

Transformer Semantic Genetic Programming for d -dimensional Symbolic Regression Problems

PHILIPP ANTHES, Johannes Gutenberg University, Germany

DOMINIK SOBANIA, Johannes Gutenberg University, Germany

FRANZ ROTHLAUF, Johannes Gutenberg University, Germany

Transformer Semantic Genetic Programming (TSGP) is a semantic search approach that uses a pre-trained transformer model as a variation operator to generate offspring programs with controlled semantic similarity to a given parent. Unlike other semantic GP approaches that rely on fixed syntactic transformations, TSGP aims to learn diverse structural variations that lead to solutions with similar semantics. We find that a single transformer model trained on millions of programs is able to generalize across symbolic regression problems of varying dimension. Evaluated on 24 real-world and synthetic datasets, TSGP significantly outperforms standard GP, SLIM_GSGP, Deep Symbolic Regression, and Denoising Autoencoder GP, achieving an average rank of 1.58 across all benchmarks. Moreover, TSGP produces more compact solutions than SLIM_GSGP, despite its higher accuracy. In addition, the target semantic distance SD_t is able to control the step size in the semantic space: small values of SD_t enable consistent improvement in fitness but often lead to larger programs, while larger values promote faster convergence and compactness. Thus, SD_t provides an effective mechanism for balancing exploration and exploitation.

CCS Concepts: • **Software and its engineering** → **Genetic programming**; • **Computing methodologies** → *Supervised learning by regression*; • **Theory of computation** → *Program semantics*.

Additional Key Words and Phrases: Genetic Programming, Transformer Models, Semantic Operators, Symbolic Regression

1 Introduction

Regression tasks like symbolic regression (SR) aim to identify programs that accurately capture the underlying relationships in data. In SR, programs are represented as symbolic expressions constructed from operators (e.g. addition, multiplication), variables (e.g. x_1, x_2, x_3), and constants. The resulting syntactical composition embodies a certain program behavior that determines the way predictions (output values) are generated from the available data (input values) and, consequently, how the program fits to the underlying data.

Genetic programming (GP) has developed several approaches to search for high-quality programs. For example, standard GP (stdGP) uses syntactic variation operators completely ignoring the impact of syntactic modifications on program behavior. As a result, even minor syntactic changes can produce offspring with significantly different behavior, which often disrupts promising programs, leading to low search efficiency [McPhee et al. 2008; Moraglio et al. 2012; Uy et al. 2009]. To address this limitation, semantic-aware approaches incorporate program behavior (semantics) in the variation process, leading to smoother, often uni-modal, fitness landscapes [Vanneschi 2016]. One of the best-performing semantic approaches is Geometric Semantic Genetic Programming (GSGP), which uses Geometric Semantic Operators (GSOs) to generate offspring that are semantically similar to their parents [Moraglio et al. 2012].

However, GSOs operate by performing linear combinations of program structures, independent of the syntactic form of the parent. This reliance on fixed, predefined transformation rules limits their ability to express semantic similarity through diverse and adaptive syntactic variations. Consequently, GSO-based approaches tend to produce overly complex and bloated programs [Anthes et al. 2025; Martins et al. 2018; Vanneschi 2024]. GSOs do not take advantage of the fact that the same semantics can be achieved by many different syntactic forms and that numerous structurally efficient variations yield similar program behavior. An effective search strategy should operate directly in semantic space, independently of its syntactic representation.

A recently introduced approach, Transformer Semantic Genetic Programming (TSGP), employs a transformer-based variation operator trained on millions of semantic variations of programs to generate semantically similar offspring [Anthes et al. 2025]. In an offline *Model Building* step, a transformer model is trained to learn which syntactic transformations lead to similar semantic behavior. During training, pairs of programs with similar semantics are formed: the model takes one program as input and aims to synthesize the other as the target output, conditioned on the semantic distance SD between them. After training, the transformer acts as a zero-shot semantic variation operator: given a parent program and a target semantic distance SD_t , it generates a new offspring program that is semantically similar to the parent, where SD_t controls the degree of semantic similarity. Transformers are particularly well suited for this difficult task, having demonstrated their ability in numerous domains to capture contextual relationships in sequences and generate results that correspond to the desired behaviors, regardless of syntactic structure [Briesch et al. 2023; Brown et al. 2020; Vaswani et al. 2017].

Previous work introduced TSGP for fixed-dimension problems ($d=4$) and showed its superiority in quality and compactness [Anthes et al. 2025]. This paper extends that work in two key directions:

- (1) We show that a single transformer model can be trained to handle SR problems of varying dimensions d . By conditioning the model on the dimension d of the target problem, we ensure that valid expressions are generated, using at most d input variables. This enables evaluation across a comprehensive benchmark set of 24 real-world and synthetic SR problems of various dimensions.
- (2) We study how the target semantic distance SD_t influences the search process. We analyze its impact on resulting convergence speed and program size and demonstrate that adjusting

SD_t allows accurate control over the step size in the semantic space, enabling to better control the balance between exploration and exploitation.

In a comprehensive experimental study, we compare the performance of TSGP against stdGP, SLIM_GSGP, and two model-based search approaches: Deep Symbolic Regression (DSR) [Petersen et al. 2019] and Denoising Autoencoder GP (DAE-GP) [Wittenberg et al. 2020]. We find that TSGP is significantly better than the benchmark approaches, achieving an average rank of 1.58 across all datasets. Furthermore, TSGP significantly outperforms SLIM_GSGP, a GSGP variant with a deflation operator, in solution compactness. Our study of the target semantic distance SD_t reveals that it effectively controls the step size in the semantic space. A small SD_t results in small semantic steps, enabling consistent fitness improvements in each generation. However, as expected, this often results in slower convergence and limited exploration, thus producing larger programs. In contrast, a larger SD_t encourages broader exploration, leading to smaller and more compact programs with faster initial convergence. Yet, this comes at the cost of reduced exploitation, which can have a negative impact on the quality of the final solution.

Sect. 2 reviews related work in semantic GP, model-based search approaches, and transformer-based one-shot estimators. Sect. 3 presents the TSGP approach, including *Model Building* 3.1 and *Model Inference* 3.2. Sect. 4 describes the experimental setup used to evaluate the approaches. Sect. 5 reports the results, focusing on solution quality and program size, and further analyzes how the semantic step size is affected by different values of the target semantic distance SD_t . Sect. 6 summarizes the key findings and suggests directions for future research.

2 Related Work

Semantic GP. Semantic-aware approaches in GP aim to control the semantic similarity between offspring and their parent programs [Vanneschi et al. 2014a]. We focus specifically on methods that, like TSGP, generate new offspring from a single parent program. Early approaches indirectly influenced semantics during variation by rejecting offspring based on semantic criteria. For example, Semantically Driven Mutation (SDM) prevents the offspring from being semantically identical to the parent [Beadle and Johnson 2009], while Semantic Similarity-based Mutation (SSM) discards offspring that deviate significantly from the parent’s semantics [Uy et al. 2009]. Both approaches increase the performance of GP by consistently producing promising solution candidates.

The first semantic operators that directly generate offspring of controlled and weighted semantic similarity were introduced with GSGP [Moraglio et al. 2012]. Geometric Semantic Mutation (GSM) creates a new program by appending a program structure to the parent considering the mutation step size ms . The appended structure consists of two randomly generated programs that are subtracted from each other, ensuring only minor semantic modifications to the parent. Due to the efficient implementation suggested by [Vanneschi et al. 2013], GSGP could be successfully applied to regression tasks in domains such as pharmacokinetics and financial data analysis [Castelli et al. 2013; McDermott et al. 2014; Vanneschi et al. 2014b].

However, repeated addition of program structures leads to exponential growth in program size during the search. As a result, approaches have been proposed to mitigate this growth: One such approach is GSGP-Red, which simplifies programs by merging repetitive structures and coefficients that occur after applying GSOs [Martins et al. 2018]. In contrast, SLIM_GSGP argues that semantic similarity can be achieved not only by adding randomly generated structures with minimal semantic impact but also by removing them in later stages of the search. Therefore, SLIM_GSGP introduced a new geometric semantic deflation operator that systematically reduces the size of the offspring [Vanneschi 2024].

Neural Network-Guided Search. In recent years, neural networks have been proposed as search operators for symbolic regression tasks. By using feedback from previous search steps, these search operators generate new programs auto-regressively, token by token.

Deep Symbolic Regression (DSR) uses reinforcement learning (RL) to train a recurrent neural network (RNN) during the search process [Petersen et al. 2019]. The RNN acts as a policy network that sequentially samples symbolic expressions according to its learned policy. Each completed expression is evaluated on the target problem, and the resulting error is used to optimize the model parameters of the RNN. Later variants improve performance by adding large-scale pre-training or hybrid search strategies, such as neural-guided GP at decoding time [Landajuola et al. 2022].

In contrast, DAE-GP is an estimation-of-distribution algorithm that uses a denoising autoencoder LSTM to model the distribution of high-quality programs in the latent space of the LSTM [Wittenberg et al. 2020]. Once the LSTM has learned the distribution of the selected population, the LSTM is used to create a new offspring population: parental programs are slightly perturbed and passed through the network, generating new programs with similar syntactic properties. For SR problems, DAE-GP has demonstrated its ability to produce programs of similar quality compared to standard syntactic operators, while significantly reducing the complexity and size of the resulting structures [Wittenberg and Rothlauf 2023].

Both DSR and DAE-GP operate in an online manner, where neural networks are trained and optimized during the search. In contrast, TSGP’s model-based search operates offline and does not require training during the search, enabling faster and more efficient generation of new programs.

One-shot Transformers. Recent transformer-based approaches to symbolic regression operate in a one-shot manner. Given input-output data, such models generate candidate solutions in a single forward pass without an iterative search or refinement of programs [Biggio et al. 2021; Kamienny et al. 2022]. The models are trained on large synthetic datasets consisting of numerical inputs, symbolic expressions, and their corresponding outputs. During training, the transformer learns to map input-output pairs to symbolic expressions. At inference time, it produces a solution via a static direct mapping, rather than exploring the solution space [Biggio et al. 2021; Kamienny et al. 2022].

However, recent work demonstrates that one-shot models have great difficulty to generalize beyond their pre-training distribution and are not competitive with search-based baselines [Voigt et al. 2025]. Unlike search-based methods (such as stdGP, SLIM_GSGP or TSGP), which iteratively explore the space of possible programs and adapt to novel problem structures, one-shot transformers rely entirely on learned input-output-expression mappings. Furthermore, while many of these models are publicly available, the used training datasets are not published, making it impossible to assess whether test problems were inadvertently seen during training.

