Genetic Programming Exploratory Power and the Discovery of Functions*

Justinian P. Rosca

Abstract

Hierarchical genetic programming (HGP) approaches rely on the discovery, modification, and use of new functions to accelerate evolution. This paper provides a qualitative explanation of the improved behavior of HGP, based on an analysis of the evolution process from the dual perspective of diversity and causality. From a static point of view, the use of an HGP approach enables the manipulation of a population of higher diversity programs. Higher diversity increases the exploratory ability of the genetic search process, as demonstrated by theoretical and experimental fitness distributions and expanded structural complexity of individuals. From a dynamic point of view, an analysis of the causality of the crossover operator suggests that HGP discovers and exploits useful structures in a bottom-up, hierarchical manner. Diversity and causality are complementary, affecting exploration and exploitation in genetic search. Unlike other machine learning techniques that need extra machinery to control the tradeoff between them, HGP automatically trades off exploration and exploitation.

1 INTRODUCTION

The problem of understanding and controlling the mechanism of genetic programming (GP) is challenging especially in the case of GP extensions for the discovery and evolution of functions. Such GP extensions have been designed with the goal of automating the discovery of functions that are beneficial during the search for solutions by exploiting opportunities to parameterize and reuse code. Two such techniques are *automatic definition* of functions (ADF) (Koza 1992) and adaptive representation (AR) (Rosca and Ballard 1994a). The former is a GP extension that allows the evolution of reusable subroutines. The latter is based on the discovery of useful building blocks of code. These blocks are identified by analyzing the trajectory of evolution. Blocks are generalized and transformed into new functions which extend the function set in an adaptive manner. As opposed to ADFs,

^{*}To appear in the Proceedings of the The Fourth Annual Conference on Evolutionary Programming, San Diego 1995

new AR functions rely on heuristically selected code fragments and do not evolve. Although ADF and AR approaches implement the ideas of discovery, modification, and use of new functions in different ways, both actually evolve a hierarchy of functions that greatly improve search efficiency. This paper refers to both mechanisms by hierarchical genetic programming (HGP).

No clear mathematical analysis currently exists for how either GP or HGP sample the solution space. The goal of this paper is to analyze the influence of different representational choices on the behavior of GP. This paper analyses several explanations for the improved behavior of HGP due to function discovery and proposes a *bottom-up* HGP evolution scenario: HGP discovers and exploits useful structures in a bottom-up, hierarchical manner.

Two complementary dimensions of genetic search are discussed in the paper: diversity of programs and GP causality. Discovery and use of encapsulated subroutines causes increased population diversity. Experimental evidence outlining increased program size and varied program shape is presented to explain this increased population diversity. The paper compares theoretical and practical distributions of fitness for randomly generated solutions in a test problem characterized by a finite function sample space.

The principle of *strong causality* states that small alterations in the underlying structure of an object, or small departures from the cause determine small changes of the object's behavior, or small changes of the effects, respectively (Rechenberg 1994), (Lohmann 1992). In GP small alterations of the programs may generate big changes in behavior. From this perspective GP is weakly causal. The trend of structures called *birth certificates* are presented as evidence for the way HGP inherits useful structures. Birth certificates represent types of crossover in the genealogic tree of a solution and record the evolution trajectory of that solution.

The paper outline is as follows. The next section defines the underlying principle of HGP, the resulting change in representation, and briefly presents the two HGP approaches used throughout the experiments and other related work. Section 3 introduces a test case and presents a theoretical analysis of fitness distributions for a uniform probability distribution of solutions. It analyzes the random generation of program trees and compares theoretical distributions of partial solutions with those actually obtained in GP for a varying function set. Changes in representation determine changes in the size, shape, and behavior of program trees. Section 4 presents an analysis of the ADF evolution dynamics. The evolution of functions offers a means for expressing, combining, and propagating useful building blocks and it also contributes in an essential way to the exploratory ability of GP. Discovered functions represent an adaptive control mechanism in the exploration-exploitation tradeoff. In conclusion the paper discusses the results and suggests future research.

2 HIERARCHICAL GENETIC PROGRAMMING

Genetic programming departs from the genetic algorithm (GA) paradigm by using trees to represent genotypes (Cramer 1985), (Koza 1992). Trees provide a flexible representation for creating and manipulating programs. This paper uses the denotations *tree* and *subtree* to refer to the parse tree of a program or a part of it respectively.

Problem representation in GP is defined by a set of problem-dependent primitive functions. Functions of one of more variables label internal nodes of the tree while functions of no arguments, called terminals, label leaves of the tree. The search space for GP is the space of all programs that can be built using these initial primitives. The intuition for HGP systems is that adapting the composition of these sets dramatically changes the behavior of GP. For example, the inclusion of more complex functions, known to be part of a final solution, will result in less computational effort spent during search and thus will enable a shorter time to finding a final solution.

The HGP approaches presented below, automatic definition of functions (ADF-GP) (Koza 1992) and adaptive representation (AR-GP) (Rosca and Ballard 1994a) use the above observation in different ways in order to accelerate search.

