

On the Effectiveness of Genetic Operations in Symbolic Regression

Bogdan Burlacu^{1,2}, Michael Affenzeller^{1,2}, and Michael Kommenda^{1,2}

¹ Heuristic and Evolutionary Algorithms Laboratory
School of Informatics, Communications and Media
University of Applied Sciences Upper Austria
Softwarepark 11, 4232 Hagenberg, Austria

² Institute for Formal Models and Verification
Johannes Kepler University Linz
Altenbergerstr. 69, 4040 Linz, Austria

{bogdan.burlacu,michael.affenzeller,michael.kommenda}@fh-hagenberg.at

*

Abstract. This paper describes a methodology for analyzing the evolutionary dynamics of genetic programming (GP) using genealogical information, diversity measures and information about the fitness variation from parent to offspring. We introduce a new subtree tracing approach for identifying the origins of genes in the structure of individuals, and we show that only a small fraction of ancestor individuals are responsible for the evolution of the best solutions in the population.

Keywords: Genetic programming, evolutionary dynamics, algorithm analysis, symbolic regression

1 Introduction

Empirical analysis in the context of different benchmark problems and tentative algorithmic improvements (such as various selection schemes or fitness assignment techniques) has a limited ability of explaining genetic programming (GP) behavior and dynamics. Results usually confirm our intuitions about the relationship between selection pressure, diversity, fitness landscapes and genetic operators, but they prove difficult to extend to more general theories about the internal functioning of GP.

This work is motivated by the necessity for a different approach to study the GP evolutionary process. Instead of looking for correlations between different selection or fitness assignment mechanisms and solution quality or diversity, we focus on the reproduction process itself and the effectiveness of the variation-producing operators in transferring genetic material.

* The final publication is available at
https://link.springer.com/chapter/10.1007/978-3-319-27340-2_46

Achieving good solutions depends on the efficient use of the available gene pool given its inherent stochasticity (random initialization, random crossover, random mutation). Under the effects of selection pressure, many suboptimal exchanges of genetic information will cause a decrease in the amount of genetic material available to the evolutionary engine. Measures to mitigate this phenomenon usually use various heuristics for guiding either selection or the crossover operator towards more promising regions of the search space [1,2].

Diversity is an important aspect of GP, considered to be a key factor in its performance. Multiple studies dedicated to GP diversity analyze diversity measures (based on various distance metrics, for example [3]) in correlation with the effects of genetic operators [4,5,6]. Genotype operations – crossover in particular – often have a negative (or at most, neutral) effect on individuals, leading to diversity loss in the population following each selection step. This effect is due to the interplay between crossover and selection which leads to an increase in average program size [7] (when sampling larger programs, crossover has a higher chance of having a neutral effect).

2 Methodology

In this paper we introduce a new methodology for the exact identification (“tracing”) of any structural change an individual is subjected to during evolution. We use this methodology in combination with population diversity and genealogy analysis methods to investigate the effects of the genetic operators in terms of how often they lead to a fitness improvement, how often they overlap (for example when the same area inside the tree is repeatedly targeted by crossover), and how often they cancel each other out.

2.1 Tracing of genotype fragments

This method is based on previous work on population genealogies [8,9]. During the algorithm run, every new generation is added to the genealogy graph with arcs connecting child vertices to their parents. When crossover is followed by mutation, both the results of crossover and mutation are saved in the graph (Figure 1).

We define an individual’s *trace graph* as a collection of vertices representing its ancestors from which the various parts of its genotype originated, connected by a collection of arcs representing the different genotype operations that gradually assembled those parts.

The tracing procedure uses a set of simple arithmetic rules to navigate genealogies and identify the relevant subtrees, based on the indices of the subtree to be traced and the index of the received fragment (Figure 2). The nodes in each tree are numbered according to their preorder index i such that, given two

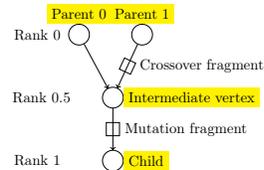


Fig. 1: Saving intermediate results in the genealogy graph

subtrees A and B , $B \subset A$ if $i_A < i_B < i_A + l_A$, where i_A, i_B are their respective preorder indices and l_A, l_B are their lengths.

