



<b>Title</b>	Zero is not a Four Letter Word: Studies in the Evolution of Language
<b>Authors(s)</b>	Stephens, Christopher R., Nicolau, Miguel, Ryan, Conor
<b>Publication date</b>	2005-04-01
<b>Publication information</b>	Stephens, Christopher R., Miguel Nicolau, and Conor Ryan. "Zero Is Not a Four Letter Word: Studies in the Evolution of Language." Springer, April 1, 2005. <a href="https://doi.org/10.1007/978-3-540-31989-4_34">https://doi.org/10.1007/978-3-540-31989-4_34</a> .
<b>Conference details</b>	Genetic Programming, 8th European Conference, EuroGP 2005, Lausanne, Switzerland, 30 March - 1 April 2005
<b>Series</b>	Lecture Notes in Computer Science
<b>Publisher</b>	Springer
<b>Item record/more information</b>	<a href="http://hdl.handle.net/10197/8275">http://hdl.handle.net/10197/8275</a>
<b>Publisher's statement</b>	The final publication is available at Springer via <a href="http://dx.doi.org/10.1007/978-3-540-31989-4_34">http://dx.doi.org/10.1007/978-3-540-31989-4_34</a> .
<b>Publisher's version (DOI)</b>	<a href="https://doi.org/10.1007/978-3-540-31989-4_34">10.1007/978-3-540-31989-4_34</a>

Downloaded 2026-05-01 23:35:10

The UCD community has made this article openly available. Please share how this access benefits you. Your story matters! (@ucd\_oa)



© Some rights reserved. For more information

# Zero is not a four letter word : Studies in the evolution of language.

Chris Stephens<sup>1</sup>, Miguel Nicolau<sup>2</sup>, and Conor Ryan<sup>2</sup>

<sup>1</sup> Instituto de Ciencias Nucleares, Universidad Nacional Autonoma de Mexico  
`stephens@nuclecu.unam.mx`

School of Theoretical Physics, Dublin Institute for Advanced Studies  
`stephens@stp.dias.ie`

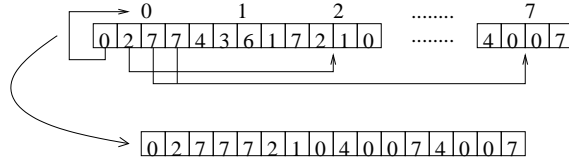
<sup>2</sup> Department Of Computer Science And Information Systems  
University of Limerick, Ireland  
{Conor.Ryan, Miguel.Nicolau}@ul.ie

**Abstract.** We examine a model genetic system that has features of both genetic programming and genetic regulatory networks, to show how various forms of degeneracy in the genotype-phenotype map can induce complex and subtle behaviour in the dynamics that lead to enhanced evolutionary robustness and can be fruitfully described in terms of an elementary algorithmic “language”.

## 1 Introduction

Evolutionary algorithms (EAs) have been applied with a remarkable degree of success to a large variety of problems. However, this is often done with little or no understanding of the dynamics of the system, and practitioners often find themselves unable to explain why a particular tweak has apparently improved their system. Features such as degeneracy and neutral evolution are generally accepted to aid evolution, but little detailed work, particularly in Genetic Programming (GP), has been carried out to investigate how GP exploits these features (with a few notable exceptions [1]). This paper presents a study in evolutionary dynamics, demonstrating how a detailed analysis detects complex and subtle structure formation. Phenomena such as competing conventions, robustness and redundancy are examined and we demonstrate how it is natural to describe these phenomena in the framework of natural language.

In section 2 we describe the representation of our system, showing how it has all the salient features of GP systems, as well as a Genotype-Phenotype map (GPM) inspired by genetic regulatory networks. We also describe the different types of degeneracy inherent in the system. Section 3 takes an in depth look at a representative run, illustrating some surprising strategies that are adopted by the system. The paper concludes with a summary and discusses areas in which our system can be used to develop a deeper understanding of bloat, degeneracy and neutrality, all of which are crucial to the development of more robust EAs.



**Fig. 1.** An example of the switchboard gene in operation. Each codon in the switchboard gene (in the current version, always located in the first position) acts as an index into the *entire* genome. Notice how this particular switchboard gene indexes itself.

