



Title	Understanding Grammatical Evolution: Grammar Design
Authors(s)	Nicolau, Miguel, Agapitos, Alexandros
Publication date	2018-09-12
Publication information	Nicolau, Miguel, and Alexandros Agapitos. "Understanding Grammatical Evolution: Grammar Design." Springer, September 12, 2018. https://doi.org/10.1007/978-3-319-78717-6_2 .
Publisher	Springer
Item record/more information	http://hdl.handle.net/10197/10529
Publisher's statement	The final publication is available at www.springerlink.com .
Publisher's version (DOI)	10.1007/978-3-319-78717-6_2

Downloaded 2026-05-02 00:25:48

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

Understanding Grammatical Evolution: Grammar Design

Miguel Nicolau · Alexandros Agapitos

Received: date / Accepted: date

Abstract A frequently overlooked consideration when using Grammatical Evolution (GE) is grammar design. This is because there is an infinite number of grammars that can specify the same syntax. There are, however, certain aspects of grammar design that greatly affect the speed of convergence and quality of solutions generated with GE. In this chapter, general guidelines for grammar design are presented. These are domain-independent, and can be used when applying GE to any problem. An extensive analysis of their effect and results across a large set of experiments are reported.

1 Introduction

One of the attractive aspects of Grammatical Evolution (GE) is how it can be easily applied to a multitude of problem domains: just design a grammar specifying the syntax of potential solutions, and supply a fitness function to evaluate them.

Easily and *just* are large over-simplifications. While specifying the syntax of solutions with a context-free grammar is a relatively simple task, depending on the problem domain (a multitude of grammars exist in the literature for symbolic regression applications, for example), there is an infinite number of grammars that can specify the same syntax. But not all of them are adequate for use with GE.

In fact, the effectiveness of a grammar is deeply tied to GE's mapping process, and its effect on the search operators. In this study, this effect is analysed, and general guidelines are provided, with the aim of improving GE's search process.

There are many aspects to consider, when designing a grammar for use with GE. Some of these include which and how many non-terminal symbols to use, recursiveness, mapping probabilities, symbol biases, length of derivation sequences,

Miguel Nicolau
College of Business, University College Dublin, Dublin, Ireland
ORCID: orcid.org/0000-0002-1981-1300
E-mail: Miguel.Nicolau@ucd.ie

Alexandros Agapitos
School of Computer Science, University College Dublin, Dublin, Ireland
E-mail: alexagapitos@gmail.com

prefix vs. infix vs. postfix notation, readability/understandability/maintenance of grammars, and many more.

These topics are analysed in this chapter, in terms of initial search space sampling (i.e. their combined effect with initialisation), search effectiveness (combined effect with the search operators), and quality of final solutions (fitness and size of final solutions). A large set of experiments are also executed, for empirical evidence. The results obtained confirm and highlight just how much grammar design can affect the performance of GE.

Based on these findings, a set of general grammar design guidelines are proposed for GE with linear genome representations. Although the resulting grammars can be substantially larger and more complex, most of these transformations can be automatically applied to well-designed base grammars.

2 Previous work

Although a large volume of work exists in the literature on grammar design, particularly in linguistics and computer science, it almost exclusively relates to their use in parsing applications, i.e. syntax verification, compiler design, text mining, etc. In GE, on the contrary, grammars are used in a constructive manner, which, when combined with linear numerical sequences (genotypes) to determine derivation sequences, creates a mostly unique role for grammars.

There is surprisingly little work in the literature on the design of grammars for GE. This is probably due to the remarkable resilient nature of the evolutionary search process: given a correct and reasonably designed grammar, GE tends to produce a working solution. This is not always ideal, however, both in terms of the search effort required, and also the quality of the final solutions produced.

One of the earliest studies of the influence of grammar design on the performance of GE [22] looked specifically at reducing the number of non-terminal symbols in grammars, and proposed an automatic process of achieving this. This resulted in a small increase in performance across all problems attempted.

Hemberg et. al [11] studied the design of grammars using prefix, infix and postfix notation, and their relative performance on a series of symbolic regression problems. The most relevant conclusion is indeed that “the choice of grammar can produce performance advantage”.

One of the most comprehensive analysis of the influence of grammar design in the performance of GE is found in Hemberg’s doctoral thesis [10]. By allowing the grammar to evolve at the same time as the linear genome structures, knowledge is uncovered about the influence of different grammars on the performance of GE. Not surprisingly, it was again concluded that the choice of grammar can influence the performance of GE, for the problems examined.

Byrne et. al [1] analysed two types of mutation events in GE, structural or nodal in nature, and showed how these can be related to exploration and exploitation, respectively. These mutation events are directly dependent on the design of the grammars used.

Harper [9] highlighted the problem of having more production rules adding non-terminal symbols to the mapping sequence than removing them, and the negative impact of this on the performance of GE, if using linear genomes. Nicolau et. al

```

<s> ::= <e>
<e> ::= <e> <o> <e>
      | ( <e> <o> <e> )
      | <f> ( <e> , <e> )
      | <v>
<o> ::= + | - | *
<f> ::= pdiv
<v> ::= x | 1.0

```

Grammar 0: Simple arithmetical expressions grammar. Division uses a protected implementation, termed `pdiv` (more details in Section 5.3).

[27] also analysed the effect of grammar design, focusing on the issue of mapping termination. Both studies highlighted the importance of good grammar design.

Grammar design, and its corresponding effect in representational bias, can also greatly influence the size of the resulting solutions. In Genetic Programming (GP), a possible outcome of this is bloat, i.e. a substantial growth in solution size, with negligible effect in performance increase [20]. Although bloat in GE is not quite as prevalent, studies have shown how grammar design directly influence the generation of very small [26] or very large [9] solutions.

More recently, work has been made on the design of grammars for sorting networks [7] and automatic program synthesis [6], although the results obtained only apply to grammar-based genetic programming systems using derivation trees.

All these studies have a common theme, which is how grammar design in GE can greatly affect its performance. The following section takes this into account, and presents a series of grammar design guidelines, or transformations for existing grammars, aimed at improving the performance of GE in different levels.