Automatic Definition of Functions

The automatic definition of functions approach (ADF-GP) assumes that parsimonious problem solutions can be specified in terms of a main program and a hierarchical collection of subroutines. The main program invokes a subset of the subroutines to perform the overall computation, while those subroutines may in turn call other subroutines computing partial results.

Genetic programming is used both to search for appropriate subroutines, and to find a way of composing discovered subroutines and primitive functions into a complete solution. In this approach, apparently, GP has to perform a more difficult search task. The problem becomes well defined if the functions and terminals that can be invoked by each subroutine and by the main program are completely specified. During evolution, only the fitness of the complete program is evaluated.

In this approach each individual program has a dual structure. The structure is defined based on a fixed number of components or *branches* to be evolved: several function branches and a main program branch. Each function branch (for instance ADF_0 , ADF_1) has a fixed number of arguments. The main program branch (Program-Body) produces the result. Each branch is a piece of LISP code built out of specific primitive terminal and function sets, and is subject to genetic operations. The set of function-defining branches, the number of arguments that each of the function possesses and the "alphabet" (function and terminal sets) of each branch define the *architecture* of a program. The references allowed between function branches determine a hierarchical organization of the set of functions. Although the number and interconnectivity of ADFs are fixed, the definition of ADFs evolve. Genetic operations on ADFs are syntactically constrained by the components on which they can operate. For example, crossover can only be performed between subtrees of the same type, where subtree type

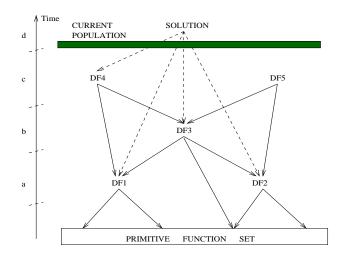


Figure 1: A hypothetical call graph of the extended function set in the AR method. The primitive function set is extended hierarchically with functions (DF1, DF2, etc.) discovered at generation numbers a,b,c.

depends on the function and terminal symbols used in the definition of that subtree. An example of a simple typing rule for an architecturally uniform population of programs is *branch typing*. Each branch of a program is designated as having a distinct type. In this case the crossover operator can only swap subtrees from analogous branches.

Adaptive Representation

In contrast to ADF's passive function definition, AR explicitly attempts to discover and use new functions. A hierarchy of automatic functions is created in a bottom-up fashion as the problem is being solved (see figure 1).

At the base of the function hierarchy lie the primitive functions from the initial function set. More complex functions are dynamically built on the primitive functions, and become stable components of the representation. The levels in the hierarchy are discovered by using heuristic information as conveyed by the environment and captured using heuristic block fitness functions, as a result of several steps:

- 1. Select blocks appearing in (good) individuals
- 2. Generalize candidate blocks to create new functions
- 3. Extend the representation with the new functions, noticing if progress is made.

The generation intervals with no function set changes represent evolutionary *epochs*. In order to make use of the newly discovered functions, at the beginning of each epoch, part of the population is replaced with random individuals built using the extended function set.

The discovery of functions in AR can be guided either by statistical information, like the fitness distribution of a newly generated subpopulation,

or by measures of the short-term influence of the new functions on the evolutionary process.

Other related work

Modularization is an approach which addresses the problems of inefficiency and scaling in GP. This issue has generated research efforts towards defining the notion of building block in GP and finding useful ways to manipulate modules of code.

A GP analogy along the lines of GA schemata theory and GA building block hypothesis has been attempted in (O'Reilly and Oppacher 1994). The main goal was understanding if GP problems have building block structure and when GP is superior to other search techniques. The approach was to generalize the definition of a GP schema from (Koza 1992) to a collection of tree fragments, that is a collection of trees possibly having subtrees removed. An individual instantiates a schema in case it "covers" (matches) all the schema fragments, overlappings between fragments not being allowed. The probability of disruption by crossover is estimated based on these definitions. The authors concluded that schema analysis is difficult and does not offer an appropriate perspective for analyzing GP.

A GP structural theory analogous to GA schemata theory fundamentally ignores the functional role of the GP representation. The analysis of building blocks in AR-GP (Rosca and Ballard 1994a) starts from this hypothesis and takes a functional approach. The ADF approach, presented earlier, is also a method of representing and using modularity in GP. Another method, *module acquisition* (Angeline 1994), (Angeline and Pollack 1994) introduced many inspirational ideas. A module is a function with a unique name defined by selecting and chopping off branches of a subtree selected randomly from an individual. The approach uses the *compression* operator to select blocks of code for creating new modules, which are introduced into a genetic library and may be invoked by other programs in the population. Two effects are achieved. First the expressiveness of the base language is increased. Second modules become frozen portions of genetic material, which are not subject to genetic operations unless they are subsequently decompressed.