## 2 Representation

The representation we use is based on two existing systems; Grammatical Evolution (GE) [2] and that of [3], where a GPM inspired by gene expression models and the phenomenon of cellular division was used. Both exhibit genetic redundancy: GE via a *degenerate* mapping scheme, where many *codons* map to the same item, while in [3] the redundancy was at the gene expression level, in that only a certain proportion of genes from an individual needed to be expressed to produce a fully specified phenotype.

The underlying representation is similar to GE in that we employ a binary string representation, but in this case, a fixed length string is used. The genotype consists of  $N_g$  genes, with each gene consisting of  $N_c$  codons, where each codon takes a symbolic value taken from an alphabet of size  $N_a$ . In our experiments we consider  $N_c = 4$  and  $N_g = N_a = 8$ , so each codon is described by three bits.

The first step is to *transcribe* the genome from the binary representation to eight genes of four codons each. Once that is done, the first gene, known as the *switchboard* gene, establishes which genes are inhibited or promoted, by using each of its four codons as indexes to the promoted genes. Then, in a manner not entirely dissimilar to cellular growth, the initial structure is replaced with a new one, consisting of four genes. Fig. 1 illustrates this process.

Once the activated genes have been identified, the system reverts to a standard GE type mapping, with each codon being used to make a choice in a grammar. Consider the grammar below, with each production rule numbered. As the codons are represented by three bits, each decodes to a value from 0 to 7. This is referred to as a *closed* grammar [2], that is, there is only one non-terminal and, hence, just a single context.

$\langle e \rangle ::= \tanh(\langle e \rangle)$	(0)	$\text{add}(\langle e \rangle, \langle e \rangle)$	(4)
$\tanh(\langle e \rangle)$	(1)	$\text{add}(\langle e \rangle, \langle e \rangle)$	(5)
$\tanh(\langle e \rangle)$	(2)	X	(6)
$\tanh(\langle e \rangle)$	(3)	0	(7)

Consider the following individual, already reduced to its activated genes:

4567 1623 0021 4401

Each codon will *always* make the same choice, e.g. codon 7 will always perform the mapping  $e \rightarrow 0$ . The mapping steps are as follows:

```
4 -> add(e,e)
5 -> add(add(e,e),e)
6 -> add(add(X,e),e)
7 -> add(add(X,0),e)
1 -> add(add(X,0),tanh(e))
6 -> add(add(X,0),tanh(X))
```

Notice that only six codons are used for the mapping; in the case where all codons have been used and still non-terminals remain, a fitness of zero is given. This straightforward mapping scheme permits us to use a more convenient notational representation to facilitate human interpretation of codons. That is, a shorthand for each codon value can be inserted as follows: 0, 1, 2, 3  $\rightarrow$  h; 4, 5  $\rightarrow$  +; 6  $\rightarrow$  X and 7  $\rightarrow$  0. Thus, the individual above can be rewritten as:

```
++X0 hXhh hhhh ++hh
```

We refer to genes described in this way as *words*. Because only six codons were used, we could, using schema notation, describe the above “sentence” as

```
++X0 hX** **** ****
```

Notice that the second word is not made up of four distinct codons/letters. More generally, the last word *used* can vary in length from one to four letters.

## 2.1 Neutrality

Degeneracy leads to the existence of neutral networks [4], where individuals from different areas in the search space have the same fitness. In this work degeneracy exists at several levels; that is, there are several ways of describing the same functionality, e.g.  $add(X, 0)$  and  $add(0, X)$ . Operators are neutral when the individual they produce is genetically different to the individual that they operated on, but phenotypically the same. Point mutation can be neutral at several of the levels above, e.g. changing a codon value such that it still selects the same rule as before, changing a value on the switchboard gene so that it generates the same word as before, but from a different gene, etc. It is also possible to perform neutral crossover. For example, crossover might only effect non-activated genes, both parents might have a copy of a required gene, etc.

## 2.2 Gene expression as a language

The kinds of words that one would expect to be produced by this grammar depend on the fitness function. Consider the function  $f(X) = 4X$ ; this can be described by using a number of different sentences. One possible solution is:

```
+++X XXX*
```

This particular sentence maps to a minimal solution of  $(+(+XX)(+XX))$ . Another solution, which maps to the same phenotype is :

+X+X +XX\*

Notice how these sentences are fundamentally different because their first words differ. This does not suggest a lack of robustness, however, as one would not expect a single evolving population to balance two such different solutions at the same time for very long, although, as described below, it is possible for a number of distinct optimal solutions to appear throughout a run, often competing with each other for dominance of the population, until one becomes extinct.