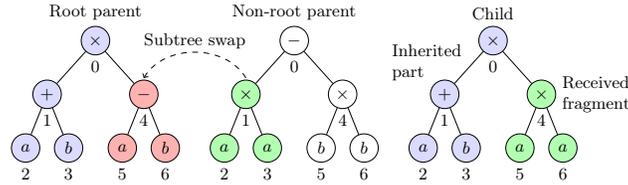


Fig. 2: Preorder arithmetics for subtree inclusion

Since some individuals within the ancestry of the traced individual may have contributed parts of their genotype to multiple offspring, there may exist multiple evolutionary trajectories in the trace graph passing through the same vertex or sequence of vertices, reflected in the graph by multiple arcs between the same two vertices, each arc representing the transmission of different genes or building blocks.

2.2 Analysis of population dynamics

The various measurements used to quantify the behavioral aspects of GP are described in more detail within the following paragraphs.

Genetic Operator Effectiveness

Operator effectiveness is calculated as the difference in fitness between the child and its root parent.

Average fitness improvement Let N be the total number of individuals in the population, t_i one individual and p_i its parent:

$$\bar{q} = \frac{1}{N} \cdot \sum_{i=1}^N (Fitness(t_i) - Fitness(p_i))$$

Best fitness improvement Return the difference between the fitness values of the best individual t_{best} and its parent p_{best}

$$q_{best} = Fitness(t_{best}) - Fitness(p_{best})$$

The average and best fitness improvements are calculated individually for crossover and mutation operations.

Average relative overlap

We define the relative overlap between two sets A_1 and A_2 using the Sørensen-

Dice coefficient³ which can also be seen as a similarity measure between sets:

$$s(A_1, A_2) = \frac{2 \cdot |A_1 \cap A_2|}{|A_1| + |A_2|}$$

The reason for using this measure is to see how much overlap exists between the trace graphs and root lineages of the individuals in the population. A high relative overlap would mean that diversity is exhausted as all the individuals have the same parents or ancestors.

Genotype and phenotype similarity

These similarity measures provide information about the evolution of diversity from both a structural (genotype) and a semantic (phenotype) perspective. Genotype similarity is calculated using a *bottom-up tree mapping* [10] that can be computed in time linear in the size of the trees and has the advantage that it works equally well for unordered trees. For two trees T_1 and T_2 and a bottom-up mapping M between them, the similarity is given by:

$$GenotypeSimilarity(T_1, T_2) = \frac{2 \cdot |M|}{|T_1| + |T_2|}$$

Phenotype similarity between two trees is calculated as the Pearson R^2 correlation coefficient between their respective output values on the training data.

Contribution ratio

While it is clear that under the influence of random evolutionary forces (such as genetic drift or hitchhiking) each of an individual's ancestors plays an equally important role in the events leading to its creation, the trace graph represents a powerful tool for analyzing the origin of genes and the way solutions are assembled by the genetic algorithm.

The size of the trace graph relative to the size of the complete ancestry can be used as a measure of the effort spent by the algorithm to achieve useful adaptation. For example, a small trace graph means that a small number of an individual's ancestors contributed to the assembly of its genotype, via an equally small number of genetic operations (crossover and mutation). The effort, seen as the ratio of effective genetic operations over the total number of genetic operations, can give an indication of how easy new and better solutions can be assembled by the algorithm.

The contribution ratio r is given by the percentage of individuals from the best solution ancestry that had an actual contribution to its structure:

$$r = \frac{|Trace(bestSolution)|}{|Ancestry(bestSolution)|}$$

³ It was also possible to use the *Jaccard index* $J(A_1, A_2) = \frac{|A_1 \cap A_2|}{|A_1 \cup A_2|}$ as it is very similar to the Sørensen-Dice coefficient. However this choice makes no practical difference for the results presented in this publication

3 Experiments

For the experimental part, we use GP to solve two symbolic regression benchmark problems:

Vladislavleva-8

$$F_8(x_1, x_2) = \frac{(x_1 - 3)^4 + (x_2 - 3)^3 - (x_2 - 3)}{(x_2 - 2)^4 + 10}$$

Poly-10

$$F(\mathbf{x}) = x_1x_2 + x_3x_4 + x_5x_6 + x_1x_7x_9 + x_3x_6x_{10}$$

The Vladislavleva-8 problem was solved using the standard GP algorithm (SGP) with a population size of 500 individuals and 50 generations (in order to be able to compute the trace graphs of each individual in the population in feasible time). For the Poly-10 problem the offspring selection GP (OSGP) [11] was also tested with a population size of 300 individuals and gender-specific selection.