It has been conjectured that problems whose solutions present symmetry patterns or opportunities to parameterize and reuse code can be solved easier in ADF-GP (Koza 1994b) but there exists no formal explanation of why ADF-GP works better than standard GP. (Kinnear 1994) explains why ADF-GP works by introducing the notion of structural regularity. He compares ADF-GP against the module acquisition approach and points out that the module acquisition approach does not directly create structural regularity. Kinnear attributes the better performance of ADF to the repeated use of calls to automatically defined functions and to the multiple use of parameters.

The *lens effect* (Koza 1994b) is the idea that the tails of the fitness distribution for randomly generated programs are larger for the ADF-GP than for standard GP. The effect is attributed to the introduction of new functions into the representation. (Altenberg 1994) outlines that a similar property should be observed in general in order to make GP search more efficient than

random search: the upper tail of the offspring fitness distribution should be wider than that for random search.

The problem of determining the appropriate architectural choices in ADF-GP has generated work on evolution of the GP architecture. The architecture itself can be evolutionarily selected in case the initial population is architecturally diverse and care is taken when crossing over individuals having different architectures (Koza 1994b). (Koza 1994a) introduces six new genetic operations for altering the architecture of an individual program: branch duplication, argument duplication, branch deletion, argument deletion, branch creation and argument creation. These operations are causal in the sense discussed later in this paper.

A rule of thumb in GA literature postulates that population diversity is important for avoiding premature convergence. A comparison of research on this topic is provided in (Ryan 1994). Ryan shows that maintaining increased diversity in GP leads to better performance. His algorithm is called "disassortative mating" because it selects parents for crossover from two different lists of individuals. One list of individuals is ranked based on fitness while the other is ranked based on the sum of size and weighted fitness. The goal is to evolve solutions of minimal size that solve the problem. However, by using directly the size constraint the GP algorithm is prevented from finding solutions. The algorithm improves convergence to a better optimum while maintaining speed.

Exploration and exploitation are recurring themes in search and learning problems (Holland 1992), (Kaelbling 1993). Exploitation takes place when search proceeds based on the action prescribed by the current system knowledge. Exploration is usually based on random actions, taken in order to experiment with more situations. For example, in learning classifier systems, roulette wheel action selection is a means of choosing exploratory actions. In the reinforcement phase of the control loop of a classifier system (Wilson 1994), matching classifiers that do not get activated are weakened. This lowers the chances of choosing unpromising actions in the near future. The weakening magnitude is usually controlled by an explicit parameter, although more elaborate schemes are possible (Wilson 1994). In contrast GP is a search technique that implicitly balances exploration and exploitation, as will be showed later.

3 THE ROLE OF THE GP REPRESENTATION

One goal of the paper is to analyzes the influence of different representational choices on the behavior of GP both theoretically and experimentally using a standard GP algorithm and HGP. The test case chosen is the parity problem. Parity is an attractive problem for several reasons. First it operates on a finite sample space, the space of Boolean functions with a given number of inputs. This enables the computation of distributions of interest for random choices of an initial population. Second, parity is difficult to learn because every time an input bit is flipped, the output also changes.

The ODD-n-PARITY problem is to find a logical composition of primitive Boolean functions that computes the sum of input bits over the field of integers modulo 2. EVEN-n-PARITY can be defined by flipping the result of

ODD-n-PARITY. The ODD-n-PARITY and EVEN-n-PARITY functions appear to be difficult to learn in GP, especially for values of n greater than five (Koza 1992).

The initial function set for the parity problem in GP is defined by the set of primitive Boolean functions of two variables:

$$\mathcal{F}_0 = \{AND, OR, NAND, NOR\} \tag{1}$$

The terminal set is defined by a set of Boolean variables:

$$\mathcal{T}_0 = \{D_0, D_1, D_2, ..., D_{n-1}\}$$

Any Boolean function of n variables is defined on the set of 2^n combinations of input values. Given a program implementing a Boolean function, its performance is computed on all possible combinations of Boolean values for the input variables and is compared with a table defining the EVEN-n-PARITY function. Each time the program and the EVEN-n-PARITY table give the same result, the program records a hit. The task is to discover a program that achieves the maximum number (2^n) of hits.

Theoretical analysis of uniform random sampling

The efficiency of a GP algorithm depends on the computational effort needed to generate a solution with a given probability. An important reference case is the random case. The probability of randomly generating a problem solution depends both on the initial function set, and on the method of generating random individuals. Our goal is to understand the influence of the function set composition, and consequently of a function discovery mechanism on this probability.

Let us consider the sample space of all functions

$$\mathcal{S} = \{ f : \mathcal{B}^n \longrightarrow \mathcal{B} \}$$

where $\mathcal{B} = \{0, 1\}$. Note that $||\mathcal{S}|| = 2^{2^n}$, thus we can obtain random elements of \mathcal{S} by flipping 2^n distinct fair coins.

Consider the random variable X mapping the finite sample space S onto the set of positive integer numbers N defined as follows: X is the number of hits of a randomly generated Boolean function $s \in S$.