The role of modularity in most complex problem solving systems, including nature, cannot be over estimated. Difficult problems are often best solved by decomposing them into a set of smaller ones, each of which can be solved more easily than the whole. Similarly, simple modules or strategies which can be reused several times, either on different problems or while solving a single problem are likely to be preferred over more complex ones. Consider the individual:

+++X +++X X000 X000

This maps to  $(+(+X(+XX)(+00))(+0X))$ , and can be reduced to  $4X$ .

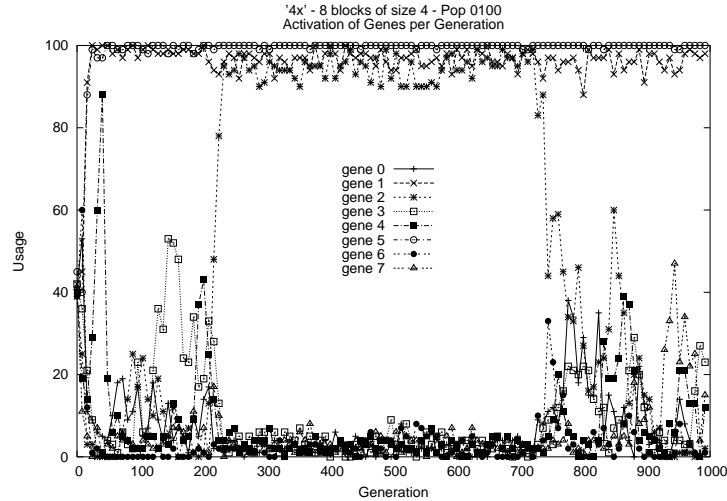
### 3 Results

We applied the system to the problem of performing symbolic regression on the function  $f(X) = 4X$ . By normal GP standards, this is a trivial problem which one would expect to appear in a reasonably sized population with a good initialisation scheme. However, we are concerned with making a detailed analysis of the dynamics, and a simple function like this keeps the analysis tractable.

In total 30 runs were conducted, all of which discovered an optimal solution. Typically, the solutions initially consisted of three or four expressed genes, but shorter solutions almost always appeared, reducing the length to usually two or sometimes three genes. Repeated activation of the same gene was ubiquitous. Typically, each run discovered several ways to represent an optimal solution. A population of 100 individuals was used with a mutation rate, implemented at the bit level, of 0.01. One-point recombination was used with probability 0.9 and restricted to occur only at the boundaries between genes. For selection, a rank-based method was used, where the ranking was applied only to individuals that successfully mapped onto syntactically correct expressions.

#### 3.1 Description of Algorithmic Language

In this section we consider a detailed description of a particular run in order to show the complexity of the dynamics associated with the GPM, even in the case of our very simple search problem. With the production rules specified in Section 2, starting off with a random population one finds, as expected, that the

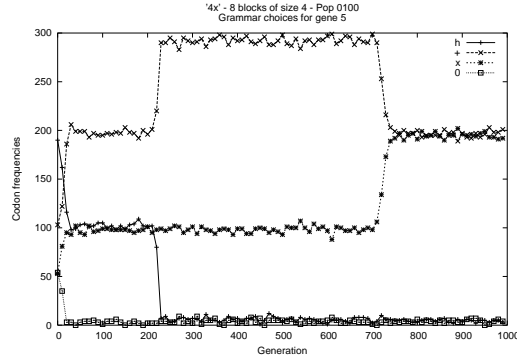


**Fig. 2.** Gene activity per generation. The graph plots the percentage of individuals per generation that have a particular gene activated.

initial codon frequencies are approximately:  $h = 50\%$ ,  $+$  = 25%,  $X = 12.5\%$  and  $0 = 12.5\%$ . Later on, however, one sees structure begins to emerge - for instance, usage of codon  $h$  is significantly less over the majority of the run, while usage of codon  $X$  is significantly greater, as it plays a useful role in fit solutions. With codon 0 there is a greatly increased usage over the middle part of the run, for a reason that will become apparent shortly.