We analyzed the algorithm dynamics using the genealogy graph, the ancestry of the best solution and the trace history of its genotype. Other additional measurements such as diversity, size and quality distributions were included for a more complete picture. All the results were averaged on a collection of 20 algorithmic runs for each problem and algorithm configuration.

In the case of SGP, we see in Figure 3 that the genetic operators produce negative improvement on average, meaning that in most cases the fitness of the child is worse than the fitness of the parent. The light-colored curves filled with green in Figure 3 represent the best improvement while the dark-colored once filled with red represent the average improvement. As average fitness improvement produced by genetic operators tends to be negative, the increase in average population fitness can be attributed to the interplay between recombination operators and selection. OSGP operator improvement is always small but positive due to the requirement that offspring are better than their parents.

The ability to produce useful genetic variation (leading to fitness improvements) is directly related to the structural diversity of the population which cannot be controlled through fitness-based selection. Results in Figures 4a and 4b reveal the relationship mediated by the selection mechanism between the structural similarity between two trees and the degree to which their root lineages and their trace graphs overlap. The high correlation (calculated as the Pearson R^2 coefficient) between the three curves corresponds intuitively to the fact that similar individuals come from similar (partially overlapping) lineages, with the important difference that trace graphs do not represent lineages in the strictest sense, as they only include those ancestors whose genes survived in the structure of the traced individual. In Figures 4c and 4d we show the correlation between semantic similarity and quality of the best solution. We see that SGP does not suffer from loss of semantic diversity. With offspring selection, as children are required to outperform their parents, the semantic similarity increases rapidly to a value close to 1.

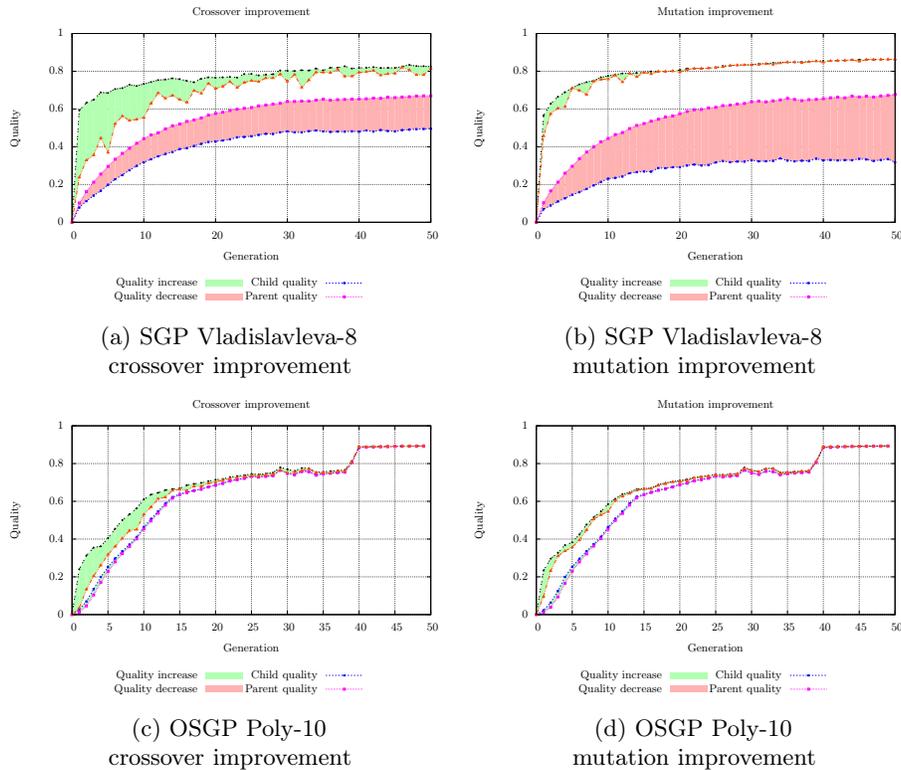


Fig. 3: SGP Vladislavleva-8 and Poly-10 best (above) and average (below) operator improvement

Another aspect of GP search is illustrated in Figure 5a, where we can observe the exploratory behavior of the OSGP algorithm in the beginning of the run, when the building blocks representing the terms of the formula are gradually discovered, and the exploitative behavior towards the end, when no big jumps in quality are produced, but the solution is incrementally improved through small changes of the tree constants and variable weighting factors.