We are interested in analyzing the probability mass function of X.

$$Prob\{X = x\} = \sum_{s \in \mathcal{S}: X(s) = x} Prob\{s\}$$

Consider a random (here and in turn the term random refers to structures generated randomly according to a uniform probability distribution) Boolean function with k hits. The k hits are due to i 1-hits and to (k-i) 0-hits. EVEN-n-PARITY takes an equal number (i.e. $\frac{n}{2}$) of 0 and 1 values over the set of input binary strings. Thus, the number of Boolean functions that coincide with EVEN-n-PARITY for a fixed set of k input strings is :

$$\sum_{i=0}^{k} \left(\begin{array}{c} \frac{n}{2} \\ i \end{array} \right) \left(\begin{array}{c} \frac{n}{2} \\ k-i \end{array} \right) = \left(\begin{array}{c} n \\ k \end{array} \right) \tag{2}$$

which implies that X has a binomial distribution, with $p=q=\frac{1}{2}$. In order to prove equality 2 use Newton's binomial on both sides of the identity $(1+x)^n=(1+x)^{\frac{n}{2}}(1+x)^{\frac{n}{2}}$ and identify coefficients. It follows that

$$Prob\{X=k\} = \frac{1}{2^n} \cdot \binom{n}{k}$$
 (3)

The expected value of X is $\frac{n}{2}$ and its variance is $\frac{n}{4}$ (Cormen *et al.* 1990).

Program diversity

In order to understand the role of representation and the effect of dynamically changing it we designed a set of experiments for estimating qualitative measures of diversity such as fitness distributions and program size in GP and HGP.

A straightforward definition of diversity in GP is the percentage of structurally distinct individuals at a given generation. Two individuals are structurally distinct if they are not isomorphic trees. However, such a definition is not practically useful. It is computationally expensive to test for tree isomorphisms. Moreover, associativity of functions is extremely difficult to take into account.

An easily observable type of variation in the population is fitness diversity. Two individuals are different if they score differently.

Another useful qualitative measure of diversity is program size. In HGP, a true measure of the size of an individual is obtained by counting all the nodes in the tree resulting after an "inline" expansion of all the called functions down to the primitive functions. This complexity measure is called "expanded structural complexity" in (Rosca and Ballard 1994b) and is based on the structural complexity (i.e. the number of tree nodes) of all the functions in the hierarchy which are called directly or indirectly by the individual. The expanded structural complexity of a program F, denoted IC(F), can be computed in a bottom-up manner starting with the lowest functions in the call graph of F. For each subfunction G, called directly or indirectly by F, IC(G) can be defined using a recursive formula (Rosca 1995).

Three experiments are reported next. First, a uniform random generation of parity tables is compared to GP random generation of program trees. Second, we vary the composition of the primitive function set and analyze again the fitness distribution of randomly generated GP programs. Third, we analyze the expanded structural complexity of GP and HGP solutions. The method of generating GP individuals in the second experiment, borrowed from (Koza 1992), is the ramped-half-and-half method. In order to create an initial population of increased diversity this method generates trees of depth varying modulo the initial maximum size (taken to be six) and of either balanced or random shape.

Diversity experiments

If elements $s \in \mathcal{S}$ are generated uniformly then the probability of generating EVEN-n-PARITY is $\frac{1}{2^{2n}}$. For the EVEN-3-PARITY problem (Koza 1994b) reports

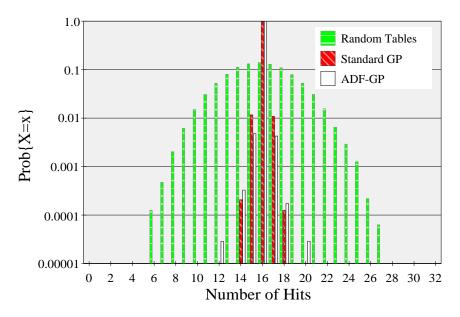


Figure 2: Probability mass function of the random variable X representing the number of hits with the EVEN-5-PARITY function in three cases (a) Random generation of Boolean tables; (b) Random generation of standard GP EVEN-5-PARITY functions; (c) Random generation of ADF-GP EVEN-5-PARITY functions, with two automatically defined functions and two arguments each.

that no solution is discovered after the random generation of 10 million parity functions. However, the above analysis implies that for n=3 it should be considerably easier (one in 256 trees) to find a solution if the random generation of trees in GP results in a uniform distribution of functions. About four billion GP trees would need to be generated in order to find one that computes EVEN-5-PARITY (n=5).

Figure 2 compares the distribution of hits obtained for a population of tables (ideal case), GP functions and ADF-GP functions. The mean and standard deviation of the distribution of randomly generated tables compares closely to the theoretical results outlined above for n=5, although only 16,000 random tables were generated. The distribution of GP functions in the EVEN-5-PARITY problem, with the function set defined in (1), shows that for h<12 or h>20 the probability of having h hits is practically zero. The GP random distribution is much narrower than might be anticipated.

It is worth examining what happens when automatically defined functions are used. Figure 2 shows that a random population of ADF-GP trees generated using the ramped-half-and-half method has a wider distribution of hits than standard GP. The effect is called the lens effect in (Koza 1994b).