As on the switchboard gene the distribution of codons is random, one finds initially a random distribution of activated genes, i.e. that any gene is active on average in 50% of the population. Later on, however, gene activity patterns emerge with more structure, as can be seen in Fig. 2. Very early on in the evolution an ordering convention is arrived at, whereby the content of the switchboard gene is largely fixed, with more than 90% of the individuals activating genes 5 and 1. Interestingly, there is a redundant usage, whereby these codons are very frequently repeated in the switchboard. As we shall see this leads to enhanced robustness. At generation 25 more than 95% of switchboard genes in the population are of the form  $5 **$  while over 60% are of the form  $55 **$ . Similarly, 92% possess  $** 1*$  and 83% possess  $** 11$ .

In Fig 3 we show the different codon frequencies in gene 5 as a function of time. Before an optimal solution is found gene 5 consists of almost 50% of codons of type  $+$ , i.e. producers. This is about double what would be expected with a random distribution. Similarly, block 1 contained nearly 90% more codons of type  $X$  than would be expected by chance. Gene 5 is clearly the most important, or “core”, gene as its codon content is so stable that there must be strong selection pressure in order to maintain this stability. Gene 1 by contrast showed more



**Fig. 3.** Grammar choices for gene 5 per generation.

variation, though the important role played by codons  $X$  and  $0$  was evident. Note that when expressed, gene 5 precedes gene 1, hence, as expected, the system puts more producers to the left and more terminals to the right. It is interesting to note that optimal solutions are detected from time to time but it is not until after generation 200 that they become established in the population. Also, fitter solutions tend to use only two different active genes - 5 and 1 - and three expressed ones - 551\* - the final non-coding tail not being expressed. The most common phenotype is  $f(X) = 2X + 2 \tanh(X)$ , that results from a combination of a repeated producing gene 5 = ++hX and a consuming gene 1 = XX00.

At generation 213 optimal solutions finally begin to successfully propagate through the population. At this point the number of + codons in gene 5 increases by 50% while the number of 0 codons in gene 1 doubles. The first solution found is shown below. Notice that the switchboard block appears twice, once in decimal form to facilitate reading, and once in the normal form, ++hh in this case.

5533. ++hh. XX00. 0+h+. X000. ++hX. +++X. ++Xh. +Xhh

(+(+(+X(+(+X X ) 0) 0)) 0)) X) = 4X

Note that this founder does not activate gene 1 as the majority of the population. It is based on the same switchboard template 55\*\* of previous suboptimal solutions, but is achieved through a single mutation of the 5 gene ++hX → +++X which, repeated, combines with gene 3 = X000 to give 4X. 5533 however, is not the dominant switchboard gene, that role being kept by 5512, and therefore the dominant optimal solution had 5512 as a founder switchboard. Moreover, with this switchboard, in general all four genes 5, 5, 1 and 2 are expressed, as can be seen in the huge increase in use of gene 2 after generation 200 in Fig. 2.

Given that all these solutions are optimal one might expect that, on average, evolution preserves their structure, apart from the effect of neutral drift. However, this is not the case. For the next 500 generations this solution spreads and evolves. Early on, among the optimal solutions, in gene 1 over 60% of the

individuals have a + codon. Later on this percentage has dropped to zero! There is no direct selection pressure for this, as we are talking about the structure of this gene only for optimal individuals. In terms of effective fitness [5] however, there is a clear explanation: any + in gene 1 means that in order to maintain an optimal solution more codons from gene 2 must be expressed. As these are subject to mutational damage there is an effective selection pressure to make the solution more robust within the context of a repeated “core” gene, +++X, that needs a minimum of 5 more terminal codons. The elimination of the + codon reduces the total number of codons that are expressed in the optimal solutions and therefore increases robustness against mutational damage.

At generation 704 a new optimal solution appears - still based on the switchboard gene 5512 but with a mutated core gene 5 of the form +X + X. An immediate advantage of this solution is that it uses fewer expressed codons and therefore will be more mutationally robust. However, another complementary and more subtle effect appears: 55\*\* in the switchboard gene requires a final terminating 0 codon in the first position of the following gene (or some other codon combination, such as hh0\*, that evaluates to zero) in order to provide an optimal solution. In the initial population, when +X + X solution is first found, this is provided by gene 1, associated with the switchboard 551\*, as more than 50% of the individuals that used the + + +X core gene already had a 0 codon in the first position there. However, of the 67 optimal individuals at generation 704 only 5 had a 0 first position codon. 300 generations later, an examination of genes 4, 6 and 7, which are neither activated by the switchboard nor expressed, show that, from the 246 such genes associated with the 82 optimal individuals, 101, i.e. 41% have a 0 in the first position. This is almost three times the percentage (15%) expected if the distribution were random! What is the reason for this extraordinary self-organization - the origin of our somewhat tongue-in-cheek title? The answer is that from 551\* a mutation on the switchboard of the third codon to activate any gene that has a 0 codon in first position would result in an optimal string. In this sense many of the non-activated genes are acting as a genetic “reserve” to protect against mutations of the switchboard gene.