3 The TSGP Approach

TSGP is a semantic search approach that uses a pre-trained transformer model as a variation operator to generate semantically similar offspring. In a single, offline *Model Building* step, the model captures structural patterns between symbolic expressions (referred to as mathematical functions) that generate high semantic similarity. During *Model Inference*, the transformer model can transfer these learned patterns to mathematical functions for unknown datasets without being fine-tuned to the task.

3.1 Model Building

Model Building in TSGP consists of three steps: First, a diverse set of mathematical functions is generated. The generated functions should well cover the possible input space. Second, semantically

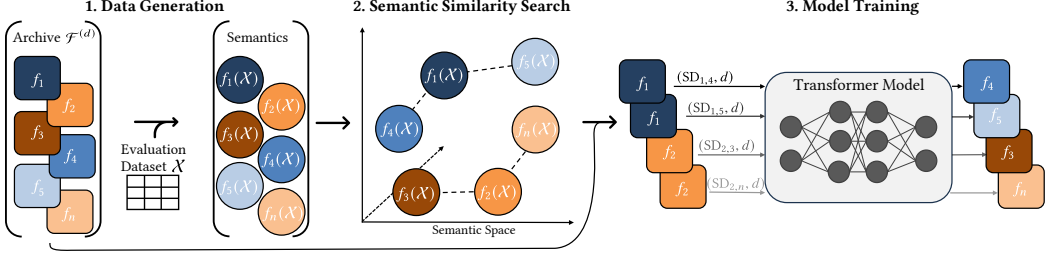


Fig. 1. *Model Building* of TSGP: (1) Diverse functions are generated and their semantics are approximated; (2) Semantically similar pairs are identified through a k -NN search in the semantic space; (3) These pairs are used as input-output examples to train a transformer model, conditioned on their semantic distance SD and the problem dimensionality d .

similar function pairs are identified and grouped to input-output training examples. Third, a transformer model is trained in a sequence-to-sequence task to generate a semantically similar output sequence based on the given input sequence. Each step is illustrated in Figure 1 and is described below.

Data Generation. In the first step of *Model Building*, we generate the data basis used to pre-train the transformer model. Unlike [Anthes et al. 2025], which trained a model for fixed-dimension problems, we build a single transformer model capable of generalizing across SR problems of varying dimensionality. Thus, instead of relying on a single archive of fixed dimension, we now create multiple archives $\mathcal{F}^{(d)} = \{f_1, f_2, \dots, f_{n_d}\}$, where each archive consists of n_d unique symbolic functions valid for d -dimensional SR problems. The functions in each archive $\mathcal{F}^{(d)}$ are obtained by applying stdGP to a diverse set of d -dimensional, synthetically generated SR problems. During each GP run, numerous programs valid for d -dimensional problems are generated, which are all collected in $\mathcal{F}^{(d)}$. The underlying SR problems are constructed from linear regression models with added Gaussian noise and inputs are sampled from a standardized distribution with zero mean and unit variance. To reduce redundancy and algebraic complexity in $\mathcal{F}^{(d)}$, we apply a simplification step using SymPy [Meurer et al. 2017]. This step normalizes expressions symbolically, minimizing syntactic variation that does not affect semantics.

We approximate the semantics of each function $f \in \mathcal{F}^{(d)}$ by evaluating it on a randomly sampled, but fixed dataset $\mathcal{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m]$, $\mathbf{x}_i \in \mathbb{R}^d$, where m is the number of evaluation points and d the input dimensionality. The dataset \mathcal{X} is constructed by drawing m input points independently from a zero-mean, unit-variance Gaussian distribution across all dimensions. With m sufficiently large, \mathcal{X} densely covers the space of plausible inputs under standardized conditions. Thus, the resulting output $f(\mathcal{X}) = (f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_m))$ represents a semantic vector $\mathbf{s}(f) \in \mathbb{R}^m$, which serves as a comprehensive numerical approximation of f 's true behavior [Krawiec and Lichocki 2009; McPhee et al. 2008; Uy et al. 2011].

Semantic Similarity Search. Based on the computed semantic vectors, semantically similar functions are identified through a semantic similarity search. Two functions $f_i, f_j \in \mathcal{F}^{(d)}$ are considered semantically similar if the distance between their semantic vectors $\mathbf{s}(f_i)$ and $\mathbf{s}(f_j)$ is small [Uy et al. 2011]. Since semantics is represented using vectors in \mathbb{R}^m , similarity can be quantified using a distance measure [Moraglio et al. 2012]. We use the Euclidean distance and define the semantic distance (SD) between two semantic vectors f_i and f_j as $SD(f_i, f_j) = \|\mathbf{s}(f_i) - \mathbf{s}(f_j)\|_2$. To identify similar functions, we perform a k -nearest neighbors (k -NN) search on the semantics for each

archive $\mathcal{F}^{(d)}$. For each function $f_i \in \mathcal{F}^{(d)}$, the k -NN search returns a set of k neighbors $N_k(f_i)$ with minimal semantic distance:

$$N_k(f_i) = \underset{\{f' \in \mathcal{F}^{(d)}, |f'|=k\}}{\operatorname{argmin}} \sum_{f \in f'} \operatorname{SD}(f_i, f). \quad (1)$$

From these sets of neighbors, we generate k input–output training pairs (f_i, f_o) per function, where f_i is the input, and each $f_o \in N_k(f_i)$ serves as a semantically similar target output.

The similarity search is performed separately for each archive $\mathcal{F}^{(d)}$. This yields $|\mathcal{F}^{(d)}| \times k$ input–output training pairs per archive. Given D archives of comparable size, the total number of training pairs is roughly $D \times |\mathcal{F}^{(d)}| \times k$.

Model Training. In the third step, a single encoder–decoder transformer model is trained on all input–output pairs (f_i, f_o) collected from the semantic similarity searches across all archives $\mathcal{F}^{(d)}$. The model uses the architecture proposed by [Vaswani et al. 2017], with randomly initialized parameters θ .

The main training objective is to generate the semantically similar output function f_o given the input function f_i . To achieve this, the input function is first encoded into the Transformer’s latent representation. The output function is then generated step by step in an auto-regressive manner: at each step, the model predicts the next token based on the latent representation of the input and the tokens already generated from the current sequence. The model parameters θ are updated in batches by minimizing the cross-entropy loss between the predicted token and the target token of the output sequence f_o . Upon convergence, the optimized parameters θ^* are stored for inference.

During training, two additional inputs are provided alongside f_i to control the semantic similarity and the dimensional validity of the output. The first is the semantic distance $\operatorname{SD}(f_i, f_j)$ between the input–output pair, which allows the model to quantitatively learn semantic similarity. In inference, conditioning the model on a target semantic distance SD_t controls how closely the generated offspring matches the parent’s behavior. The second input is the dimensionality d of the archive $\mathcal{F}^{(d)}$ from which the pair was drawn. Since the model only receives function pairs of dimensionality less than or equal to d , it learns to generate expressions that depend on not more than d variables. In inference, setting d to the dimensionality of the problem guarantees that the outputs are valid for the input space.

The result of *Model Building* is a single transformer model that can generate expressions with controlled semantic similarity to a given input without further fine-tuning, across problems of varying dimension.

3.2 Model Inference

After *Model Building*, the pre-trained transformer model with parameters θ^* is used during the evolutionary search as a semantic variation operator, replacing the standard GP variation operators.

During variation, the entire population is processed as a single batch and independently generates for each parent function f_i a new offspring f_o , conditioned on the target semantic distance SD_t and the problem dimensionality d . The parameter SD_t controls how closely the behavior of the offspring matches that of the parent, thus allowing the step size in the semantic space to be regulated. Conditioning the model on the dimensionality of the problem d ensures that the generated expression uses only variables x_1 through x_d , guaranteeing compatibility with the d -dimensional input of the target problem.

Function generation proceeds auto-regressively: the model simultaneously predicts the next token for every function in the batch, conditioned on the respective parent function and the context of the tokens generated so far. Inference stops when either an end-of-sequence (EOS) token is

created for every function in the batch, or the maximum sequence length is reached. To ensure syntactic validity of the generated output, we employ syntax control as proposed by [Wittenberg 2022], which enforces syntactic correctness during auto-regressive generation. At every decoding step, token probabilities are modified to assign zero probability to syntactically invalid options, such as operators with incorrect arity.

4 Experimental Settings

Datasets. We compare the performance of TSGP with stdGP, SLIM_GSGP, DSR, and DAE-GP on 24 SR problems taken from Penn Machine Learning Benchmarks (PMLB) [Romano et al. 2021], a widely used benchmark suite in SR research, which is also included in large-scale studies such as SRBench [La Cava et al. 2021]. For our experiments, we select all available real-world datasets with a dimensionality ranging from 2 to 5 dimensions where a sufficient number of samples (> 100) are available. We complement these 12 real-world datasets with 12 synthetic datasets, including the Pollen dataset used in [Anthes et al. 2025] and randomly sampled Feynman equations spanning various dimensionalities. On the Feynman datasets, solution quality is evaluated using a fixed set of 10,000 samples to ensure a reliable evaluation. The datasets are listed in Table 1.

Before training and evaluation, the feature and target variables in the datasets are standardized to have zero mean and unit variance. This is in accordance with the standardized input space used during TSGP’s *Model Building* step, ensuring the transferability of the pre-trained model to datasets of different value ranges. Standardization also improves solution quality and reduces bloat in GP-based methods by making the search process invariant to the scale of input features and constants [Dick et al. 2020; Owen et al. 2018].

Search Settings. For TSGP, stdGP and DAE-GP, we use the evolutionary computation framework DEAP [Fortin et al. 2012], with neural networks being implemented with Keras [Chollet et al. 2015]. For SLIM_GSGP and DSR, we use publicly available Python implementations [Petersen et al. 2019; Vanneschi 2024]. GP parameters that are common in all methods are described in Table 2.

Initial populations are generated using Ramped Half-and-Half (RHH). For TSGP, stdGP and DAE-GP, which are based on the DEAP framework, we choose an initialization depth between 2 and 5. For SLIM_GSGP, we rely on the default initialization values specified by the framework. We set the maximum allowed tree depth during a run to 17 [Koza 1993].