Figure 3 shows similar hit distributions for GP when the composition of the function set is varied. In the *all-functions* plots, all 16 Boolean functions are included in \mathcal{F} while in the *some-functions* plots a random selection of half of the 16 functions, including the primitive ones in \mathcal{F}_0 , are included in the initial function set \mathcal{F} . Figure 3 shows the same experiment performed with ADF-GP. Two defined functions have been used. ADF0 has the function

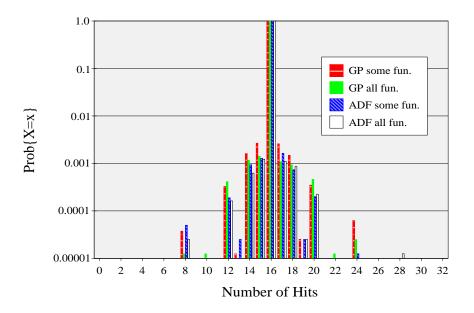


Figure 3: Probability mass function of the number of hits when all or some (a random selection of) Boolean functions of two variables are used in generating EVEN-5-PARITY programs. The primitive set include the primitive functions AND, OR, NAND, NOR.

set \mathcal{F}_0 . ADF1 can additionally invoke ADF0. The program body can additionally invoke both ADF0 and ADF1.

Table 1 presents a sample of complexity results obtained with the standard GP algorithm and with ADF-GP. The rows having 0 in the "Generation" column correspond to an initial random generation of programs. The other two rows are the results at the end of successful runs. The ADF-GP rows include the structural complexity values obtained for the two evolved subfunctions (ADF0 and ADF1) and the main program body (Body). The table shows that expanded complexity in ADF-GP is several orders of magnitude higher than the structural complexity of programs in standard GP.

Table 1: Complexity results for EVEN-5-PARITY tested with the standard GP and ADF-GP algorithms. Average values are determined from 12 runs.

	Generation	Structural Complexity			Expanded Complexity	
Method		ADF0	ADF1	Body	Best	Average
Std.GP	0	-	-	-	15	6.53
Std.GP	28	-	-	-	180	241.2
ADF-GP	0	15	15	45	423	439.9
ADF-GP	30	41	13	95	5497	6429.3

Comparison of results

The narrow GP hit distribution suggests a low population diversity. A solution by means of GP will be difficult to obtain, because it would require more generations, and thus an increased computational effort, to create diverse individuals. Moreover, search may be successful provided that fitness-proportionate selection and the genetic operators used do not narrow the population diversity even more. This change in the hit distribution for HGP is a direct result of the introduction of higher level functions into the representation. It is one of the hypotheses explaining why HGP approaches work better than standard GP.

When the function set is varied an even wider distribution will result (see the *GP-some-functions* and *GP-all-functions* distributions from figure 3). When defined functions are used the hit distribution does not become much wider. However the ADF-GP method still generates larger standard deviations and thus increased diversity. Randomly generated programs with the highest number of hits (28) were obtained using this method.

Note that \mathcal{F}_0 is complete in the sense that any Boolean function can be written just using functions from \mathcal{F}_0 . The new Boolean functions of two variables that are added to the function set enable a different organization of the search space. The effect of apparently non-useful functions, initially included in the function set, is beneficial. All new functions, either ADFs or initial extensions of the function set, are based on the initial primitive functions and terminals. Theoretically, the search space remains the same, the space of all programs that can be built based on \mathcal{F}_0 and \mathcal{T}_0 . The sampling of the search space by means of the crossover operator is changed in ADF-GP. Still, any solution that can be obtained by ADF-GP could theoretically be found by GP although the time to find the program would be significantly larger. From the static point of view of creating an initial population, using ADFs is equivalent to considering a larger initial function set.

A more formal interpretation of this remarks can be stated by considering the *closure* requirement in GP (Koza 1992). Closure requires that any function be well defined for any combination of arguments (terminals or results of other function calls) that it may encounter. Suppose that any subtree returns a value from a domain, call it \mathcal{D} , and that the result returned by a program depends on a subset of variables from \mathcal{T} which defines the input space. Define \mathcal{F}_{total} to be the set of functions mapping the input space onto \mathcal{D} . Then an ADF is a function from \mathcal{F}_{total} . ADF-GP may simply be interpreted as GP over an enlarged function set \mathcal{F}_{total} . Over generations, the use of ADFs is equivalent to a dynamic sampling of various functions from this much larger function set.

Naturally, $\mathcal{F} \subset \mathcal{F}_{total}$. It may be difficult to determine the appropriate functions from \mathcal{F}_{total} necessary to solve a given problem. It is unrealistic to consider huge functions sets in either GP or ADF-GP. However, GP can be used to select primitives that can be better combined to yield candidate solution improvements (Koza 1994b).

The increased standard deviation of program hits can be correlated with the increased size and more diverse structure of individuals obtained by using ADF-GP or, similarly, AR-GP. This results in an increased GP exploration of the space of programs. The use of an HGP approach enables the manipulation of a population of higher diversity programs, which positively affects the efficiency of an HGP algorithm for complex problems.