As in [5] one can summarize the algorithmic “language” that has emerged. This is an evolving language, in the sense that the system continually finds “fitter” genes (the “words” of the language) and “fitter” ways of expressing them through the switchboard (the “syntax” of the language), where “fitter” means an effective fitness that also measures evolutionary robustness. In Table 1 we give a description of the algorithmic language that emerged after 1000 generations.

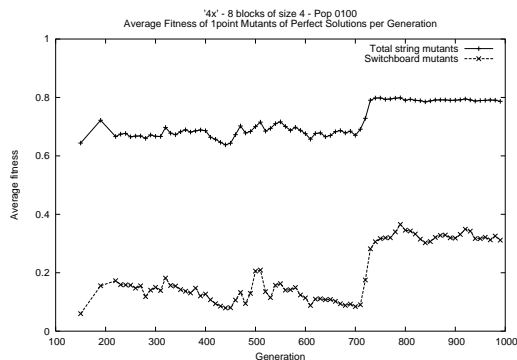
It might seem surprising that the system based itself on (+X(+X(+X(+X 0))), which actually uses more expressed codons - 9 - than the minimal solution (+X(+X(+X X))), which uses only 7 and therefore might be surmised to be more evolutionarily robust. The answer is, of course, that because the first solution uses a repeated building block - +X + X - this block may be expressed twice. In this case the effective number of codons that are subject to mutational damage of (+X(+X(+X(+X 0))) is only 5 compared to the 7 used by (+X(+X(+X X))).

**Table 1.** Description of the algorithmic language that has emerged by generation 1000

Words of the Language		
Gene	Codon content	Logic
Gene 5	+X + X	Gives possibility of at least 2X when expressed once and 4x when expressed twice
Gene 1	0 * **	Terminates the twice expressed 00 * * 5 gene with 0 codons that do not affect the 4X expression
Genes 4,6,7	0 * **	Genetic Reserve - predominantly not expressed but backup in case of mutation in third switchboard codon
Genes 2,3	* * **	“Non-coding” regions - either do not form part of the language or have unknown functionality
Syntax of the Language		
Switchboard	551*	Puts producing gene 5, expressed twice, before consuming gene 1 to given optimal ordering

We have discussed the evolution of robustness and believe that the above description of the results of the experiment offers unequivocal evidence for it. However, one can determine more rigorous and quantitative measures. If one can think of one optimal solution as being more robust than another then this should imply that the more robust state’s “neighbours” are on average fitter, where the notion of neighbour depends on the operator we are thinking of robustness with respect to; for mutation, it is natural to use Hamming distance as a measure of neighbourhood. To this end, for every optimal solution we consider the average fitness of its 96 one-mutant neighbours. In Fig. 4 we see a graph of this quantity as a function of time. The key observation is: upon discovery of the solution (+X(+X(+X(+X 0)))) at generation 704, the system now has a solution that uses in the genotype only five distinct codons, other than the switchboard gene, whereas the previous solution, based on + + +X, used 9 distinct codons. This implies that, as the + + +X solution has  $(32 - 13) = 19$  non-expressed codons, the +X + X-based solution, which uses 4 less should have  $4/19 \sim 20\%$  more optimal one-mutant neighbours which is roughly the increase seen in Fig 4.

The evolution of robustness is even more pronounced if we examine the fitness of the one-mutant neighbours of optimal solutions where we consider only the switchboard gene, i.e. only 12 neighbours. Fig 4 also shows the temporal evolution of the average fitness of the one-mutant neighbours of the switchboard. We have also examined the evolution of the ratio of one-point mutants of both the entire genotype and the switchboard that generate optimal and valid solutions respectively. The behaviour is similar to that for fitness, showing marked increases when a more robust solution is found. Naturally, in the case of the switchboard, due to the important role this gene plays in generating the syntax, and the lack of neutral mutations due to the fact that all four codons are used to activate other genes, one notes that the average fitness of the one-mutant neighbours is small. This is a sign that the switchboard is more brittle than many of the other genes. However, even in the case where the optimal solution uses 4



**Fig. 4.** Average fitness of 1-point mutants of perfect individuals, and of 1-point mutants of each perfect individual’s switchboard.

expressed genes, there is still some degree of robustness. As the average number of optimal mutants can be as high as 13% this means that up to 2 switchboard codons could be mutated and still leave an optimal individual. After generation 704 the new class of optimal solution uses only three activated genes hence the fourth codon on the switchboard loses its importance. Therefore a minimum of 25% of the one-mutant neighbours should be optimal. However, it was observed that 27 – 28% is the norm and hence, the system had evolved robustness above and beyond just finding a solution that uses fewer expressed codons.