The terminal set includes all input variables $V = \{v_1, v_2, \dots, v_d\}$ that correspond to the dimensionality d of the data set and ephemeral random constants (ERCs) sampled from $[-0.5, 0.5]$ in steps of 0.1. The function set consists of addition, subtraction, multiplication, protected division (returning 1 on division by zero) and pow. We include pow because Sympy represents division as exponentiation [Meurer et al. 2017], which leads to having pows in the training data of TSGP. Although the pow operator is not available in the SLIM_GSGP or DSR frameworks, its inclusion in stdGP did not affect performance compared to previous results [Anthes et al. 2025].

The population size is set to 100 and the search is conducted for 100 generations. Parents are selected via tournament selection with a tournament size of 5 and no elitism is applied. Each dataset is divided into training and test sets in a 75% to 25% ratio. The best program (i.e., the one with the lowest training RMSE) is extracted as the final solution and evaluated on the test set to measure its prediction quality. We use Root Mean Squared Error (RMSE) as the evaluation metric, which quantifies the Euclidean distance between the function’s semantics $f(X)$ and the target semantics.

TSGP-Setup. To create training data for the transformer model of TSGP, we apply stdGP to 50 synthetic regression problems for each dimensionality $d \in \{2, 3, 4, 5\}$. Each stdGP run uses a population of 2,000 and continues until 150,000 unique functions are collected. All other parameters match those specified in Table 2. To balance accuracy and simplicity and reduce bloat during

Table 1. The benchmark datasets consist of 12 real-world and 12 synthetic problems with dimensionality ranging from 2 to 5 and sample sizes from 100 to 10,000.

| Data set | # Samples | Dimensionality | Type |
|---------------------------|-----------|----------------|------------|
| Analcatdata Apnea1 | 475 | 3 | real-world |
| Analcatdata Apnea2 | 475 | 3 | real-world |
| Analcatdata Neavote | 100 | 2 | real-world |
| Chscase Geyser1 | 222 | 2 | real-world |
| Cloud | 108 | 5 | real-world |
| ERA | 1000 | 4 | real-world |
| ESL | 488 | 4 | real-world |
| Feynman_I_18_4 | 10000 | 4 | synthetic |
| Feynman_I_24_6 | 10000 | 4 | synthetic |
| Feynman_I_25_13 | 10000 | 2 | synthetic |
| Feynman_I_29_4 | 10000 | 2 | synthetic |
| Feynman_I_43_43 | 10000 | 4 | synthetic |
| Feynman_II_34_2 | 10000 | 3 | synthetic |
| Feynman_II_38_3 | 10000 | 4 | synthetic |
| Feynman_II_4_23 | 10000 | 3 | synthetic |
| Feynman_III_14_14 | 10000 | 5 | synthetic |
| Feynman_test_3 | 10000 | 4 | synthetic |
| Feynman_test_4 | 10000 | 5 | synthetic |
| Galaxy | 323 | 4 | real-world |
| LEV | 1000 | 4 | real-world |
| Pollen | 3848 | 4 | synthetic |
| Rabe | 120 | 2 | real-world |
| Vinnie | 380 | 2 | real-world |
| Visualizing Environmental | 111 | 3 | real-world |

Table 2. Configuration parameters and values used for the evaluated methods.

| Parameter | Value |
|-------------------|--------------------------------|
| Initialization | Ramped Half-and-Half |
| Primitive Set | {V, ERC, +, -, ×, %, pow} |
| ERC Range | [-0.5, 0.5], stepsize 0.1 |
| Population Size | 100 |
| Generations | 100 |
| Selection | Tournament selection of size 5 |
| Evaluation Metric | Root Mean Squared Error (RMSE) |
| Runs | 30 |

evolution, we employ double tournament selection [Luke and Panait 2002]. Generated expressions are simplified using SymPy, except for numerical constants, which are preserved in their original form to be compatible with the ERC values of the search and to ensure comparability with the other benchmarked algorithms. This results in a vocabulary for TSGP consisting of the primitive set and ERCs (as in Table 2), variables up to dimensionality 5, integers from -5 to 5 (produced by

SymPy combining sub-structures), additional token placeholder for the target semantic distance SD_t and dataset dimensionality d . The full vocabulary is provided in Appendix A

Pairs of programs with high semantic similarity are obtained using a k -nearest neighborhood search ($k=3$) based on the Euclidean distance between semantic vectors, each computed on a fixed, standardized dataset of 500 input points. To enable an efficient similarity search, we use the FAISS library [Douze et al. 2024], which reduces computational costs by clustering semantic vectors and computing distances only within each cluster. An input–output pair (f_i, f_o) is created for each function f_i and its neighbor f_o if $0 < SD(f_i, f_o) < 100$ and if both functions are within the maximum length of 100 tokens. This yields approximately 5 million pairs per dimensionality, resulting in a total training set of 20 million function pairs, which is 4 times higher than used in [Anthes et al. 2025], where TSGP was used for SR problems of only one size ($d = 4$).

The transformer follows the architecture of [Vaswani et al. 2017] with 8 attention heads, a hidden dimension of 128, and 2 encoder–decoder layers, resulting in 3.8 million parameters. We use the AdamW optimizer [Loshchilov and Hutter 2017] and implement a cosine learning rate scheduler with a learning rate of 10^{-3} , which gradually reduces the learning rate over the training duration of 8 epochs. After training, the same trained model is used in all experiments.

In our experiments, we evaluate TSGP under different settings of the target semantic distance SD_t . For clarity, we denote the variant with $SD_t = 1$ as TSGP1, which serves as our default configuration. Other variants are labeled accordingly (e.g. TSGP0.1 for $SD_t = 0.1$, TSGP5 for $SD_t = 5$). Unless otherwise specified, “TSGP” in the results refers to TSGP1.

Baseline Configurations. StdGP uses subtree crossover with a probability of 90% and a bias toward selecting terminals of 10%, along with subtree mutation applied with a probability of 10%, where newly generated subtrees have a depth ranging from 0 to 2 [Koza 1993].

For SLIM_GSGP, we use the default configuration of the framework. Thus, the SLIM+SIG2 variant is employed and no geometric crossover is performed, resulting in a one-parent-based semantic search that makes it comparable to TSGP. The default inflation mutation rate is set to 0.2, defining the probability of selecting the inflation mutation over the deflate operation during the mutation process of a program.

DSR follows the default configuration of the framework. The reward is defined as the negative root mean squared error (RMSE) to match the metric to measure solution quality, with a batch size of 100 and trained for 100 iterations to align with the settings used in GP.

DAE-GP is implemented as suggested by [Wittenberg et al. 2023], using an LSTM-based autoencoder with a single hidden layer whose dimensionality is dynamically adjusted. Training is performed using the Adam optimizer with a learning rate of 0.001 until convergence. Input corruption is applied using Levenshtein tree edit at a high edit percentage of 95%, which promotes extensive structural perturbations and encourages exploration in the search space, helping to counteract premature convergence [Wittenberg et al. 2023]

5 Experiments

5.1 Prediction Quality & Convergence

We start our analysis with a performance comparison of TSGP against the benchmark methods stdGP, SLIM_GSGP, DSR, and DAE-GP. For TSGP, we set the target semantic distance to $SD_t = 1$ and denote this setting as TSGP1. In each run, the best-performing program on the training set after 100 generations is selected, and its generalization performance is evaluated using the root mean squared error (RMSE) on the held-out test set. Table 3 reports the median test RMSE over 30 independent runs for each dataset. Additionally, we compute the average rank of each algorithm across all datasets, along with the standard deviation. The best results are highlighted in bold

Table 3. Median test RMSE of the best programs (solutions) identified within 100 generations for TSGP with $SD_t = 1$, stdGP, SLIM_GSGP (SLIM), DSR, and DAE-GP (DAE) for the 24 analyzed datasets. Bold values indicate the best prediction quality (lowest RMSE). Significant differences of the best results are indicated by the label symbols.

| Data set | <i>a</i> TSGP1 | <i>b</i> stdGP | <i>c</i> SLIM | <i>d</i> DSR | <i>e</i> DAE |
|------------------------------|-----------------------------|---------------------------|--------------------------|---------------------------|--------------|
| Analcatdata Apnea1 | 0.9592 | <i>e</i> 0.8176 | 0.9643 | 1.0705 | 1.0662 |
| Analcatdata Apnea2 | 0.9277 | <i>acde</i> 0.6785 | 1.0287 | 1.2835 | 1.0234 |
| Analcatdata Neavote | <i>c</i> 0.2358 | 0.2599 | 0.3536 | 0.2383 | 0.2505 |
| Chscase Geyser1 | 0.4988 | 0.4906 | 0.4969 | <i>c</i> 0.4732 | 0.4860 |
| Cloud | 0.4601 | 0.5653 | 0.5180 | <i>bce</i> 0.3557 | 0.6248 |
| ERA | 0.8006 | 0.8250 | <i>bde</i> 0.7943 | 0.9003 | 0.8979 |
| ESL | <i>bde</i> 0.3781 | 0.5192 | 0.4138 | 0.4613 | 0.5665 |
| Feynman I 18 4 | <i>bcd</i> 0.1679 | 0.3124 | 0.2059 | 0.4354 | 0.5344 |
| Feynman I 24 6 | <i>bcd</i> 0.1617 | 0.3752 | 0.2838 | 0.6824 | 0.7271 |
| Feynman I 25 13 | 0.1893 | 0.2078 | <i>de</i> 0.1891 | 0.4736 | 0.5152 |
| Feynman I 29 4 | <i>cde</i> 0.2843 | 0.3261 | 0.3405 | 0.5770 | 0.5922 |
| Feynman I 43 43 | <i>bcd</i> 0.3703 | 0.5182 | 0.4312 | 0.8031 | 0.8103 |
| Feynman II 34 2 | <i>bcd</i> 0.0881 | 0.2253 | 0.2014 | 0.6618 | 0.6854 |
| Feynman II 38 3 | <i>bcd</i> 0.3141 | 0.5260 | 0.4123 | 0.7659 | 0.8289 |
| Feynman II 4 23 | <i>bcd</i> 0.2836 | 0.3926 | 0.3936 | 0.7245 | 0.7713 |
| Feynman III 14 14 | <i>bcd</i> 0.4377 | 0.5384 | 0.5405 | 0.9377 | 0.8753 |
| Feynman Test 3 | <i>bcd</i> 0.1776 | 0.2391 | 0.2329 | 0.4765 | 0.5511 |
| Feynman Test 4 | <i>cde</i> 0.2211 | 0.2864 | 0.2490 | 0.4698 | 0.5779 |
| Galaxy | <i>cde</i> 0.2824 | 0.3145 | 0.3435 | 0.3874 | 0.4491 |
| LEV | <i>bcd</i> 0.6568 | 0.7166 | 0.6629 | 0.8223 | 0.8265 |
| Pollen | <i>bcd</i> 0.4700 | 0.5038 | 0.5139 | 0.7259 | 0.8005 |
| Rabe | <i>bcd</i> 0.0929 | 0.1427 | 0.2243 | 0.2454 | 0.3005 |
| Vinnie | 0.5127 | 0.5154 | 0.5121 | <i>abce</i> 0.4862 | 0.5235 |
| Visualizing Environmental | 0.8213 | 0.8122 | 0.8391 | <i>abce</i> 0.6871 | 0.8174 |
| Rank (Mean & Std) | <i>bcd</i> 1.58±1.04 | 2.71±0.84 | 2.67±1.03 | 3.54±1.29 | 4.50 ±0.87 |

font and are tested for statistical significance using the Mann–Whitney U test with $\alpha = 0.05$ and Bonferroni correction for multiple comparisons. The method labels (*a,b,c,d,e*) highlight statistically significant differences to the worse performing methods.