3 Grammar Design

A series of grammar transformations are presented in this section. To illustrate their application, Grammar 0 (G0) is used as a starting point. This is a typical grammar for symbolic regression applications with GE, slightly simplified (no unary operators) to illustrate the design techniques presented.

3.1 Balanced Grammars

In order to generate variable-length, unbounded phenotype solutions for any problem, GE makes use of two components:

- Variable-length genotype structures¹;
- Recursively-defined grammar rules.

Both conditions are required, in order to generate phenotype solutions of any size. If genotypes are unbounded but the grammar has no recursively defined symbols, a phenotype solution is only as large as the largest sequence of terminal

¹ Technically, genotypes used with GE are length-bounded, in the sense that they cannot be smaller than zero, or larger than what the memory of the machine running the experiments can hold. This maximum size is, however, a technical limitation, rather than a conceptual bound.

```

<s> ::= <e>
<e> ::= <e> <o> <e> | <v>
      | ( <e> <o> <e> ) | <v>
      | <f> ( <e> , <e> ) | <v>
<o> ::= + | - | *
<f> ::= pdiv
<v> ::= x | 1.0

```

Grammar 1: Balanced recursion grammar.

symbols generated by the longest derivation path through the grammar (some studies [18] use this as both a means to limit solution size, and also to ensure validity of mapping, by always using genotype strings sufficiently long to terminate any derivation sequence).

If on the other hand the grammar is recursive, but the genotype is a fixed-length structure (of length l), the maximum phenotype solution length is the longest derivation path through the grammar smaller or equal to l mapping steps (unless wrapping is used).

The use of recursiveness in grammars with GE appears right from its first publication [38], although some interesting GE applications make use of non-recursive grammars, defining either fixed-length solutions (such as the design of a genetic algorithm using GE [28]), or maximum-length bounded solutions (such as the design of an ant-colony optimisation algorithm using GE [40]).

Recursiveness should be applied with care, however. Single levels of recursiveness are easy to implement and understand (see e.g. the second line of G0), but multiple recursive symbols or multiple line recursiveness easily become hard to design and/or understand. See for example the history of attempts at solving the Santa Fe Ant Trail Problem [15] with GE, from the original incorrect grammar [29], to a first [35] and then second [9] correction, and its analysis [27]. And yet, recent publications [42, 18, 17, 21] still use incorrect grammars, not respecting the original problem syntax.

The influence of grammar recursiveness in the ability to terminate mapping, and thus in the effectiveness of GE, has been studied as early as 2003 [39], where productions were categorised based on whether they added, maintained or reduced the number of non-terminal symbols left to map, when applied.

A way to label productions as recursive or not was proposed by Ryan and Azad [37], within the context of initialisation. Subsequently, Harper [9] labelled grammars as *explosive* or *balanced*, depending on whether there is a higher probability of adding non-terminal symbols during mapping over adding terminal symbols.

These analyses are all related. In this chapter, we define as *explosive* grammars where at least one symbol has more recursive productions than non-recursive. An example is grammar G0: the symbol <e> has three recursive productions associated, and a single non-recursive production. This kind of grammar is explosive, and has a low probability of generating a fully mapped phenotype string, when used in conjunction with a randomly-generated genotype string.

For a grammar to be *balanced*, each recursive production should have a corresponding "consuming" production. Grammar 1 (G1) achieves this, by having a non-recursive production for every recursive production associated with <e>.

Note that this does not alter any biases, other than that of replacing $\langle e \rangle$ using a recursive or non-recursive production. There is also a downside: G1 now has a 50% probability of generating expressions consisting solely of either x or 1.0 .

3.2 Unlinked Productions

When using GE with a linear genome, if the grammar has several non-terminal symbols, the function of a codon (the production it will choose) is dependent on the non-terminal symbol to be mapped, at a given stage of the mapping process. This means that, when crossover occurs, the function of a codon might change.

The potential destructive nature of such changes is addressed in Section 3.3, by reducing the number of non-terminal symbols. But if several non-terminal symbols are needed, production rules associated with different symbols can be *functionally linked*, in the sense that all codon values that choose a specific production for a symbol, if used to make a choice for a different symbol, will always choose another specific production.

Take grammar G1 as an example. Symbol $\langle e \rangle$ has six associated productions, which is a multiple of both the number of productions associated with $\langle o \rangle$ (three) and $\langle v \rangle$ (two). So if codon values transforming $\langle e \rangle$ into $\langle e \rangle \langle o \rangle \langle e \rangle$ are used to choose a production for $\langle o \rangle$ or $\langle v \rangle$, they will always transform them into $+$ and x , respectively. Table 1 illustrates this.

Table 1: Functionally-linked productions of grammar G1: codon values transforming $\langle e \rangle$ into $\langle e \rangle \langle o \rangle \langle e \rangle$ will always transform $\langle o \rangle$ into $+$ and $\langle v \rangle$ into x ; by contrast, $\langle o \rangle$ and $\langle v \rangle$ are unlinked, as codon values transforming $\langle o \rangle$ into $+$ can transform $\langle v \rangle$ into either x or 1.0

codon values	Parity	$\langle e \rangle$	$\langle o \rangle$	$\langle v \rangle$
0,6,12,...	even	$\langle e \rangle \langle o \rangle \langle e \rangle$	$+$	x
1,7,13,...	odd	$\langle v \rangle$	$-$	1.0
2,8,14,...	even	$(\langle e \rangle \langle o \rangle \langle e \rangle)$	$*$	x
3,9,15,...	odd	$\langle v \rangle$	$+$	1.0
4,10,16,...	even	$\langle f \rangle (\langle e \rangle , \langle e \rangle)$	$-$	x
5,11,17,...	odd	$\langle v \rangle$	$*$	1.0

This can introduce biases in the exploration of the search space. As individuals grow in size (only achievable in standard GE through the use of the crossover operator), a larger proportion of even codon values are required at the start of the genome (see Table 1), and a larger proportion of odd values towards the end (to terminate the recursion of the $\langle e \rangle$ symbol). However, this will also result in a higher proportion of symbol x at the start of the genome, and of 1.0 towards the end. Only through later specific mutation events can suitable proportions of x and 1.0 be evolved, assuming the unbalanced solution will survive until then. In other words, the structural function of the $\langle e \rangle$ symbol, and the nodal content function of the $\langle v \rangle$ symbol are linked.