Finally, the contribution ratio for SGP and OSGP was calculated at 13% and 4%, respectively, showing a high degree of interrelatedness between individuals which leads to low genetic operator efficiency. Fit individuals contribute multiple times, but selection pressure exceeds their variability potential. Offspring selection improves efficiency by adapting selection pressure.

4 Conclusion and outlook

Our results show that in most cases GP operators do not lead to fitness improvement. The tracing of the best solution indicates that a few critical operations

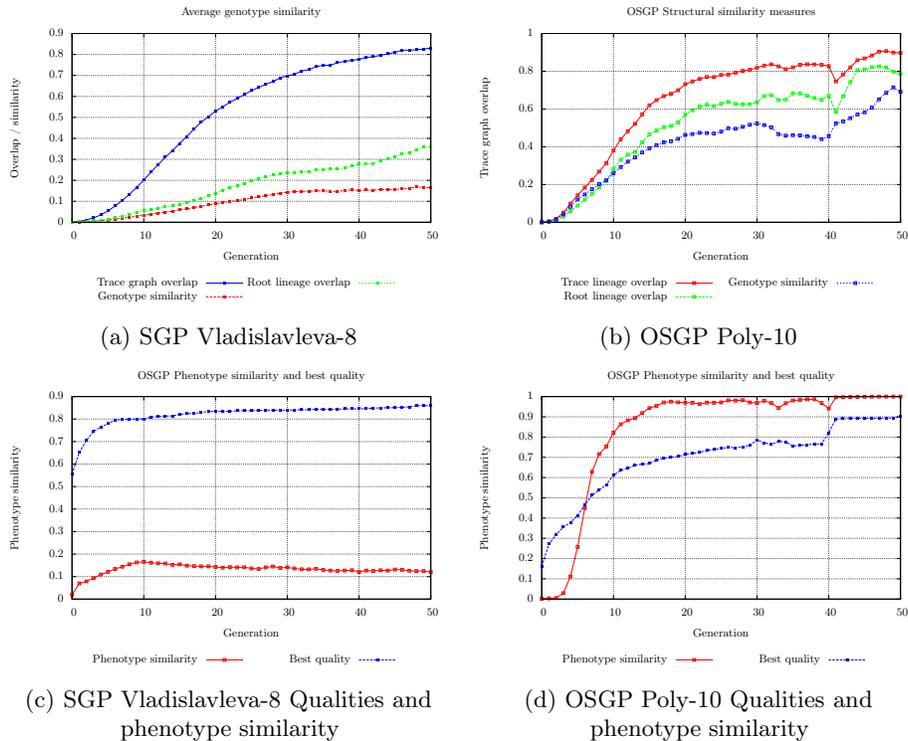


Fig. 4: Relationship between root lineage/trace graph overlap and genotype similarity

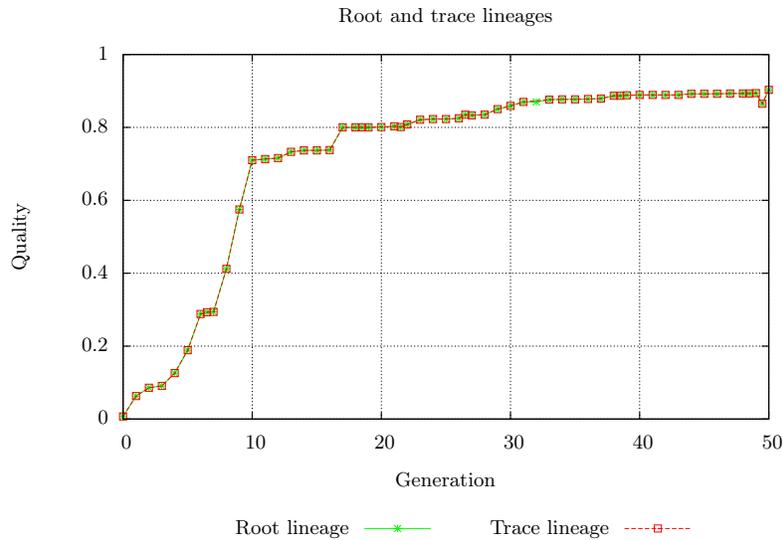
when the algorithm is able to assemble high fitness solution elements out of pre-existing, disparate genes are responsible for the performance of the entire run. A significantly small fraction (around 13% for SGP and 4% for OSGP) of all ancestors of the best individual have an actual contribution to its final structure.