The results show that TSGP achieves the best overall performance, significantly outperforming all baselines. Across all datasets, TSGP achieves the lowest average rank (1.58), with statistically significant improvements over stdGP (rank 2.71), SLIM_GSGP (rank 2.67), DSR (rank 3.54), and DAE-GP (rank 4.50).

TSGP’s strong performance is robust across dataset characteristics. TSGP consistently produces the best solutions regardless of dimensionality, data type (synthetic or real-world), or sample size. In contrast, the other model-based methods DSR and DAE-GP perform overall poorly, with average ranks of 3.54 and 4.50, respectively. Both rely on online model training during the evolutionary search, which requires a large amount of training samples within the search to generate promising candidate solutions. TSGP, on the other hand, works offline and uses a pre-trained model, making it significantly more efficient at finding high-quality solutions. Given the inferior performance of

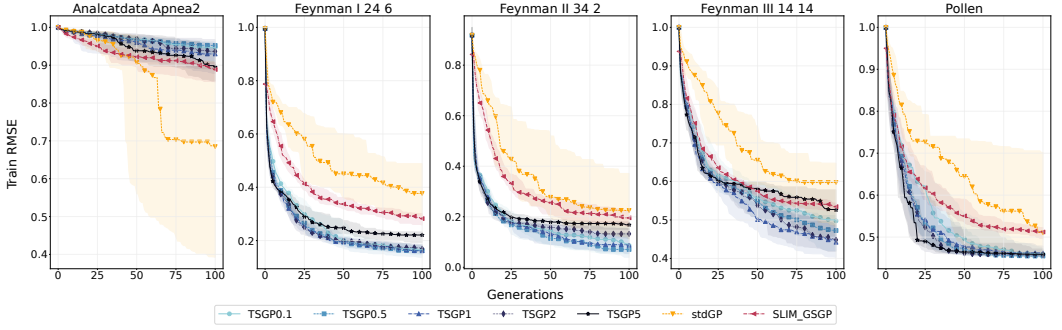


Fig. 2. Median training RMSE of the best programs (solutions) of TSGP, stdGP, SLIM_GSGP over generations on a subset of the analyzed datasets.

DSR and DAE-GP, we focus the remainder of our analysis on the methods that perform the best, TSGP, stdGP, and SLIM_GSGP.

We study the convergence speed of TSGP, stdGP, and SLIM_GSGP to assess how quickly each algorithm discovers programs with low error during the search. This reflects the effectiveness of their variation operators in improving candidate solutions. Consequently, Figure 2 plots the median training RMSE of the best programs over the number of generations. The interquartile range (IQR) indicates performance variability across 30 runs. For clarity, we show five representative datasets; results for the remaining datasets are provided in Appendix B. In addition to the standard setting (TSGP1), we analyze TSGP under varying SD_t , which controls the degree of semantic similarity between parent and offspring. Small $SD_t = 0.1$ encourages high semantic similarity, while larger values (up to $SD_t = 5$) promote greater dissimilarity. These variants are labeled in the legend and visualized using a continuous color scale from light blue (low SD_t) to black (high SD_t).

We start the analysis with the standard variant TSGP1, setting $SD_t = 1$. TSGP1 converges significantly faster than stdGP and SLIM_GSGP for almost all datasets and quickly identifies high-quality solutions, which are continuously improved throughout the search. For example, on Feynman I 24 6, TSGP1 finds better solutions than stdGP within just 9 generations and surpasses the termination performance of SLIM_GSGP by generation 21. In contrast to semantic-based methods, stdGP exhibits slower convergence, reflecting the inefficiency of syntactic variation operators that ignore the semantic effects.

Furthermore, the IQR of TSGP and SLIM_GSGP is notably smaller than that of stdGP, indicating greater stability and consistency in solution quality across runs. This stability arises because semantic-aware methods explore directly in the semantic solution space and are not based on random structural changes that have uncertain behavioral effects, leading to more consistent and reliable results over multiple runs.

The analysis of TSGP under varying SD_t values reveals a trade-off between convergence speed and final solution quality. TSGP with small SD_t often converges more slowly but achieves the highest final performance. In contrast, high SD_t often leads to faster initial progress (especially on datasets like Pollen and Feynman III 14 14), but frequently results in early plateaus (e.g. on Feynman I 24 6), limiting further improvement. TSGP with $SD_t = 1$ offers a robust balance, achieving consistently strong convergence behavior and high prediction accuracy across a wide range of problems.

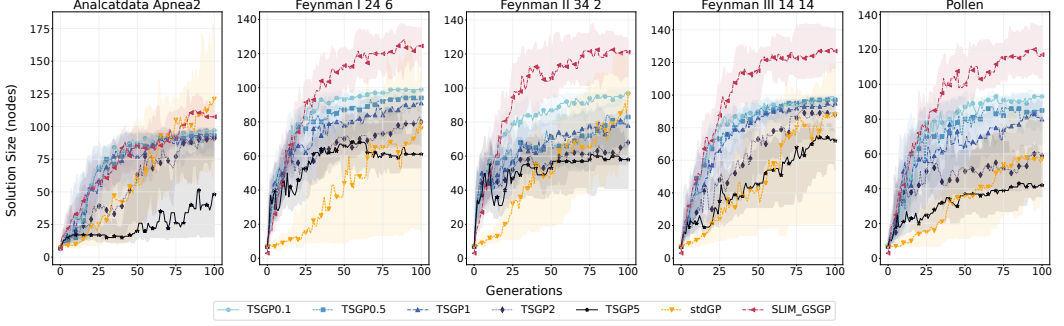


Fig. 3. Median size of the solutions of TSGP, stdGP, SLIM_GSGP over generations on a subset of the analyzed datasets.

5.2 Solution Size

Beyond prediction accuracy, program size is a relevant property of solutions, as smaller ones tend to be more interpretable [Javed et al. 2022; Wittenberg and Rothlauf 2023]. Usually, the size of symbolic expressions is measured by counting the nodes in the expression tree. Table 4 reports the median size of the solutions after 100 generations. First, we compare the solution sizes (sizes of the best programs) generated by the default variant TSGP1 with those produced by stdGP and SLIM_GSGP. The results are listed in Table 4. As before, these values represent the median of 30 independent runs. The best results (smallest solutions) are shown in bold and are tested for significant differences with the other solution sizes of the other methods. Significant differences are marked with representative labels (*a, b, c*).

The results show that TSGP1 generates the smallest solutions on average, with a mean rank of 1.58, followed by stdGP (1.75) and SLIM_GSGP (2.67). This is particularly noteworthy because TSGP1 achieves this compactness while also achieving the best predictive performance.

Furthermore, there are significant differences when comparing the semantic methods. TSGP1, which uses a transformer model to generate semantic similarities, produces significantly smaller solutions than SLIM_GSGP, which uses linear combination-based GSOs for variation. Being trained on millions of diverse syntactic variations that lead to semantically similar offspring, TSGP discovers solutions that are not only accurate but also more compact than those generated by GSO-based approaches.

Figure 3 plots the median size of the best programs over the number of generations. TSGP1 and its variants with varying SD_t are compared with stdGP and SLIM_GSGP. The plot displays the median values and IQR across 30 independent runs. For consistency and clarity, the same subset of datasets shown in the previous section is used; the results for the remaining datasets are provided in the Appendix C.

We first compare TSGP1 with the baseline methods. In early generations, where TSGP1 quickly identifies high-quality solutions, the solution size increases significantly, comparable to the growth observed with SLIM_GSGP. However, this growth decreases in later generations, even though the quality of the solutions continuously improves. This indicates that the transformer model of TSGP successfully learned a variety of structural semantic modifications, allowing it to create semantically similar offspring of similar size to the parent. SLIM_GSGP also shows reduced growth in later generations, due to its deflate operator, but the effect is much weaker. This is because SLIM_GSGP relies on fixed syntactic rules for variation, which limits its ability to adaptively modify the structure of the program while preserving semantics.