The link between $\langle e \rangle$ and $\langle o \rangle$ is less obvious, but biases also exist: a solution requiring only the use of the $+$, $-$ and $*$ operators will be biased towards $+$ and $*$, when solutions grow in size (first and third productions associated with $\langle e \rangle$).

```

<s> ::= <e>
<e> ::= <e> <o> <e> | <v>
      | ( <e> <o> <e> ) | <v>
      | <f> ( <e> , <e> ) | <v>
<o> ::= + | + | + | + | + | + | +
      | - | - | - | - | - | - | -
      | * | * | * | * | * | * | *
<f> ::= pdiv
<v> ::= x | x | x | x | x | x | x
      | x | x | x | x | x | x | x
      | x | x | x | x | x | x | x
      | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0
      | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0
      | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0

```

Grammar 2: Unlinked productions grammar.

Note that these biases can be either beneficial or detrimental, depending on the problem domain. But for a black-box approach, with no domain knowledge, it is desirable to explore the search space in the most unbiased way.

An example is the typical grammar used with GE for the Max problem [8]. Inconsistent results were obtained with different GE mapping orders [5], and further analysis [4] produced no clear explanations for this issue; it was subsequently shown [27] that the grammar used in those experiments suffered from functionally linked productions.

The problem of linked productions was identified as early as 2002 [13]. The solution proposed then was the adoption of a different mapping procedure, called the *Bucket Rule*. However, it was later shown [27] that these biases can also be removed without modifying the standard GE mapping process, through careful grammar design.

Grammar 2 (G2) illustrates how to achieve this, through production rule repetition. In this case, six copies of each production associated with <o> are introduced (for a total of $6 * 3 = 18$ productions), whereas 18 copies of each production associated with <v> are used, for a total of $18 * 2 = 36$ productions. Note that sufficiently large codon value ranges are required, when using this technique, to ensure minimal production choice biases, but this is good GE practice anyway [32].

3.3 Reduced Non-Terminals

The use of linear genomes and a one point crossover with GE leads to what has been termed the *ripple effect* [31]. This means that, when viewed at a derivation tree level, crossover removes several sub-trees from each parent, which are filled with genetic material from the other parent. Given that the exchanged genetic material consists of a numerical sequence, the actual phenotypic material received may or may not correspond to the original phenotypic material from the other parent: codons reinterpreted under different derivation tree nodes (i.e. non-terminal symbols) will generate different, potentially never-seen before phenotypic material.

This change of interpretation of genetic material can be very damaging to the already fragile locality of crossover in GE [36]. However, many grammars in the

```

<e> ::= <e> + <e> | <e> - <e> | <e> * <e>
      | <e> + <e> | <e> - <e> | <e> * <e>
      | x | x | x
      | 1.0 | 1.0 | 1.0
      | ( <e> + <e> ) | ( <e> - <e> ) | ( <e> * <e> )
      | ( <e> + <e> ) | ( <e> - <e> ) | ( <e> * <e> )
      | x | x | x
      | 1.0 | 1.0 | 1.0
      | pdiv ( <e> , <e> ) | pdiv ( <e> , <e> )
      | pdiv ( <e> , <e> ) | pdiv ( <e> , <e> ) | pdiv ( <e> , <e> )
      | x | x | x
      | 1.0 | 1.0 | 1.0

```

Grammar 3: Single non-terminal grammar.

```

<e> ::= ( <e> + <e> )
      | <e> + <e>
      | ( <e> - <e> )
      | <e> - <e>
      | ( <e> * <e> )
      | <e> * <e>
      | <e> pdiv <e>
      | ( <e> pdiv <e> )
      | x | x | x | x
      | 1.0 | 1.0 | 1.0 | 1.0

```

Grammar 4: Corrected-biases grammar.

literature define more non-terminal symbols than strictly required, creating further cases of reinterpretation of exchanged genetic material.

A solution is thus to reduce the number of non-terminal symbols as much as possible [22]. In fact, for symbolic regression (the most common application domain of GP-like systems [43]), grammars with a single non-terminal symbol can be used (effectively single-type languages, corresponding to GP's closure requirement).

Grammar 3 (G3) shows a single non-terminal symbol version of G1 (using G1 or G2 as a base is irrelevant, given that a single non-terminal symbol remains, <e>, so no linked productions will occur). This process works by replacing non-recursively defined non-terminal symbols by all of their productions, wherever they are used, while keeping the biases of the original grammar. A more detailed description and step by step illustration can be found in the relevant publication [22].

3.4 Grammar Biases

Grammars such as G0 are common in the literature. However, it has a 66.666% bias towards the use of one of the operators (+, -, *), resulting in a 22% bias for each, and a 33.333% bias towards the use of pdiv, which may not be desired.

To ensure an unbiased exploration of the search space, all four operators should have the same biases. This also makes the search space more comparable to that of GP. Grammar 4 (G4) shows a transformation of G3 to take this into account.

```

<e> ::= + <e> <e>
      | - <e> <e>
      | * <e> <e>
      | pdiv <e> <e>
      | x | x
      | 1.0 | 1.0

```

Grammar 5: Prefix-notation grammar.

```

<e> ::= + <e> <e>
      | - <e> <e>
      | * <e> <e>
      | pdiv <e> <e>
      | <v>
      | <v>
      | <v>
      | <v>
<v> ::= x | x | x | x | x | x | x | x | x
      | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0

```

Grammar 6: Compromise transformations for a compact, understandable grammar. Note that unlinking of productions associated with <e> and <v> was required.

3.5 Infix/Prefix Notation

From a mathematical point of view, using a single non-terminal symbol grammar, a prefix or postfix notation will essentially produce the same performance (subject to the stochastic nature of the search process), as they explore the same (inverted) syntax space. However, infix will not, if used both with and without parenthesised expressions. For example, a prefix expression $*xx$ may become $*x+xx$, if the codon encoding the second argument is mutated; with infix notation, however, $x*x$ can become either $x*(x+x)$, which is equivalent, or $x*x+x$, which is not.

