu-Marian Udrescu and Max Tegmark. 2020. AI Feynman: A physics-inspired method for symbolic regression. *Science Advances* 6, 16 (2020), eaay2631.

- [36] Marco Virgolin and Solon P Pissis. [n. d.]. Symbolic Regression is NP-hard. *Transactions on Machine Learning Research* ([n. d.]).
- [37] Max Willsey, Chandrakana Nandi, Yisu Remy Wang, Oliver Flatt, Zachary Tatlock, and Pavel Panchekha. 2021. Egg: Fast and extensible equality saturation. *Proceedings of the ACM on Programming Languages* 5, POPL (2021), 1–29.