Table 4. Median program size of the solutions identified within 100 generations by TSGP with $SD_t = 1$, stdGP, and SLIM_GSGP (SLIM). Values in bold indicate the smallest solutions. Significant differences with respect to all other methods are indicated by the label symbols.

| Data set | $_a$ TSGP1 | $_b$ stdGP | $_c$ SLIM |
|------------------------------|--------------------------------------|----------------|-----------------|
| Analcata Apnea1 | $_c$ 90 | 94 | 120 |
| Analcata Apnea2 | $_c$ 92 | 121 | 108 |
| Analcata Neavote | 86 | 105 | 84 |
| Chscase Geyser1 | 87 | 113 | 71 |
| Cloud | $_c$ 74 | 97 | 107 |
| ERA | 85 | 72 | 90 |
| ESL | 79 | ac 7 | 100 |
| Feynman I 18 4 | 64 | $_c$ 25 | 100 |
| Feynman I 24 6 | 91 | $_c$ 76 | 124 |
| Feynman I 25 13 | $_c$ 83 | 108 | 122 |
| Feynman I 29 4 | bc 81 | 101 | 128 |
| Feynman I 43 43 | 93 | $_c$ 57 | 126 |
| Feynman II 34 2 | $_c$ 78 | 96 | 121 |
| Feynman II 38 3 | 90 | 86 | 126 |
| Feynman II 4 23 | $_c$ 97 | 108 | 122 |
| Feynman III 14 14 | 95 | $_c$ 88 | 127 |
| Feynman Test 3 | 83 | $_c$ 76 | 106 |
| Feynman Test 4 | $_c$ 69 | 73 | 104 |
| Galaxy | $_c$ 86 | 92 | 98 |
| LEV | $_c$ 85 | 98 | 109 |
| Pollen | 80 | ac 57 | 117 |
| Rabe | 83 | 63 | 80 |
| Vinnie | 78 | 99 | $_a$ 50 |
| Visualizing Environmental | $_c$ 93 | 99 | 127 |
| Rank (Mean & Std) | $_c$ 1.58\pm0.57 | 1.75 \pm 0.9 | 2.67 \pm 0.69 |

For the different TSGP variants, Figure 3 reveals that the solution size depends on SD_t . TSGP with larger SD_t produces significantly smaller solution structures than TSGP with small SD_t . This effect is particularly evident on the Pollen dataset, where all variants of TSGP achieve similar final predictive performance at termination, yet their solution sizes differ drastically, ranging from 42 nodes with $SD_t = 5$ to 93 nodes with $SD_t = 0.1$. This illustrates that larger semantic steps lead to the discovery of more compact expressions, while small steps lead to higher solution quality at the cost of increased structural complexity.

5.3 Step Size in the Semantic Space

The step size measures the extent of the behavioral changes introduced during variation. It is relevant for the interplay between exploration and exploitation during search. Large step sizes favor exploration, enabling search to explore new areas in the solution space leading to a higher probability of finding the actual global optima [Rothlauf 2011]. In contrast, smaller step sizes lead to stronger exploitation (local search) and a gradual improvement in solution quality. Therefore, a successful search strategy must balance both aspects. In TSGP, this balance is controlled by the

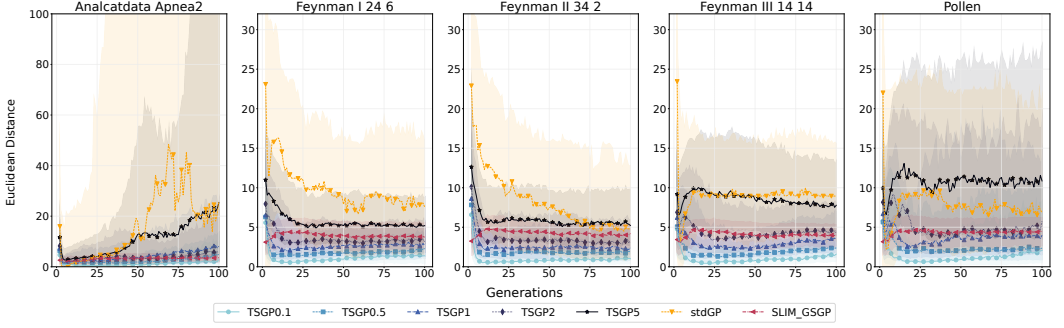


Fig. 4. Median Euclidean distance between the semantics $s(f_i)$ and $s(f_o)$ for TSGP with varying SD_t , stdGP, SLIM_GSGP over the number of generations on a subset of the analyzed datasets.

target semantic distance SD_t , which conditions the transformer model to generate offspring with a specified degree of semantic distance to the parent.

We analyze step size from two perspectives: First, we analyze how SD_t influences the actual semantic step size, and second, we study how frequently new best solutions are discovered during a run. We focus on the same datasets as before; results for the remaining datasets are provided in the Appendix D.

We measure the semantic distance $SD(f_i, f_o)$ between a parent solution f_i and its offspring f_o by using the Euclidean distance between their semantic vectors $s(f_i)$ and $s(f_o)$. Small Euclidean distances indicate a small step in semantic space and high semantic similarity; larger distances correspond to larger steps and lower similarity. We calculate the semantic vectors for a fixed, standardized input matrix of randomly sampled points (zero mean, unit variance), analogously to the approach used in *Model Building*. This allows us to study semantic similarities between datasets. Only successful variations are considered in the analysis, where the offspring differs structurally from its parent (which is almost always the case).

Figure 4 plots the median semantic distance and IQR over the number of generations for TSGP (with varying SD_t), stdGP, and SLIM_GSGP. TSGP variants are color-coded on a continuous scale from light blue (small SD_t) to black (large SD_t). We show results averaged over 30 runs. A smaller Euclidean distance corresponds to higher semantic similarity.

TSGP1 generates offspring that are semantically much closer to their parental program than those produced by stdGP or SLIM_GSGP. This is evident from the consistently lower median Euclidean distance across generations, confirming that the transformer has successfully learned to express semantic similarity through syntactic transformations and can generalize this capability to unseen datasets.

In contrast, stdGP produces large semantic modifications with high median distances and particularly high 75th percentiles. This is an expected behavior, as standard mutation and crossover operators ignore the behavioral change of programs. SLIM_GSGP maintains a relatively stable semantic similarity with low IQR across generations because its geometric semantic operators explicitly preserve the program behavior. Thus, both semantic methods, TSGP and SLIM_GSGP, exhibit more stable and controlled variations in program behavior compared to a strict syntactic search like stdGP.

Across all datasets, including those in Figure 8 in Appendix D, the influence of the target SD_t on the actual semantic distance is clearly visible. With small SD_t , TSGP produces offspring with low semantic distance (high similarity), while increasing SD_t leads to larger distances and lower

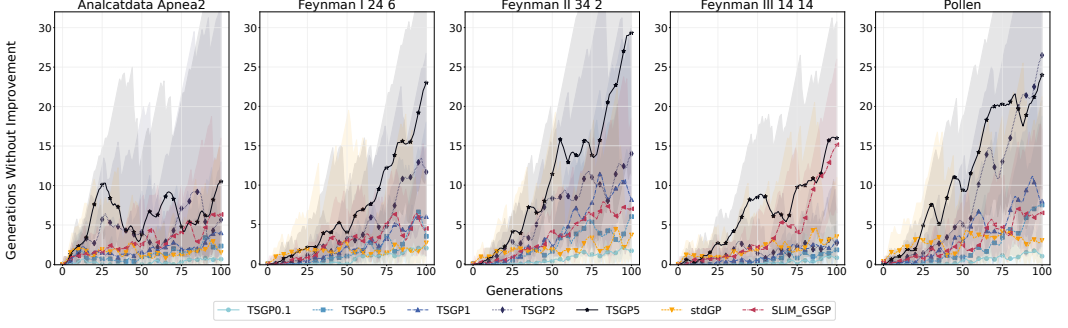


Fig. 5. Median number of generations without improving the training RMSE over the number of generations. Results are for TSGP with varying SD_t , stdGP, and SLIM_GSGP.

similarity. This shows that SD_t serves as an effective control parameter for the semantic step size. In particular, in many cases, the actual semantic distance closely matches the target SD_t , indicating that the transformer has learned to precisely and adaptively modulate the semantic similarity between parent and offspring.

Compared to [Anthes et al. 2025], we can better control the target semantic distance between offspring and parent. We believe that this is an effect of the larger training data as well as the more extensive training of the transformer, which enables better generalization and more reliable semantic variations over a variety of datasets.

For those datasets (e.g. Analcatdata Apnea2) where stdGP outperforms the semantic approaches including TSGP, stdGP performs particularly large semantic search steps. We believe that those problems are very rugged and require extensive exploration of the solution space to find high-quality solutions. Thus, the lower and much better controlled step sizes used in TSGP (and SLIM_GSGP) do not match the large step sizes necessary to solve such problems. Instead, lower step sizes lead to premature convergence. For solving such problems, a stronger randomization of search (higher exploration) using large search steps would be beneficial.

The setting of SD_t also strongly affects the frequency with which algorithms discover new candidate solutions during the search. Figure 5 plots the median number of generations without finding a new best solution over the number of generations. Lower values indicate more frequent improvements. To improve clarity, the curves are smoothed, and the y-axis range is fixed.

We find that TSGP with a small SD_t (light blue) finds new best solutions in almost every generation, as indicated by the consistently low values. This reflects strong exploitation of the candidate solutions: small semantic adjustments to the parental programs cause a large proportion of the advantageous behavior to be passed on to the offspring. However, these small changes also limit progress, as only minor fitness improvements are possible, resulting in slower convergence (see Figure 2). Furthermore, it results in larger solution sizes and increased structural complexity (see Figure 3), as high exploitation reduces the ability of TSGP to explore various structural compositions. In contrast, larger SD_t lead to more exploration. As the semantic distance between parent and offspring is larger, improvements occur less frequently, however, successful search steps often lead to much higher improvements. Especially in the first generations, strong exploration often leads to strong fitness improvements (Figure 2) and the discovery of more compact solutions (Figure 3). In summary, TSGP with $SD_t = 1$ achieves a robust balance between exploration and exploitation, leading to strong and consistent performance across diverse problems.

6 Conclusions and Future Work

This work substantially extends a previous conference paper introducing Transformer Semantic Genetic Programming (TSGP) [Anthes et al. 2025], a semantic search method that employs a pre-trained transformer model as a zero-shot variation operator to generate offspring with controlled semantic similarity to a parent program.

We find that a single transformer model can effectively learn and generalize semantic similarity across SR problems of varying input dimensionality. In a comprehensive evaluation of 24 real-world and synthetic datasets, TSGP achieves an average rank of 1.58, significantly outperforming all other baseline methods, including standard GP (rank 2.71), SLIM_GSGP (rank 2.67), Deep Symbolic Regression (rank 3.54), and DAE-GP (rank 4.50). Notably, despite its superior performance, TSGP produces solutions of similar compactness to standard GP (average rank 1.58 vs. 1.75) and significantly more compact than those of SLIM_GSGP (rank 2.67), the bloat-mitigating variant of GSGP.

Further analysis confirms that conditioning the transformer model on a target semantic distance SD_t enables precise control of the actual semantic similarity between parent and offspring. Larger SD_t values enhance exploration, leading to faster initial convergence and smaller solution structures. Smaller SD_t values promote exploitation, enabling gradual fitness improvements but often at the cost of increased program size and slower convergence. When using fixed semantic step sizes, TSGP with $SD_t = 1$ leads to a robust balance between exploitation and exploration and high performance across diverse problems.