There is no obvious choice to make here. Infix provides a more connected search space, but at the expense of further loss of locality for the genetic operators; prefix/postfix do the opposite. They also provide a more comparable search space to that of tree-based GP. Grammar 5 (G5) shows a prefix version of G4.

3.6 Compromise Grammars

Although the transformations presented can be achieved through an algorithmic process (and thus automated), the resulting grammars can easily grow exponentially in both size and complexity, and become very hard to understand or modify. The addition of a carefully chosen single non-terminal symbol (<v> in this case) is often enough to maintain the readability of a grammar, at the expense of a slight worsening of crossover locality. Grammar 6 (G6) illustrates this. Although not apparent when compared to G5, in problems with large numbers of operators and variables, this can drastically reduce the number of productions in a grammar (see Table 5, grammars G5 (5625 productions) and G6 (17 productions)).

4 Transformations Analysis

A series of detailed experiments were performed, using grammars G0-G6, to analyse the effect of grammar design on search space biases, size, termination, repetition, and performance. Table 2 shows the experimental setup used.

Table 2: Experimental Setup

Population Size	500*
Number of generations	50
Random initialisation genome length	31
PTC2 max expansions	31
Maximum genome length	-
Selection Tournament Size	1%
Elitism (for generational replacement)	1%
Crossover Ratio	50%
Average Mutation Events per Individual	1
Max Wrapping events	0

* 200 for Shape Match (Easy and Medium) (see Section 5)

* 1000 for V4, K12, Housing, EPar5 and Mux11 (see Section 5)

* 2000 for Dow (see Section 5)

Populations were initialised using Random Initialisation (RND) (random integer strings), or using a depth-less variant of Probabilistic Tree Creation 2 (PTC2) [19, 23]. These were chosen as many recent publications still use RND initialisation [17, 21], whereas PTC2 was chosen for its proven performance [23]. Depending on the grammar used, the specified number of codons/expansions for initialisation is sufficient to generate expressions of up to a corresponding syntax tree depth of 5.

4.1 Initialisation Biases

This first experiment analyses the initial populations generated using the two initialisation methods. The proportion of each phenotypic symbol (+, -, *, pdiv, 1.0 (const) and x (var)) in all successfully mapped individuals was recorded, along with measures related to the mapping process: average phenotype length, number of invalid (non-mapping) individuals generated, and number of repeated (valid) phenotypic solutions. Fig. 1 shows the results obtained, for 100 independent runs.

The top half of the figure illustrates how grammar design affects the initial sampling of the search space. G0-G3 exhibit a bias towards the use of division. In the first three grammars, this is because 2/3 of the recursive definitions of the <e> symbol use a function of the set (+, -, *), whereas 1/3 use division; this results in a biased sampling of $2/3 \div 3 = 2/9$ for each of (+, -, *), and 3/9 for division. These biases are held with the reduction of non-terminal symbols used in G3.

Other biases are seen in this figure. G0-G2 and G6, due to their larger number of non-terminal symbols, generate smaller expressions than the other grammars, when provided with the same number of genes (RND) or expansions (PTC2); this explains their slightly smaller frequency of functions, and higher frequency of terms (1.0 and x). In any case, all grammars exhibit a larger proportion of terms in their phenotypes, when using RND; this is due to the very high probability

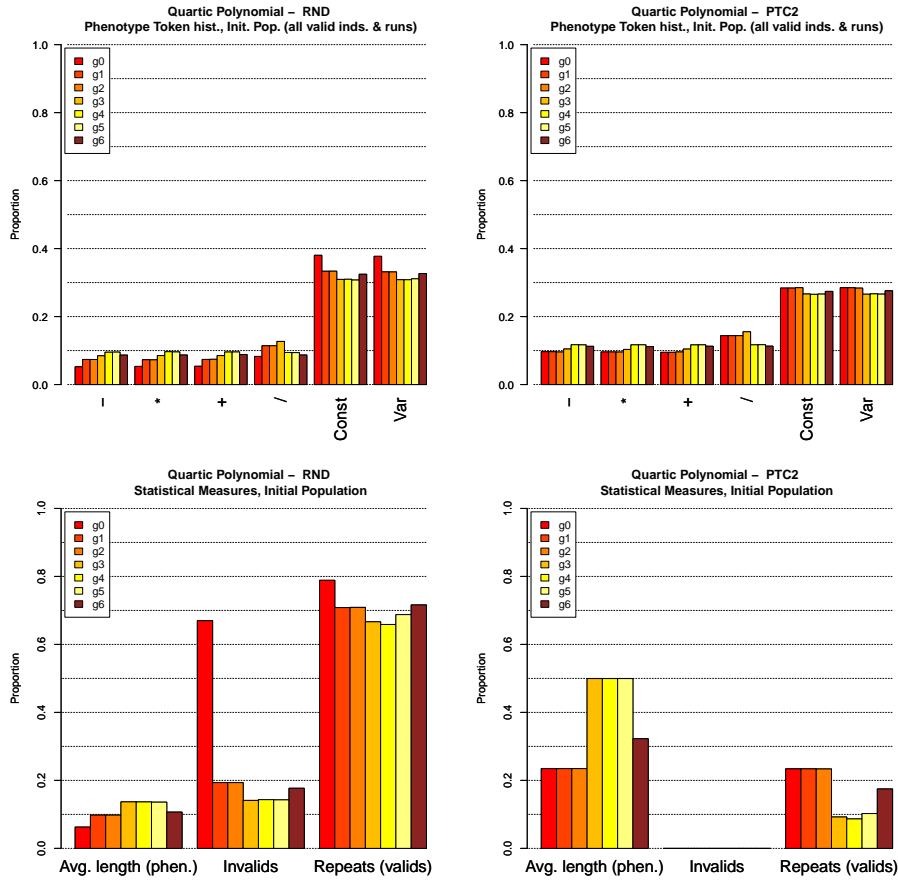


Fig. 1: Symbol frequency proportions (top) in the initial populations (measured over all successfully mapped solutions), and mapping process related statistics (bottom), using RND (left) and PTC2 (right). Results obtained from 100 independent runs.

of transforming an $\langle e \rangle$ symbol into a term ($1/4$ for G_0 , and $1/2$ for all the other grammars), meaning many solutions will consist of a single term. The use of PTC2 substantially reduces the appearance of such solutions.