4 HGP EVOLUTION DYNAMICS

The arguments presented so far have analyzed a static picture. The wider hit distributions and the increased expanded structural complexity suggested an increased exploration rate. The focus of attention in this section moves to an analysis of the HGP evolution dynamics through the crossover operator. A simple analysis of the effects of the crossover operator suggests that GP structures are highly unstable. However, a more careful analysis of the way HGP discovers useful structures reveals that selection gradually favors changes with small effects on the individual behavior.

Causality in GP and ADF-GP

The main problem in identifying how HGP works is determining the effects of the crossover operation as reflected in the variation of fitness from parents to offspring. The intuition is that most crossover operations have a harmful effect. In particular, offspring of individuals that are already partially adapted to the "environment" and already have a complex structure are more likely to have a worse fitness. This is close to the conclusions of the role of mutation in natural evolution (Wills 1993).

Consider a partial solution to a hypothesis formation problem obtained using standard GP and represented by a tree T. Consider that T is selected as a parent and it is possible to obtain a solution by modifying T in such a way that a certain subtree T_i is not changed. Consider also that crossover points are chosen with uniform probability over the set of m nodes of T. The probability of choosing a crossover point v that does not lie within T_i is:

$$Prob(Select(v)|v \notin T_i) = 1 - \frac{Size(T_i)}{Size(T)}$$

The bigger T_i is (and this is true in the case of a hypothetical convergence to a solution) the smaller the probability of keeping it unchanged. The GP dynamics of program trees changes shows the phenomenon of *instability* or poor causality of GP structures.

In spite of the apparent non-causality, GP works. Two explanations could be given to this apparent paradox. The first explanation is the *exploitation* of structures preserved in the population. This would be true for a GA, where the crossover operator maintains fixed positions for exchanged alleles. In contrast, GP crossover is non-homologous in the sense that it does not preserve the position of the subtree on which it operates, being allowed to paste a subtree at any tree level. The probability of choosing homologous crossover points in two structurally similar parents in order to transmit functionality from parents is inversely proportional to the square of the average size of the trees. Thus the first explanation has low probability.

The second explanation is that most changes selected for in later stages of evolution are "small" changes, most probably changes at higher tree depths.

Explaining this paradox for ADF-GP follows the same pattern but is even more apparent. The ADF approach attacks the search problem at *different structural levels* simultaneously. GP has to discover both the definitions for a fixed set of sub-functions, each with a predefined number of parameters, and how to combine calls to the automatically defined functions within the main body. This corresponds roughly to discovering a way to decompose the problem and solve the subproblems given only the maximum number of subproblems and the general structure of the subproblems (i.e. the number of parameters). Due to the imposed hierarchical ordering of ADFs we consider ADFs as different structural levels.

During GP search, modifications are alternatively made at each of the structural levels. A code fragment brought from another individual changes its function entirely if it contains calls to other ADFs. For example, consider a piece of code involving calls to lower order ADFs that is pasted into a higher-order function or a main body as a result of a crossover operation. Also suppose the definitions of the ADFs in the two parents are entirely different. Lexical scope dictates the definition to be used when computing a call to a sub-function, so that the calls to ADFs from the piece of code transplanted will refer to the definitions in the new scope, determining a totally different function. This situation is quite frequent because even when parents are similar the probability of a homologous change is very low. It demonstrates the *non-causality* of ADF-GP.

In general, GP and ADF-GP exhibit poor causality. It is useful to visualize how the search for a solution may generally proceed in ADF-GP. Each of the ADF functions represents a different sub-function. Consider the last modification imposed on a program tree before it becomes an acceptable solution. It is very unlikely, but not impossible that this last change has been a change with a large influence, for example a change in one of the functions at the basis of the hierarchy. This situation represents a lucky change. Most probably, though, it was a change at the highest level, in the program body.

We hypothesize as a general principle of GP dynamics that selection most often favors small changes. These changes have the biggest chance of being successful, therefore respecting the causality principle. The effect of this principle is a stabilization on lower level ADFs that will be useful. Evolution freezes such subroutines, looking for changes at higher levels. The selection of potentially useful subroutines generalizing blocks of code drives the adaptation of the problem representation in the AR approach.

Birth certificates

In order to test the above hypothesis we have studied the most recent part of the genealogy tree for EVEN-n-PARITY parity problems. This was done by assigning to each individual a *birth certificate* that specifies its parents and the method of birth (one of ADF0 crossover, ADF1 crossover, main program body crossover or reproduction). We hoped that an analysis of the birth certificates starting with the final solution and tracing its origin backwards, would shed light on the GP dynamics, as hypothesized above. In order to determine the effect of the different types of birth operations we compute a temporally discounted frequency factor for a given solution tree T for each type of birth:

$$frequency(T, type, d) = \left(\frac{1 - \gamma}{1 - \gamma^d}\right) \sum_{i=0}^{k_T} \chi_{\{type\}}(i) \cdot \gamma^{depth(T_i)}$$
 (4)

where k_T is the number of programs in the genealogy tree of T down to a depth d, and $\chi_{\{type\}}(T_i)$ is the characteristic function of ancestor T_i of T, returning 1 if T_i has a birth certificate of birth-type type and 0 otherwise. The scaling factor $\frac{1-\gamma}{1-\gamma^d}$ is a normalizing constant that makes each type of discounted frequencies for a fixed tree T add up to 1. We used a discounted formula to reflect the higher importance of crossover operations from more recent generations.