In future work, we will analyze the syntactic mechanisms underlying TSGP’s ability to generate solutions that are both effective and compact. Furthermore, we will explore adaptive scheduling of SD_t during search and investigate how optimizing and diversifying the transformer’s training data affects solution quality and generalization.

Acknowledgments

We would like to thank David Wittenberg for the valuable insights into DAE-GP and for providing his framework. We also like to thank the entire team in Mainz for the inspiring discussions and thoughtful contributions.

References

- Philipp Anthes, Dominik Sobania, and Franz Rothlauf. 2025. Transformer Semantic Genetic Programming for Symbolic Regression. In *Proceedings of the Genetic and Evolutionary Computation Conference*. 952–960.
- Lawrence Beadle and Colin G Johnson. 2009. Semantically driven mutation in genetic programming. In *2009 IEEE Congress on Evolutionary Computation*. IEEE, 1336–1342.
- Luca Biggio, Tommaso Bendinelli, Alexander Neitz, Aurelien Lucchi, and Giambattista Parascandolo. 2021. Neural symbolic regression that scales. In *International Conference on Machine Learning*. Pmlr, 936–945.
- Martin Briesch, Dominik Sobania, and Franz Rothlauf. 2023. Large language models suffer from their own output: An analysis of the self-consuming training loop. *arXiv preprint arXiv:2311.16822* (2023).
- 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. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- Mauro Castelli, Leonardo Vanneschi, and Sara Silva. 2013. Prediction of high performance concrete strength using genetic programming with geometric semantic genetic operators. *Expert Systems with Applications* 40, 17 (2013), 6856–6862.
- François Chollet et al. 2015. Keras. <https://keras.io>.
- Grant Dick, Caitlin A Owen, and Peter A Whigham. 2020. Feature standardisation and coefficient optimisation for effective symbolic regression. In *Proceedings of the 2020 Genetic and Evolutionary Computation Conference*. 306–314.
- Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The faiss library. *arXiv preprint arXiv:2401.08281* (2024).
- Félix-Antoine Fortin, François-Michel De Rainville, Marc-André Gardner Gardner, Marc Parizeau, and Christian Gagné. 2012. DEAP: Evolutionary algorithms made easy. *The Journal of Machine Learning Research* 13, 1 (2012), 2171–2175.

- Noman Javed, Fernand Gobet, and Peter Lane. 2022. Simplification of genetic programs: a literature survey. *Data Mining and Knowledge Discovery* 36, 4 (2022), 1279–1300.
- Pierre-Alexandre Kamienny, Stéphane d’Ascoli, Guillaume Lample, and François Charton. 2022. End-to-end symbolic regression with transformers. *Advances in Neural Information Processing Systems* 35 (2022), 10269–10281.
- John R Koza. 1993. *On the programming of computers by means of natural selection*. MIT press.
- Krzysztof Krawiec and Paweł Lichocki. 2009. Approximating geometric crossover in semantic space. In *Proceedings of the 11th Annual conference on Genetic and evolutionary computation*. 987–994.
- William La Cava, Bogdan Burlacu, Marco Virgolin, Michael Kommenda, Patryk Orzechowski, Fabrício 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.
- Mikel Landajuela, Chak Shing Lee, Jiachen Yang, Ruben Glatt, Claudio P Santiago, Ignacio Aravena, Terrell Mundhenk, Garrett Mulcahy, and Brenden K Petersen. 2022. A unified framework for deep symbolic regression. *Advances in Neural Information Processing Systems* 35 (2022), 33985–33998.
- Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101* (2017).
- Sean Luke and Liviu Panait. 2002. Fighting bloat with nonparametric parsimony pressure. In *International conference on parallel problem solving from nature*. Springer, 411–421.
- Joao Francisco BS Martins, Luiz Otavio VB Oliveira, Luis F Miranda, Felipe Casadei, and Gisele L Pappa. 2018. Solving the exponential growth of symbolic regression trees in geometric semantic genetic programming. In *Proceedings of the genetic and evolutionary computation conference*. 1151–1158.
- James McDermott, Alexandros Agapitos, Anthony Brabazon, and Michael O’Neill. 2014. Geometric semantic genetic programming for financial data. In *Applications of Evolutionary Computation: 17th European Conference, EvoApplications 2014, Granada, Spain, April 23–25, 2014, Revised Selected Papers* 17. Springer, 215–226.
- Nicholas Freitag McPhee, Brian Ohs, and Tyler Hutchison. 2008. Semantic building blocks in genetic programming. In *European Conference on Genetic Programming*. Springer, 134–145.
- Aaron Meurer, Christopher P. Smith, Mateusz Paprocki, Ondřej Čertík, Sergey B. Kirpichev, Matthew Rocklin, AMiT Kumar, Sergiu Ivanov, Jason K. Moore, Sartaj Singh, Thilina Rathnayake, Sean Vig, Brian E. Granger, Richard P. Muller, Francesco Bonazzi, Harsh Gupta, Shivam Vats, Fredrik Johansson, Fabian Pedregosa, Matthew J. Curry, Andy R. Terrel, Štěpán Roučka, Ashutosh Saboo, Isuru Fernando, Sumith Kulal, Robert Cimrman, and Anthony Scopatz. 2017. SymPy: symbolic computing in Python. *PeerJ Computer Science* 3 (Jan. 2017), e103. <https://doi.org/10.7717/peerj-cs.103>
- Alberto Moraglio, Krzysztof Krawiec, and Colin G. Johnson. 2012. Geometric Semantic Genetic Programming. In *Parallel Problem Solving from Nature - PPSN XII*, Carlos A. Coello Coello, Vincenzo Cutello, Kalyanmoy Deb, Stephanie Forrest, Giuseppe Nicosia, and Mario Pavone (Eds.). Springer, Berlin, Heidelberg, 21–31. https://doi.org/10.1007/978-3-642-32937-1_3
- Caitlin A Owen, Grant Dick, and Peter A Whigham. 2018. Feature standardisation in symbolic regression. In *Australasian Joint Conference on Artificial Intelligence*. Springer, 565–576.
- Brenden K Petersen, Mikel Landajuela, T Nathan Mundhenk, Claudio P Santiago, Soo K Kim, and Joanne T Kim. 2019. Deep symbolic regression: Recovering mathematical expressions from data via risk-seeking policy gradients. *arXiv preprint arXiv:1912.04871* (2019).
- Joseph D Romano, Trang T Le, William La Cava, John T Gregg, Daniel J Goldberg, Praneel Chakraborty, Natasha L Ray, Daniel Himmelstein, Weixuan Fu, and Jason H Moore. 2021. PMLB v1.0: an open source dataset collection for benchmarking machine learning methods. *arXiv preprint arXiv:2012.00058v2* (2021).
- Franz Rothlauf. 2011. *Design of Modern Heuristics: Principles and Application* (1st ed.). Springer Publishing Company, Incorporated.
- Nguyen Quang Uy, Nguyen Xuan Hoai, and Michael O’Neill. 2009. Semantics based mutation in genetic programming: The case for real-valued symbolic regression. In *15th international conference on soft computing, Mendel*, Vol. 9. 73–91.
- Nguyen Quang Uy, Nguyen Xuan Hoai, Michael O’Neill, Robert I McKay, and Edgar Galván-López. 2011. Semantically-based crossover in genetic programming: application to real-valued symbolic regression. *Genetic Programming and Evolvable Machines* 12 (2011), 91–119.
- Leonardo Vanneschi. 2016. An introduction to geometric semantic genetic programming. In *NEO 2015: Results of the Numerical and Evolutionary Optimization Workshop NEO 2015 held at September 23–25 2015 in Tijuana, Mexico*. Springer, 3–42.
- Leonardo Vanneschi. 2024. SLIM_GSGP: The non-bloating geometric semantic genetic programming. In *European Conference on Genetic Programming (Part of EvoStar)*. Springer, 125–141.
- Leonardo Vanneschi, Mauro Castelli, Luca Manzoni, and Sara Silva. 2013. A new implementation of geometric semantic GP and its application to problems in pharmacokinetics. In *European Conference on Genetic Programming*. Springer, 205–216.
- Leonardo Vanneschi, Mauro Castelli, and Sara Silva. 2014a. A survey of semantic methods in genetic programming. *Genetic Programming and Evolvable Machines* 15 (2014), 195–214.

- Leonardo Vanneschi, Sara Silva, Mauro Castelli, and Luca Manzoni. 2014b. Geometric semantic genetic programming for real life applications. *Genetic programming theory and practice xi* (2014), 191–209.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- Henrik Voigt, Paul Kahlmeyer, Kai Lawonn, Michael Habeck, and Joachim Giesen. 2025. Analyzing Generalization in Pre-Trained Symbolic Regression. *arXiv preprint arXiv:2509.19849* (2025).
- David Wittenberg. 2022. Using denoising autoencoder genetic programming to control exploration and exploitation in search. In *European Conference on Genetic Programming (Part of EvoStar)*. Springer, 102–117.
- David Wittenberg and Franz Rothlauf. 2023. Small solutions for real-world symbolic regression using denoising autoencoder genetic programming. In *European Conference on Genetic Programming (Part of EvoStar)*. Springer, 101–116.
- David Wittenberg, Franz Rothlauf, and Christian Gagné. 2023. Denoising autoencoder genetic programming: strategies to control exploration and exploitation in search. *Genetic Programming and Evolvable Machines* 24, 2 (2023), 17.
- David Wittenberg, Franz Rothlauf, and Dirk Schweim. 2020. DAE-GP: denoising autoencoder LSTM networks as probabilistic models in estimation of distribution genetic programming. In *Proceedings of the 2020 Genetic and Evolutionary Computation Conference*. 1037–1045.