The bottom half of the figure shows statistics related to the mapping process. The average length is a ratio of the average number of functions/terms in the generated phenotypes (e.g. $1.0 + x$ and $(1.0 + x)$ both have length 3) over 31. RND creates much shorter solutions than PTC2, due to the 50% probability of creating very short solutions. As for G_0 , although it is biased towards larger solutions, it usually cannot create such solutions, due to its reduced probability of terminating the mapping process. In any case, G_0 - G_2 and G_6 create shorter phenotypes, due to the intermediate symbols used during the mapping process.

When using RND, the proportion of non-mapping individuals is alarmingly high for explosive G_0 . Although far lower, G_1 , G_2 and G_6 also have some difficulty to map integer strings, due to the number of derivation steps required to terminate

the mapping process (a problem shared with G0). G3-G5 map around 85% of the random integer strings. Naturally, all phenotypes generated by PTC2 are valid.

Within the successfully-mapped individuals, there is a very high proportion of repeated solutions with RND, again a consequence of the high probability of generating single term solutions. This is slightly higher for G0, due to the increasing difficulty in successfully mapping long individuals. PTC2, ramping by size, generates longer solutions, and thus generates less repetition.

4.2 Random Walk Biases

This second experiment tests the effect of recombination and mutation on the initial populations seen in the previous experiment. To this end, 50 generations were performed in each run, but all individuals (mapping or not) were assigned the same flat fitness. This effectively means random walks are performed, biased both by grammar design, and by the genetic operators used (linear crossover and integer mutation). Fig. 2 shows the results obtained, at the last generation.

The top half of the figure shows no major differences between RND and PTC2. There is an even stronger bias towards using + and 1.0, due to shorter solutions being generated. However, the biases introduced by grammar design are still quite present, particularly towards division, in G0-G3.

The bottom plot shows that grammar design affects the dissemination of illegal individuals by a large amount: by generation 50, over 90% of individuals in runs using G0 are invalid, due to its non-terminating bias. G1, G2 and G6 exhibit $\approx 70\%$ invalids, suffering from their complex, many non-terminal symbols mapping (a problem shared with G0). G3-G5 have an expected 50% chance of generating mapping individuals under random search. This is the case using both RND and PTC2, which shows that this bias exists irrespective of the starting population.

Finally, there is a huge amount of repetition in the (valid) solutions generated. The majority are short, single term solutions, as they are more likely to survive unharmed (and unchanged) genetic operators applied with no fitness pressure.

4.3 Termination Biases

Mapping termination has always been a hotly debated topic in GE. This third experiment investigates how it is influenced by grammar design. The experimental setup is the same as in the previous sub-section, except that non-mapping individuals are assigned a very bad fitness score (the usual way of dealing with non-mapping solutions in GE [30]), whereas all others are assigned a fixed (good) fitness. The results obtained are shown in Fig. 3.

The top half of the figure shows once again that, irrespectively of using RND or PTC2, symbol biases are practically the same. The previously analysed bias towards division is still visible in G0-G3. Particularly worrying is the high bias of 1.0 over x in G0 and G1. This is a direct reflection of the propagation and recombination of mapping individuals, in combination with linked grammar productions. In G0, a large amount of codons with a value $c_i \% 4 = 3$ are required, to choose the recursion stopping production “<e> becomes <v>”; but those same codons, when interpreted under the context of <v>, will choose the term 1.0. The

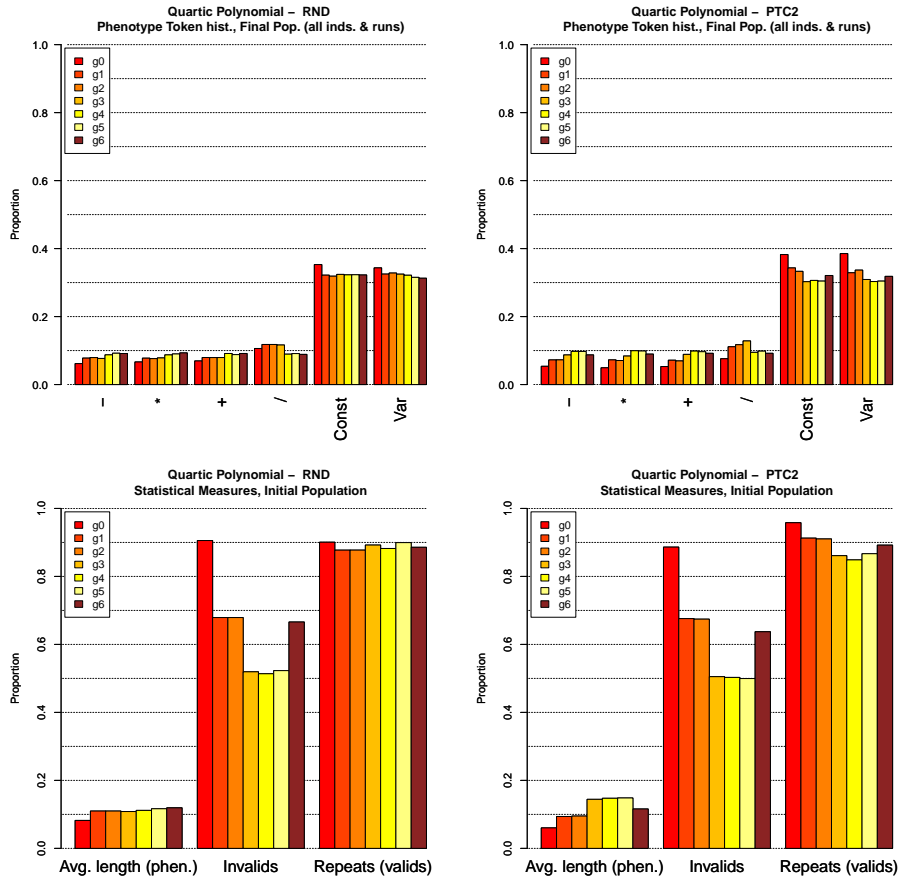


Fig. 2: Symbol frequency proportions (top) in the final populations (measured over all successfully mapped solutions), and mapping process related statistics (bottom), using RND (left) and PTC2 (right). Results obtained from 100 independent runs.

same effect is seen in G1: odd codon values are required to stop recursion of the `<e>` symbol, but these will choose the symbol `1.0` over the symbol `x` (see Table 1). These problems are removed by unlinking production choices (grammars G2-G6).