The tracing methodology can reveal interesting and previously unexplored aspects of GP evolution regarding genetic operators and their effects on population dynamics. In contrast to other methods and techniques, our approach provides a more accurate and complete description of the evolutionary process.

Acknowledgments The work described in this paper was done within the COMET Project Heuristic Optimization in Production and Logistics (HOPL), #843532 funded by the Austrian Research Promotion Agency (FFG).

References

1. Burke, E.K., Gustafson, S., Kendall, G., Krasnogor, N.: Is increased diversity in genetic programming beneficial? an analysis of the effects on performance. In



(a) OSGP Poly-10 best solution (the term x_3x_4 was already present in the initial formula)

- Sarker, R., Reynolds, R., Abbass, H., Tan, K.C., McKay, B., Essam, D., Gedeon, T., eds.: Proceedings of the 2003 Congress on Evolutionary Computation CEC2003, Canberra, IEEE Press (2003) 1398–1405
2. Burke, E.K., Gustafson, S., Kendall, G.: Diversity in genetic programming: An analysis of measures and correlation with fitness. *IEEE Transactions on Evolutionary Computation* **8** (2004) 47–62
 3. Mattiussi, C., Waibel, M., Floreano, D.: Measures of diversity for populations and distances between individuals with highly reorganizable genomes. *Evolutionary Computation* **12** (2004) 495–515
 4. Ekárt, A., Németh, S.Z.: Maintaining the diversity of genetic programs. In Foster, J.A., Lutton, E., Miller, J., Ryan, C., Tettamanzi, A.G.B., eds.: *Genetic Programming, Proceedings of the 5th European Conference, EuroGP 2002*. Volume 2278 of LNCS., Kinsale, Ireland, Springer-Verlag (2002) 162–171
 5. Nguyen, T.H., Nguyen, X.H.: A brief overview of population diversity measures in genetic programming. In Pham, T.L., Le, H.K., Nguyen, X.H., eds.: *Proceedings of the Third Asian-Pacific workshop on Genetic Programming, Military Technical Academy, Hanoi, VietNam (2006)* 128–139
 6. Jackson, D.: Phenotypic diversity in initial genetic programming populations. In Esparcia-Alcazar, A.I., Ekart, A., Silva, S., Dignum, S., Uyar, A.S., eds.: *Proceedings of the 13th European Conference on Genetic Programming, EuroGP 2010*. Volume 6021 of LNCS., Istanbul, Springer (2010) 98–109
 7. Dignum, S., Poli, R.: Crossover, sampling, bloat and the harmful effects of size limits. In O’Neill, M., Vanneschi, L., Gustafson, S., Esparcia Alcazar, A.I., De Falco, I., Della Cioppa, A., Tarantino, E., eds.: *Proceedings of the 11th European Conference on Genetic Programming, EuroGP 2008*. Volume 4971 of Lecture Notes in Computer Science., Naples, Springer (2008) 158–169

8. Burlacu, B., Affenzeller, M., Kommenda, M., Winkler, S.M., Kronberger, G.: Evolution tracking in genetic programming. In Jimenez, E., Sokolov, B., eds.: The 24th European Modeling and Simulation Symposium, EMSS 2012, Vienna, Austria (2012)
9. Burlacu, B., Affenzeller, M., Kommenda, M., Winkler, S., Kronberger, G.: Visualization of genetic lineages and inheritance information in genetic programming. In: GECCO '13 Companion: Proceeding of the fifteenth annual conference companion on Genetic and evolutionary computation conference companion, Amsterdam, The Netherlands, ACM (2013) 1351–1358
10. Valiente, G.: An efficient bottom-up distance between trees. In: Proceedings of the 8th International Symposium of String Processing and Information Retrieval, Press (2001) 212–219
11. Affenzeller, M., Winkler, S., Wagner, S., Beham, A.: Genetic Algorithms and Genetic Programming: Modern Concepts and Practical Applications. Numerical Insights. CRC Press, Singapore (2009)