A TSGP Vocabulary

The vocabulary of the transformer model is divided into four main categories: special tokens, arithmetic operations, input variables, and constant values. Special tokens include the empty string, [UNK] (unknown token), [start], and [end], which serve as standard sequence delimiters commonly used in transformer architectures. The set of arithmetic operations consists of standard binary functions: add (addition), sub (subtraction), mul (multiplication), div (division), and pow (exponentiation). Input variables are represented as x_1 through x_5 , allowing up to five distinct features, though the actual variables used are dynamically adjusted by the model hyperparameter d , which matches the problem’s dimensionality. Constant values include both integers in the ranges $[-5, -1]$ and $[1, 5]$ in steps of 1, as well as fractional values from $[-0.5, 0.5]$ in steps of 0.1, matching the Ephemeral Random Constant (ERC) range used during the evolutionary search. Integer constants are included to maintain consistency with the training data of TSGP, where such values frequently emerge from the simplification of symbolic subexpressions using SymPy.

| Special Tokens | |
|---------------------------|--|
| "" | Empty string |
| [UNK] | Unknown token |
| [start] | Start token |
| [end] | End token |
| Arithmetic Operations | |
| add | Addition ($a + b$) |
| sub | Subtraction ($a - b$) |
| mul | Multiplication ($a \times b$) |
| div | Division (a/b) |
| pow | Power (a^b) |
| Input Variables | |
| x_1, x_2, x_3, x_4, x_5 | Available input features |
| Constant Values | |
| $[-5, -1]$ | Negative integers (step size 1) |
| $[-0.5, 0.5]$ | Floating-point constants (step size 0.1) |
| $[1, 5]$ | Positive integers (step size 1) |

B Training RMSE

Table 5. Median training RMSE of the best programs (solutions) found within 100 generations for TSGP with $SD_t = 1$, stdGP, SLIM_GSGP (SLIM), DSR, and DAE-GP (DAE) for the 24 analyzed datasets. Bold values indicate the best prediction quality (lowest RMSE). Significant differences of the best results are indicated by the label symbols.

| Data set | <i>a</i> TSGP1 | <i>b</i> stdGP | <i>c</i> SLIM | <i>d</i> DSR | <i>e</i> DAE |
|---------------------------|-----------------------------|---------------------------|-------------------------|--------------|--------------|
| Analcatdata Apnea1 | 0.9359 | <i>d</i> 0.6768 | 0.8788 | 0.9316 | 0.9915 |
| Analcatdata Apnea2 | 0.9298 | <i>abce</i> 0.6848 | 0.8886 | 0.9416 | 0.9915 |
| Analcatdata Neavote | <i>bde</i> 0.2122 | 0.2277 | 0.2155 | 0.2512 | 0.2561 |
| Chscase Geyser1 | <i>de</i> 0.4530 | 0.4657 | 0.4544 | 0.4888 | 0.4859 |
| Cloud | <i>bcd</i> 0.2835 | 0.4099 | 0.3431 | 0.4649 | 0.5166 |
| ERA | <i>bcd</i> 0.7775 | 0.8077 | 0.7856 | 0.8349 | 0.8901 |
| ESL | <i>bcd</i> 0.3567 | 0.5324 | 0.3846 | 0.4694 | 0.5720 |
| Feynman I 18 4 | <i>bcd</i> 0.1630 | 0.3145 | 0.2024 | 0.4371 | 0.5300 |
| Feynman I 24 6 | <i>bcd</i> 0.1616 | 0.3768 | 0.2825 | 0.6633 | 0.7275 |
| Feynman I 25 13 | 0.1934 | 0.2132 | <i>de</i> 0.1817 | 0.4831 | 0.5175 |
| Feynman I 29 4 | <i>cde</i> 0.2815 | 0.3224 | 0.3383 | 0.5796 | 0.5845 |
| Feynman I 43 43 | <i>bcd</i> 0.3809 | 0.5353 | 0.4260 | 0.7755 | 0.8029 |
| Feynman II 34 2 | <i>bcd</i> 0.0876 | 0.2246 | 0.1957 | 0.6490 | 0.6767 |
| Feynman II 38 3 | <i>bcd</i> 0.3154 | 0.5196 | 0.4030 | 0.7690 | 0.8307 |
| Feynman II 4 23 | <i>bcd</i> 0.2879 | 0.3978 | 0.3880 | 0.6999 | 0.7417 |
| Feynman III 14 14 | <i>bcd</i> 0.4421 | 0.5971 | 0.5346 | 0.8361 | 0.8752 |
| Feynman Test 3 | <i>bcd</i> 0.1718 | 0.2328 | 0.2301 | 0.4769 | 0.5491 |
| Feynman Test 4 | <i>cde</i> 0.2237 | 0.2797 | 0.2479 | 0.4717 | 0.5803 |
| Galaxy | 0.2739 | 0.3055 | <i>de</i> 0.2641 | 0.3979 | 0.4580 |
| LEV | <i>bcd</i> 0.6536 | 0.7196 | 0.6669 | 0.7564 | 0.8222 |
| Pollen | <i>cde</i> 0.4624 | 0.5102 | 0.5121 | 0.7170 | 0.7836 |
| Rabe | <i>bde</i> 0.0934 | 0.1814 | 0.0993 | 0.1891 | 0.3448 |
| Vinnie | <i>bde</i> 0.4889 | 0.4939 | 0.5018 | 0.4977 | 0.5042 |
| Visualizing Environmental | <i>de</i> 0.7364 | 0.7641 | 0.7423 | 0.8241 | 0.8193 |
| Rank (Mean & Std) | <i>bcd</i> 1.29±0.73 | 2.75±0.66 | 2.08±0.57 | 3.96±0.45 | 4.92±0.28 |

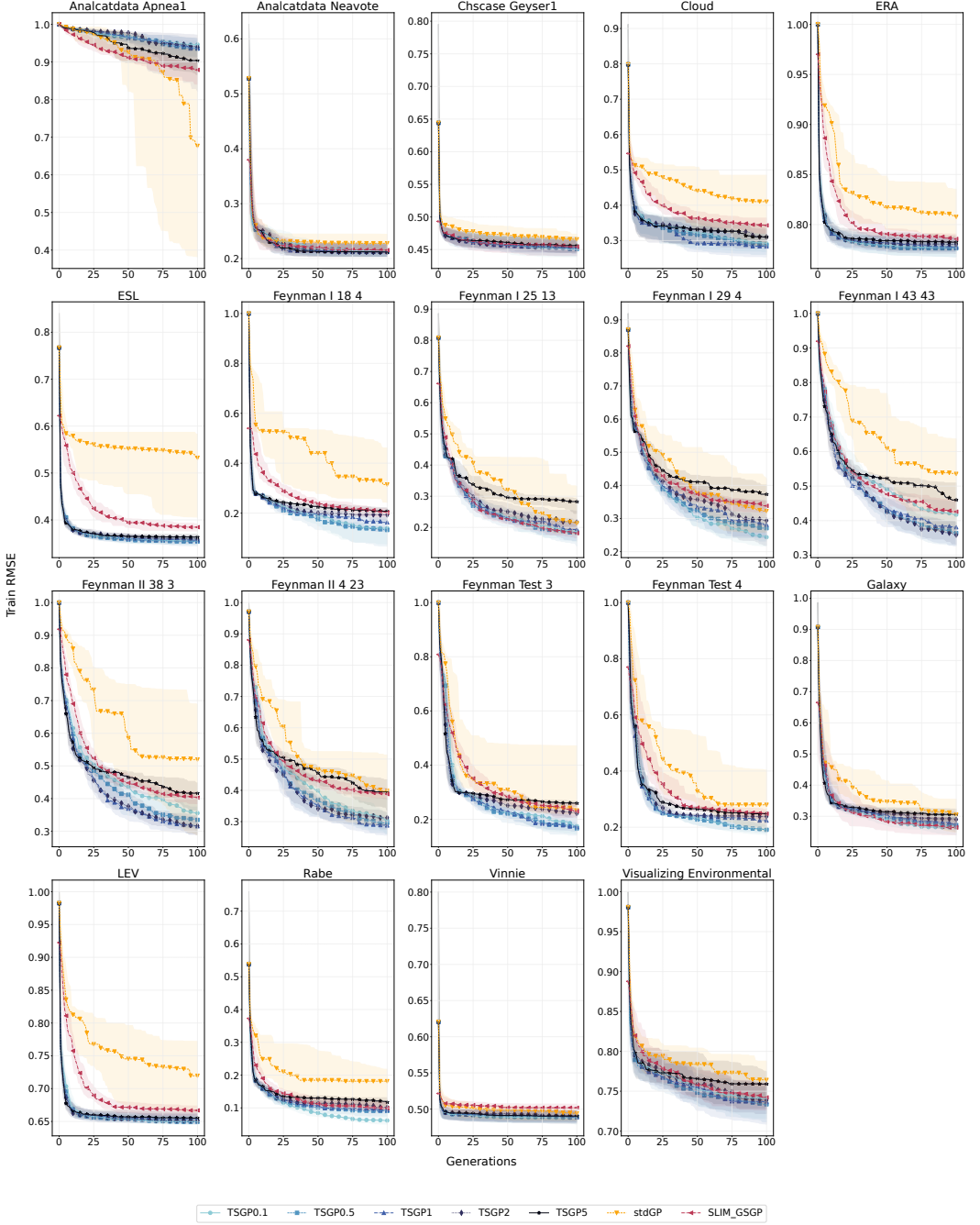


Fig. 6. Median training RMSE of the solutions of TSGP, stdGP, SLIM_GSGP over the number of generations on the remaining analyzed datasets