The use of bad fitness for non-mapping individuals is effective at reducing the number of illegal solutions, without requiring the propagation of extremely small solutions. There is still a large proportion of repetition, driven in this case by genetic drift: similar individuals undergoing crossover are less likely to produce illegal offspring, leading to increased repetition in the final population. Finally, note the smaller difference between RND and PTC2 at the end of these runs.

4.4 Performance Biases

To test how these findings affect performance, a series of symbolic regression experiments were ran with all seven grammars, to solve the following problems:

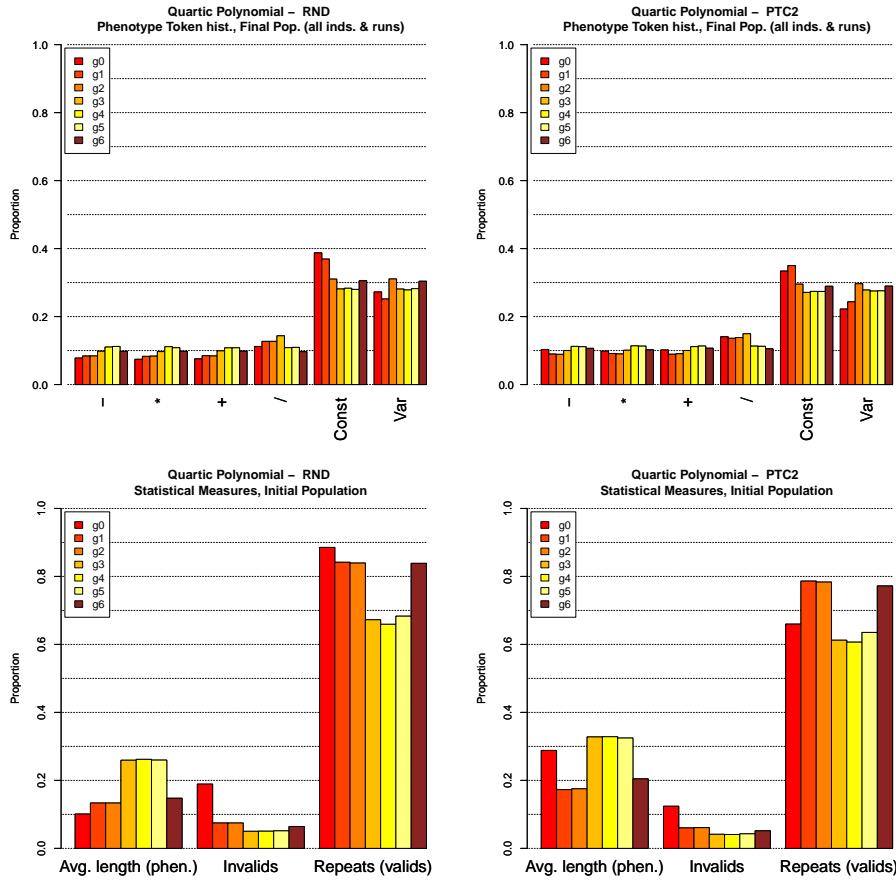


Fig. 3: Symbol frequency proportions (top) in the final populations (measured over all successfully mapped solutions), and mapping process related statistics (bottom), using RND (left) and PTC2 (right). Results obtained from 100 independent runs. Non-mapping solutions were assigned the worst possible fitness.

1. Quartic Polynomial: $x^4 + x^3 + x^2 + x$;
2. Sextic Polynomial: $x^6 + x^5 + \dots + x^2 + x$;
3. Octic Polynomial: $x^8 + x^7 + \dots + x^2 + x$;
4. Dectic Polynomial: $x^{10} + x^8 + \dots + x^2 + x$.

These problems are easy to solve, with a controlled degree of difficulty, and can be correlated to the effects of symbol biases. Note that fully configured GE runs were performed; this included the use of tails [27] at a 50% ratio at initialisation, for better mapping termination. Fig. 4 shows the measured biases and statistics for the quartic and dectic polynomial, using PTC2 (RND results were similar).

There is a clear bias towards the use of multiplication (mostly), addition, and x , the symbols required to solve the problem. The smaller uses of division, subtraction, and 1.0 are mostly due to non-effective code (bloat). It is interesting to observe that a bias towards division is still found when using G3. Solution length is short, a reflection of GE's mapping process and also the easy nature of the