While our conclusions thus far have been gleaned from an examination of a single (though representative) run, all our experiences with other runs suggested that the phenomena we observed are common across a large number of runs. It is however legitimate to ask what happens over those runs. The problem with this is that many of the observed phenomena are contingent: the switchboard structure, and subtleties such as genetic reserve, will look quite different in different runs. Two basic related phenomena associated with robustness that can be seen over many runs are: the tendency to activate more than once the same gene - especially the core gene - and a tendency to use fewer expressed genes. In fact, solutions that use two expressed genes occurred considerably more frequently than three-gene solutions, (typically 10-15 times more often) while four-gene solutions are rare indeed. This tendency is the equivalent of the more familiar phenomena of bloat in standard GP. In both case there is a tendency for the system to reach a state where the ratio of coding material to non-coding material is minimized. In GP this is achieved mainly by increasing the amount of non-coding material while here it is by minimizing the amount of coding material. Obviously, the payoff is enhanced evolutionary robustness via resistance to mutational damage.

## 4 Conclusions

In this paper we investigated how the existence of a degenerate GPM can lead to the evolution of robustness and the emergence of an algorithmic language as a result of the self-organization of the GPM.

In distinction to previous work, we concentrated on an in depth analysis of a single run, to give an idea of the tremendous subtlety and complexity of the phenomena that can occur, even in this simple scenario. We saw that the manner in which the system can build robustness can be very varied, from simply developing solutions that require fewer expressed genes, to influencing the content of non-coding parts of the genome and the pattern of gene expression, such as the creation of genetic reserves. We saw and quantified a tendency to reduce the size of the effective coding region - a phenomena analogous to bloat in GP. We saw that robustness can evolve both continuously and in a more punctuated manner, as when passing between solutions with different numbers of expressed genes.

Our study was motivated by a desire to offer a phenomenological predictive framework and description of the evolution of robustness in the context of a genetic model with some language-like features. There exists a formal mathematical framework in which to describe these phenomena - induced symmetry breaking of the genotype-phenotype map and effective fitness as a quantitative measure of this fitness [5]. We will return to a description within this framework at a later date. We believe that further studies of our model and framework will lead to a much deeper understanding of the phenomena of bloat, as well as help in the design of better genetic operators and therefore more competent EAs. A further motivation is that of [3] - to understand the origins and evolution of language.

## References

1. Banzhaf, W.: Genotype-Phenotype Mapping and Neutral Variation - A case study in Genetic Programming. In: Davidor et al., (Eds.): Proceedings of the third conference on Parallel Problem Solving from Nature. Lecture Notes in Computer Science, Vol. 866. Springer-Verlag. (1994) 322-332
2. O'Neill, M. and Ryan, C.: Grammatical Evolution - Evolving programs in an arbitrary language. Kluwer Academic Publishers. (2003)
3. Angeles, O., Stephens, C.R., Waelbroeck, H.: Emergence of Algorithmic Language in Genetic Systems. *BioSystems* **47**. (1998) 129-147
4. Van Nimwegen, E., Crutchfield, J.P., and Huynen, M.: Neutral Evolution of Mutational Robustness. *Proc. Natl. Acad. Sci. USA* **96**. (1996) 9716-9720
5. Stephens, C.R. and Mora, J.: Effective Fitness as an Alternative Paradigm for Evolutionary Computation. *Gen. Prog. Evol. Hardware* **2**. (2000) 7-32.
6. Keller, R. and Banzhaf, W. : Genetic Programming using Genotype-Phenotype Mapping from Linear Genomes into Linear Phenotypes. In: Genetic Programming 1996: Proceedings of the First Annual Conference. MIT Press. (1996)
7. Wong, M. and Leung, K. Inductive Logic Programming Using Genetic Algorithms. I.I.A.S. (1994)