C Solution Size

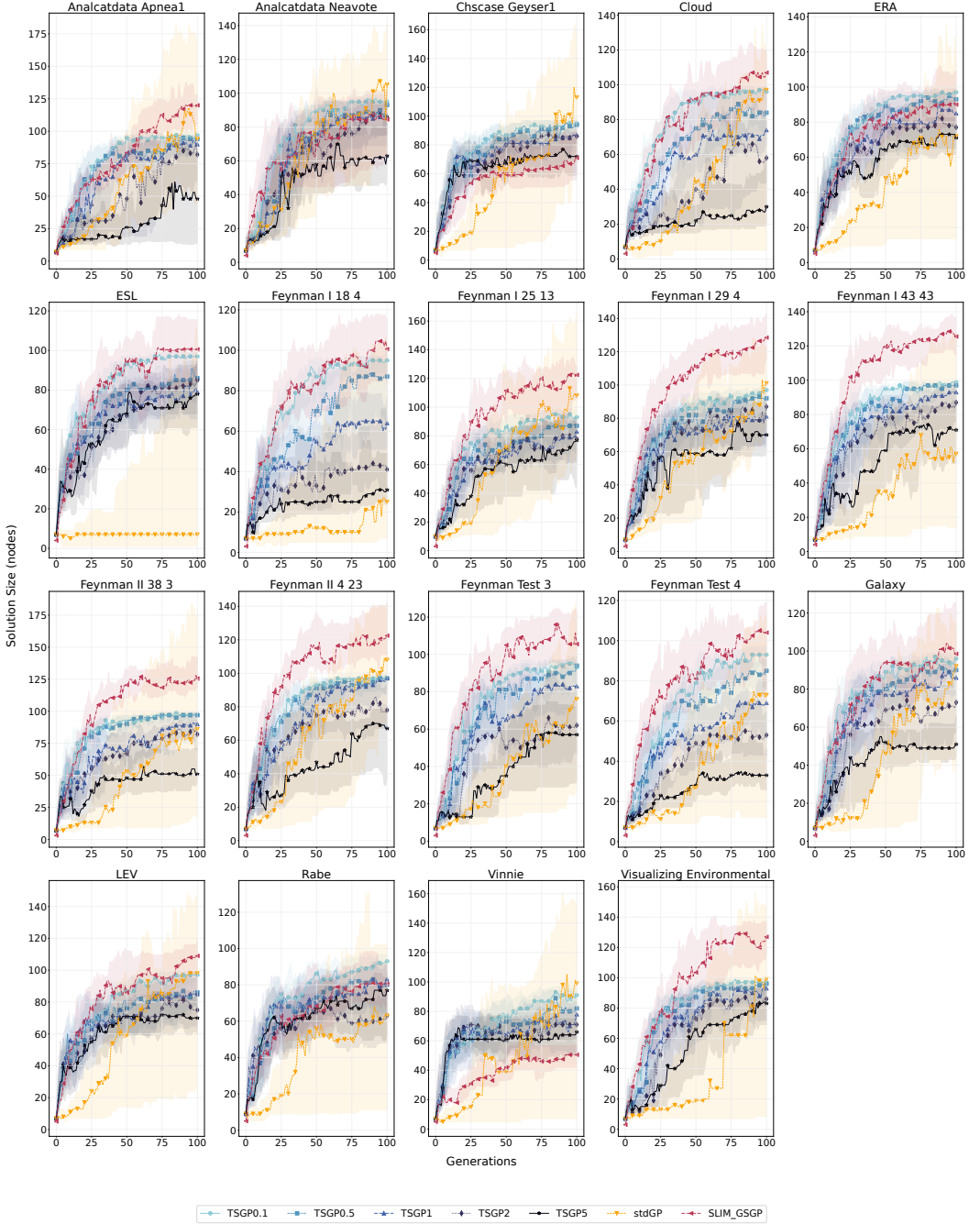


Fig. 7. Median solution size of TSGP, stdGP, SLIM_GSGP over the number of generations on the remaining analyzed datasets.

D Step Size

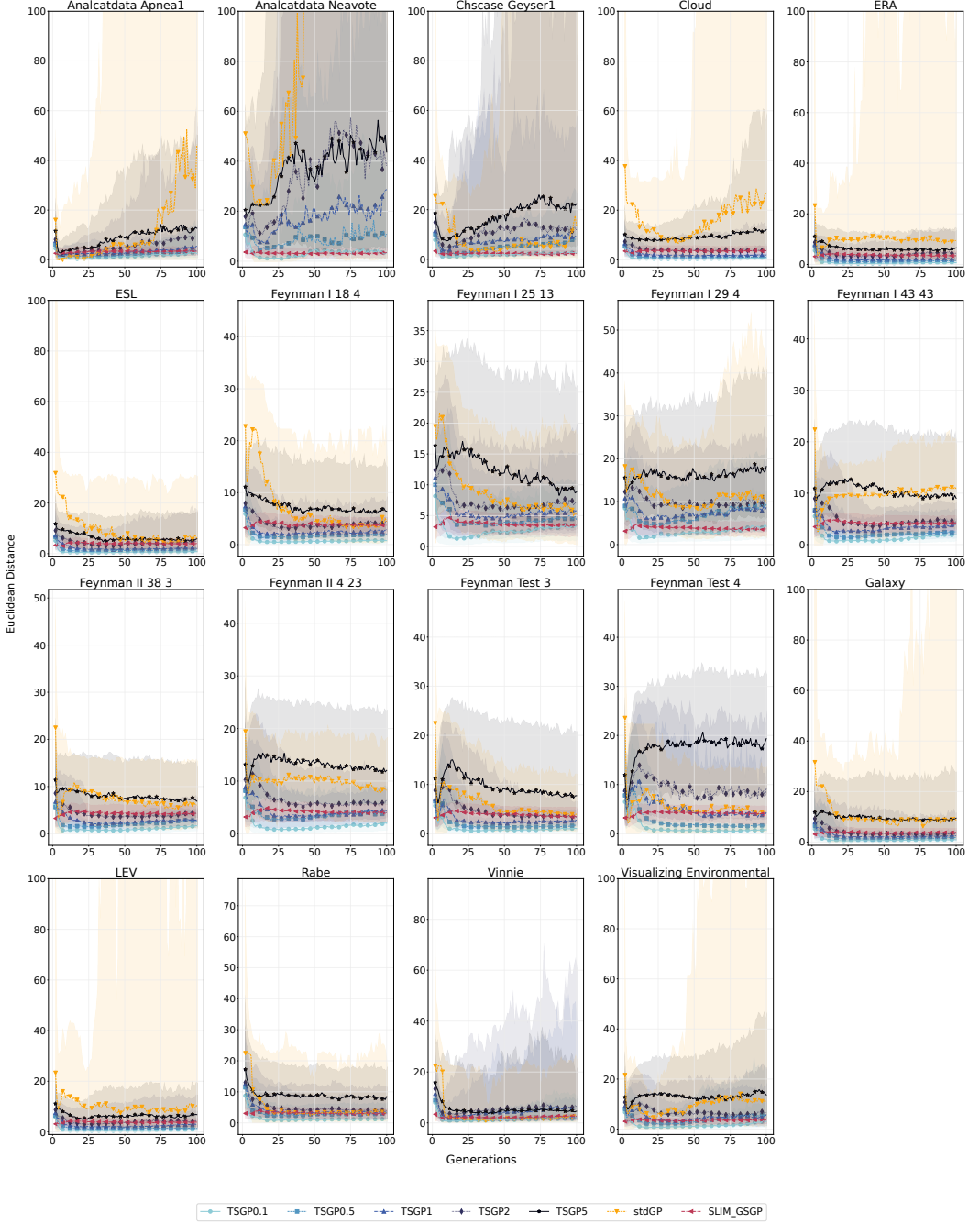


Fig. 8. Median Euclidean distance between the semantics $s(f_i)$ and $s(f_o)$ for TSGP with varying SD_t , stdGP, SLIM_GSGP over the number of generations on the remaining analyzed datasets.

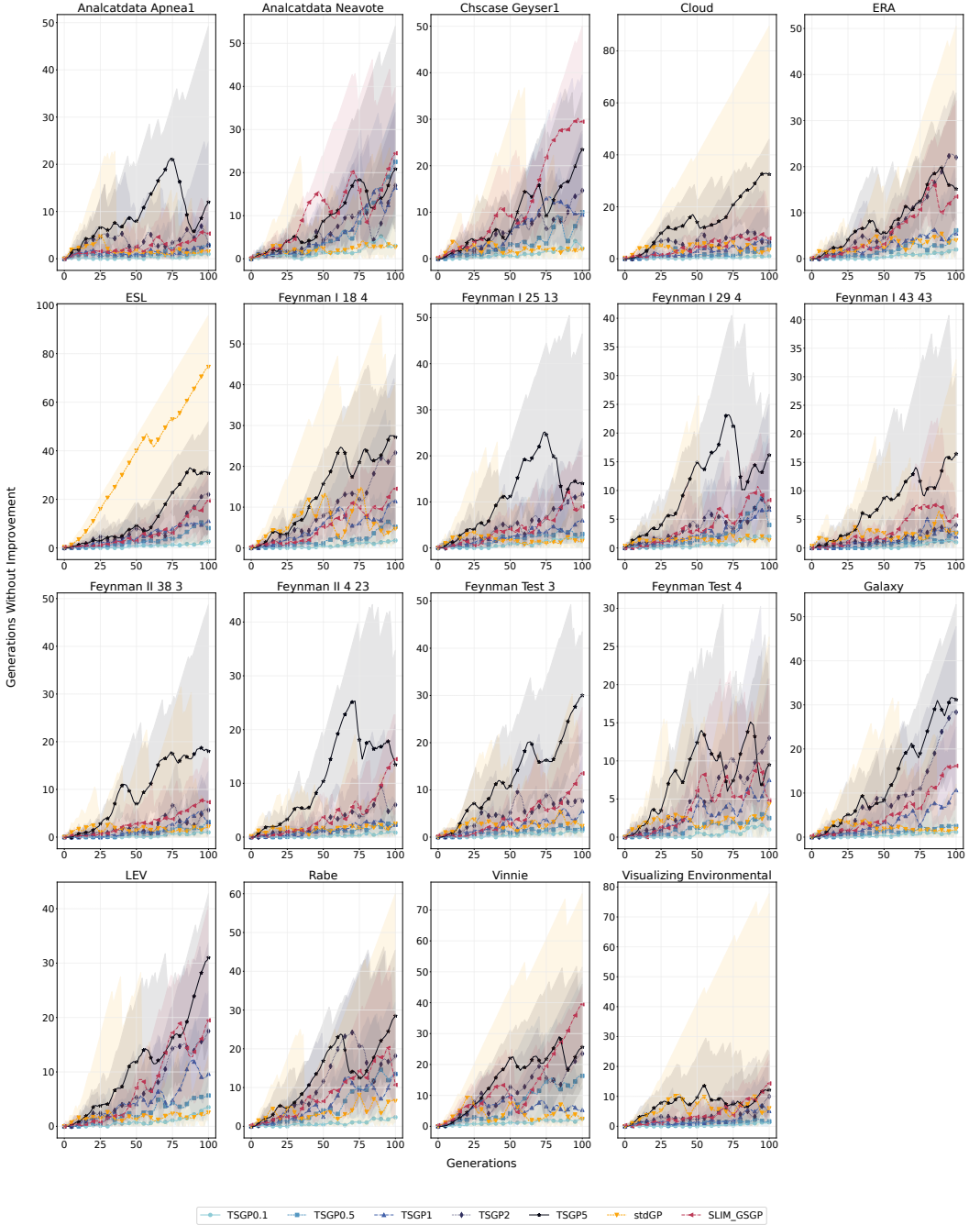


Fig. 9. Median number of generations without improving the training RMSE over the number of generations. Results are for TSGP with varying SD_t , stdGP, and SLIM_GSGP.