Table 2 presents the results for several successful runs of ADF-GP for EVEN-5-PARITY, with two ADFs and three arguments each. In most cases, the frequency factors are highest for the program-body or clearly decrease from program-body to ADF1 to ADF0. These results support the earlier conclusion that ADF-GP search relies in most cases on changes at higher and higher structural levels which make it possible to exploit good code fragments that already appeared in the population.

Table 2: Statistics of birth certificates in successful runs of EVEN-5-PARITY using ADF-GP with a zero mutation rate. Each certificate of a given type counts one unit and is temporally discounted with a discount factor $\gamma = 0.8$ based on its age. Only certificates at most 8 generations old have been considered.

GP Run	# Individuals	Birth C	ertificate	Final					
	Explored	ADF0	ADF1	Body	Generation				
1	123,009	0.295	0.0	0.704	32				
2	110,892	0.221	0.472	0.416	32				
3	62,699	0.077	0.526	0.397	17				
4	35,162	0.447	0.102	0.451	9				
5	55,748	0.1	0.214	0.685	15				
6	55,438	0.093	0.202	0.704	15				

The above numerical results have taken into account only a small time window compared to the entire number of generations. A complete picture of the importance of various types of crossover during the entire GP evolution can be constructed from a more detailed analysis of birth certificates. Such an analysis is depicted in figure 4. It suggests the overall importance of a birth certificate type from generation 0 till the solution is found. While the percentage of program-body changes increases, the percentage of ADF changes decreases.

The results of this section suggests that HGP discovers and exploits useful structures in a bottom-up, hierarchical manner. This is the *bottom-up* evolution thesis. AR-GP has an explicit policy for a bottom-up *exploitation* of discovered structures for making search more efficient, while ADF-GP neglects it.

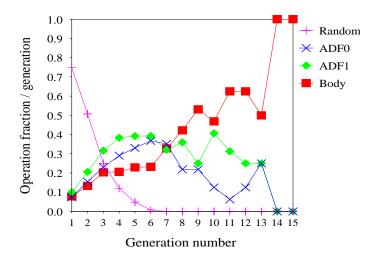


Figure 4: Variation of the fraction of crossover types over generations, while looking for a solution to EVEN-5-PARITY. *Random* indicates the propagation of random individuals from the initial population due to reproduction.

5 DISCUSSION OF RESULTS

The main difficulty of solving complex parity problems in GP is that the computational effort increases exponentially with the size of the problem. Each fitness evaluation becomes more expensive as the problem is scaled up. For instance, the number of fitness cases in the EVEN-6-PARITY problem is twice that of the EVEN-5-PARITY problem and it doubles with every unitary increase in the order of the problem. Moreover, the number of fitness evaluations necessary to find a solution with high probability increases with problem size (Koza 1994b). One explanation of the poor convergence is the inability of standard GP to exploit opportunities for code generalization and reuse. In contrast, by using ADFs or adapting the representation as in AR-GP the same problems can be solved more easily (Koza 1994b), (Rosca and Ballard 1994c). We gave a qualitative explanation of the improved behavior of HGP, based on an analysis of the evolution process on two dimensions: diversity and causality. Next we relate these ideas to the tradeoff between exploration and exploitation.

This paper shows that there exists an implicit bias in the random generation of GP solution encodings which confines population diversity. Diversity increases as a result of changes in representation (the lens effect in (Koza 1994b)). The expanded structural complexity of HGP individuals also increases diversity too and highlights the distinctive size and shape of trees as compared to standard GP. Increased diversity is related to an increased exploration of the space of programs.

The results presented confirm that as the population evolves, increasingly causal changes become more important and are selected. Low level crossover changes in the function hierarchy are highly non-causal and have an increased exploratory role. Exploitative changes are adopted later in the process as the average program fitness increases.

The stabilization of changes in the function hierarchy occurs bottom-up. The GP search process exploits the structures already discovered although it does not avoid spending unnecessary effort with state space exploration. Useful genetic operations become more and more causal. Thus, causality is correlated with search space exploitation. This suggests that GP is able to dynamically balance exploration and exploitation.

The operations of addition for the evolution of architecture in (Koza 1994a) respect the principle of strong causality. They are executed so that they preserve the behavior of the resulting programs. They also decrease the probability that a future random change will drastically alter the behavior of the program. However, the deletion operations do not possess the nice properties mentioned above. They have the antagonist role of confining the increase in size of the evolved programs.