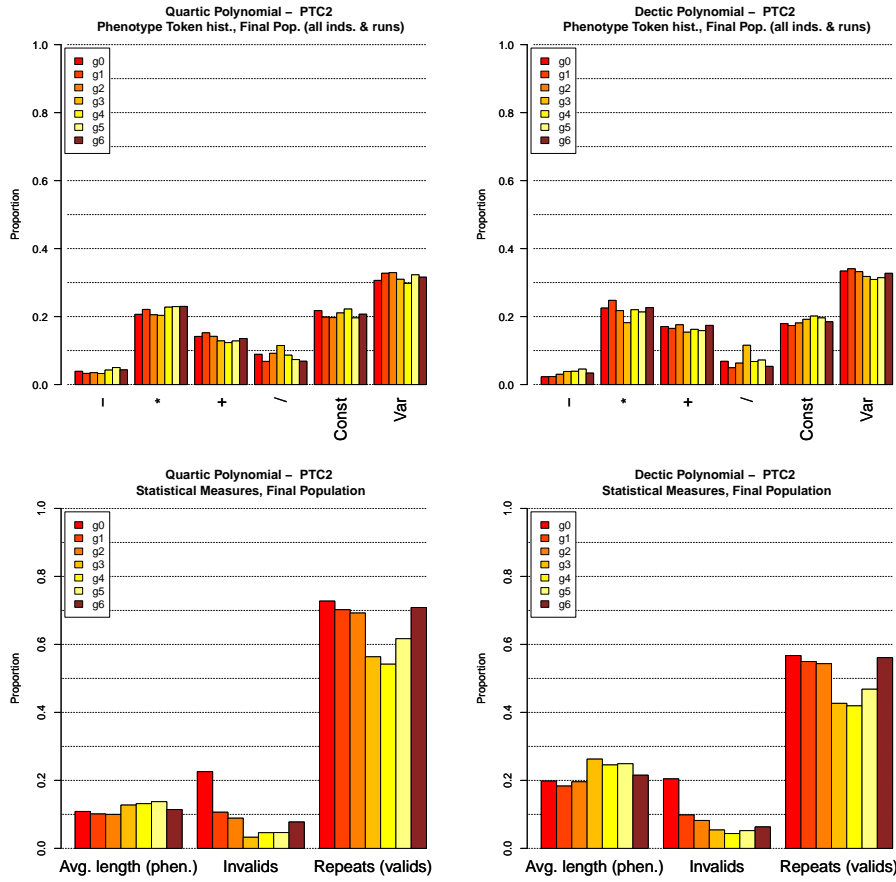


Fig. 4: Symbol frequency proportions (top) in the final populations (measured over all successfully mapped solutions), and mapping process related statistics (bottom), using PTC2, on Quartic (left) and Dectic (right) polynomial. Results obtained from 100 independent runs.

problems, with the downside of still a large amount of repetition in the last generation (due both to genetic drift and to smaller solutions being generated). Finally, experiments using G0 still generate a large amount of illegal solutions.

Fig. 5 plots the results obtained. Most configurations solve the easier problem on every run, but as the polynomial degree increases, performance slowly worsens. Results using PTC2 are better than those using RND, as expected [23].

The graphs also illustrate how grammar design can affect (or not) the performance of GE. The most obvious observation is that the reduction of grammar complexity and the associated termination biases can vastly improve performance: setups using G0-G2 are consistently worse than all other setups. Also interesting is how the bias of G2 and G3 towards the use of division has almost no effect on performance. This is because division can be used almost as effectively as multiplication to increase the degree of the polynomial ($x \times x$ versus $\frac{x}{1/x}$).

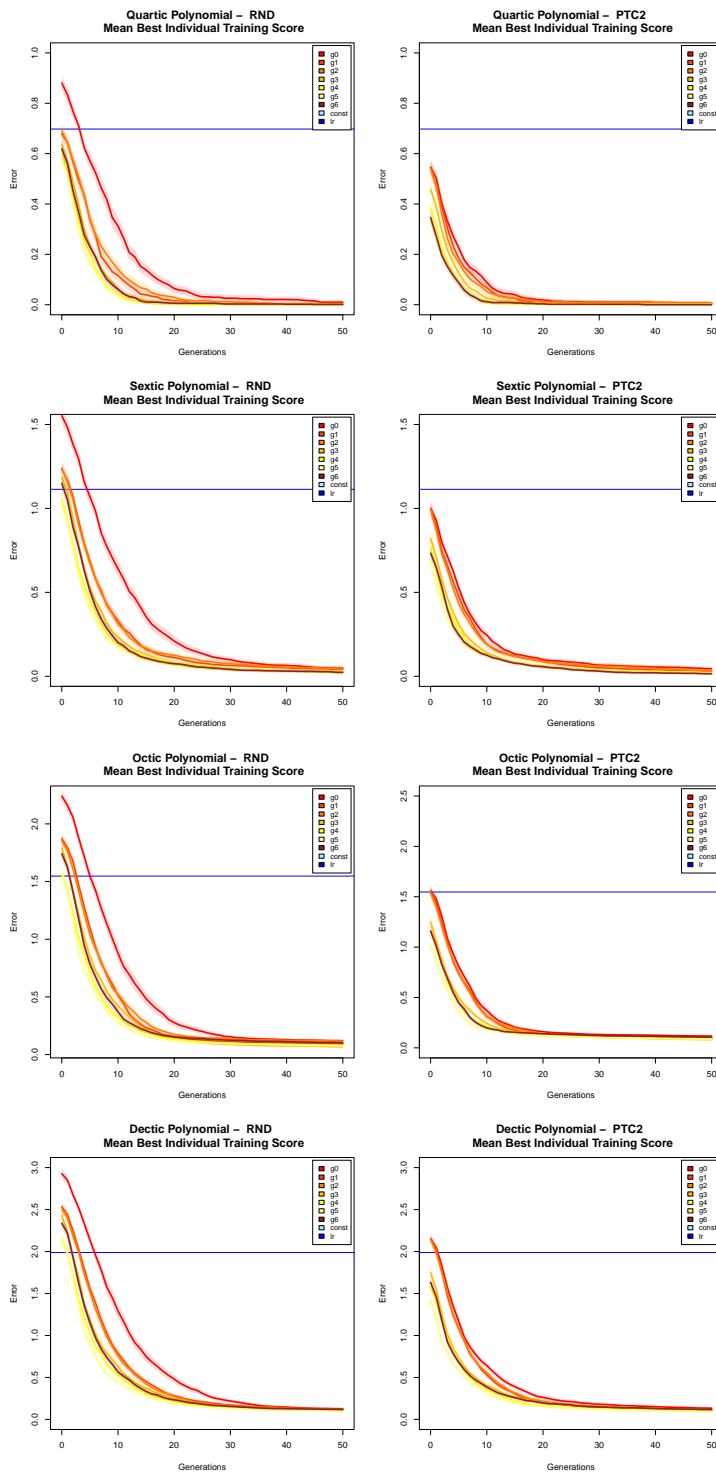


Fig. 5: Mean best individual fitness for the polynomial regression experiments, using RND (left) and PTC2 (right), averaged across 100 independent runs; error is measured as RMSE. Shades indicate 95% confidence intervals about the mean. Cyan and blue lines (where visible) show the RMSE of a constant predictor (mean of the train response variable) and a linear regression model of the training set, respectively.