6 CONCLUSIONS

This paper presents a unifying view of the two approaches to the discovery of functions, ADF-GP and AR-GP emphasizing the hierarchical structure of the resulting problem representation. It showed that the exploration of the search space in GP heavily depends on the power of the discovery and evolution of functions. This is turn, resides in the ability of an evolutionary computation to amplify the exploration power of the genetic search process and to exploit favorable changes.

A hierarchy of functions works as an amplifier of good as well as bad effects. HGP benefits in both cases. Normally, GP operates in a state space where solutions are scarce. Undesirable effects compose together resulting in a poorly rated behavior and help in eliminating a bad individual from the search frontier.

Discovery of functions is also an adaptive mechanism for trading off exploration and exploitation in GP. Most often the control structure of a search algorithm explicitly balances exploration and exploitation by means of control parameters. In contrast, GP is a search technique that implicitly balances exploration and exploitation.

The emergent structure in ADF-GP as an effect of causality is an explicit policy in AR-GP. The bottom-up evolution of HGP discussed in the paper justifies this explicit search for building blocks and the expansion of the problem representation, which was successfully used in AR-GP (Rosca and Ballard 1994a).

Future work aims at an improved version of AR-GP that will additionally allow for evolution of functions. One could use the insights about the dynamics of GP search presented in this paper to come up with a more refined and efficient GP-like system, that involves automatic adaptation of control parameters.

Acknowledgments

I would like to thank Peter Angeline for valuable suggestions made on several versions of this paper. I also thank David Fogel and Jonas Karlsson for comments on an earlier draft. Last but not least I thank Dana Ballard for many inspiring discussions. This research was sponsored by the National Science Foundation under grant numbered IRI-9406481 and by DARPA research grant numbered MDA972-92-J-1012.

References

Altenberg, L. (1994). The evolution of evolvability. In K. Kinnear, editor, *Advances in Genetic Programming*. MIT Press.

Angeline, P. J. (1994) Genetic programming and emergent intelligence. In K. Kinnear, editor, *Advances in Genetic Programming*. MIT Press.

Angeline, P. J., and J. B. Pollack (1994). Coevolving high level representations. In C. G. Langton, editor, *Artificial Life III, SFI Studies in the Sciences of Complexity*, volume XVII, pages 55–71, Redwood City, CA, Addison-Wesley.

Cormen, T. H., and C. E. Leiserson, and R. L. Rivest., (1990) *Introduction to Algorithms*. MIT Press.

Cramer, N. L. (1985). A representation for the adaptive generation of simple sequential programs. In *Proceedings of the First International Conference on Genetic Algorithms*. Morgan Kaufmann Publishers, Inc.

Holland, J. H. (1992). *Adaptation in Natural and Artificial Systems, An Introductory Analysis with Applications to Biology, Control and Artificial Intelligence*. MIT Press, 2nd edition.

Kaelbling, L. P. (1993). Learning in Embedded Systems. MIT Press.

Kinnear, K. (1994). Alternatives in automatic function definition. In K. Kinnear, editor, *Advances in Genetic Programming*. MIT Press.

Koza, J. R. (1992). *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. MIT Press, Cambridge, Massachusetts.

Koza, J. R. (1994). Architecture-altering operations for evolving the architecture of a multi-part program in genetic programming. Computer Science Department STAN-CS-TR-94-1528, Stanford University.

Koza, J. R. (1994). *Genetic Programming II*. MIT Press, Cambridge, Massachusetts.

Lohmann, P. (1992). Structure evolution and incomplete induction. In *Parallel Problem Solving from Nature* 2, pages 175–185. Elsevier Science Publishers.

O'Reilly, U.-M. and F. Oppacher, (1994). The troubling aspects of a building block hypothesis for genetic programming. In *Proceedings of the Third Workshop on Foundations of Genetic Algorithms*. Morgan Kaufmann Publishers, Inc.

Rechenberg, I. (1994). Evolution strategy. In J. M. Zurada, R. J. Marks-II, and C. J. Robinson, editors, *Computational Intelligence - Imitating Life*, pages 147–159. IEEE Press.

Rosca, J. P. (1995). An analysis of hierarchical genetic programming. Technical Report 566, University of Rochester, Computer Science Department, January.

Rosca, J. P. and D. H. Ballard (1994). Genetic programming with adaptive representations. Technical Report 489, University of Rochester, Computer Science Department.

Rosca, J. P. and D. H. Ballard (1994). Hierarchical self-organization in genetic programming. In *Proceedings of the Eleventh International Conference on Machine Learning*. Morgan Kaufmann Publishers, Inc.

Rosca, J. P. and D. H. Ballard (1994). Learning by adapting representations in genetic programming. In *Proceedings of the IEEE World Congress on Computational Intelligence*. IEEE Press, Orlando.

Ryan, C. O. (1994). Pygmies and civil servants. In K. Kinnear, editor, *Advances in Genetic Programming*. MIT Press.

Wills, C. (1993). The Runaway Brain. BasicBooks.

Wilson, S. W. (1994). Zcs: A zeroth level classifier system. *Evolutionary Computation*, 2(1):1–18.