Finally, the relative differences between grammars are consistent across the two initialisers, apart from G0: the lack of invalid solutions in the initial population of PTC2, along with its larger initial solution size (see Fig. 1), substantially improve the final performance of runs using G0.

5 Performance Analysis

5.1 Problems

To measure the effect of grammar design on the final outcome of GE’s evolutionary process, a series of experiments were ran across several problem types. These include regression, classification and design problems, across several application domains and difficulty ranges. Table 3 lists all the problems attempted.

5.2 Grammars

For every problem attempted, seven grammar versions were designed, G0-G6, as detailed in Section 3. Each of the grammars respects as much as possible the original function and terminal sets as defined in their original publications; this is important for benchmarking purposes, as attempting to solve a problem using different functions can severely alter the difficulty of the benchmarks [24]. For all problems, constants were created using digit concatenation, with 100 possible values within the range [0.0...9.9], and a step of 0.1.

Each of the problems was attempted using both RND and PTC2, over 100 independent runs. To illustrate the grammar transformation process, and its impact on the search effectiveness of GE, three problems are examined in detail: Keijzer-6, Vladislavleva-4, and Shape Match (hard).

5.2.1 Keijzer-6

The Keijzer-6 (K6) problem [12], also known as the Harmonic function, is a single-variable problem, using only addition, multiplication and three unary functions. A typical GE grammar defining possible solutions for this problem is shown in Table 4, cell G0. The symbol $\langle e \rangle$ has more productions creating new $\langle e \rangle$ symbols rather than mapping them to something else, so the G1 version addresses this, by having a “ $|\langle v \rangle$ ” production for each production replacing one $\langle e \rangle$ symbol with two. There are two sets of non-terminal symbols with linked productions: $\langle o \rangle$ with $\langle v \rangle$, and $\langle e \rangle$ with $\langle d \rangle$. This is addressed in G2, through repetition of productions.

G3 reduces the number of non-terminal symbols to a single one, while maintaining the biases of G1 and G2. This results in a much larger grammar, with 3000 production rules. Grammar G4 addresses the slightly higher bias towards use of the operators $+$ and $*$ over the functions *inv*, *neg* and *sqrt*.

Grammar G5 is a conversion to prefix notation, which considerably reduces the complexity of the grammar, by removing bracketed versions of the $+$ and $*$ operators. Finally, G6 separates the definition of functions and terminals into two symbols ($\langle e \rangle$ and $\langle v \rangle$), resulting in a very compact grammar, but without the extra complexity of grammars G0-G2.

Table 3: Benchmark problems; if specified, $E[a, b, c]$ means a grid of points evenly spaced with an interval of c , from a to b inclusive, whereas $U[a, b, c]$ means c uniform random samples drawn from a to b inclusive; the specified type is regression (R), classification (C#) (# = number of classes), or image matching (IM)

Name	Vars.	Data Source	Type	Training Set Test Set
Keijzer-6 (K6) [12]	1	$y = \sum_i^x \frac{1}{i}$	R	$E[1, 50, 1]$ $E[1, 120, 1]$
Pagie-1 (P1) [34]	2	$y = \frac{1}{1+x_1^{-4}} + \frac{1}{1+x_2^{-4}}$	R	$E[-5, 5, .4]$ $E[-5, 5, .1]$
Vladislavleva-4 (V4) [41]	5	$y = \frac{10}{5 + \sum_{i=1}^5 (x_i - 3)^2}$	R	$U[.05, 6.05, 1024]$ $U[-.25, 6.35, 5k]$
Tower [41]	5	Gas chromatography data	R	4721 points 278 points
Korns-12 (K12) [14]	5	$y = 2 - 2.1\cos(9.8x_1)\sin(1.3x_5)$	R	$U[-50, 50, 10k]$ $U[-50, 50, 10k]$
Forest [3]	12	Forest Fires data	R	414 points 103 points
Housing [16]	13	Housing values	R	354 points 152 points
Dow Chemical (Dow) *	57	Chemical process data	R	747 points 319 points
Even-parity 5 (EPar5) [15]	5	Parity of boolean inputs	C2	32 points -
Multiplexer 11 (Mux11) [15]	11	Boolean multiplexer	C2	2048 points -
Breast Cancer Wisconsin [16]	9	Diagnostic data	C2	400 points 283 points
Wine Quality: red [2]	11	Physicochemical test data	C10	1000 points 599 points
Wine Quality: white [2]	11	Physicochemical test data	C10	3000 points 1898 points
Shape Match (Easy) [33]	-	Generated 250x250 shape	IM	- -
Shape Match (Medium) [33]	-	Generated 250x250 shape	IM	- -
Shape Match (Hard) [33]	-	Generated 250x250 shape	IM	- -

* Source: <http://gpbenchmarks.org/symbolicregressioncompetition/>

5.2.2 Vladislavleva-4

The Vladislavleva-4 (V4) problem [41], also known as the UBall5D function, is a five-variable problem, and its function set, as defined in its original publication, is extensive:

- Functions: $+$, $-$, $*$, $/$, $square$, x^{real} , $x + real$, $x \cdot real$
- Terminals: x_0, x_1, x_2, x_3, x_4

This results in a complex base grammar, shown in Table 5, cell G0. The definition of the $\langle e \rangle$ symbol has a heavy bias towards growth, which is addressed in G1. It also has a more complex set of linked production, between the symbols $\langle a \rangle$ (5 productions), $\langle v \rangle$ (5 productions) and $\langle d \rangle$ (10 productions); G2 addresses this by introducing 10 copies of each original production associated with $\langle a \rangle$, and 50 copies of each original production associated with $\langle v \rangle$.

Table 5: Grammars used for the Vladislavleva-4 (V4) problem, versions G0-G6; a set of productions followed by $\{n\}$ means they are repeated n times, whereas the notation $\alpha|\dots|\delta$ is shorthand for a production for each of the elements in the sequence $[\alpha..\delta]$; constants in G6 are created using `GECodeonValue` [25]

G0	G1	G2
<pre> <s> ::= <e> <e> ::= <e> <o> <e> (<e> <o> <e>) <f1> (<e>, <e>) <f2> (<e>, 2) <a> <o> ::= + - * <f1> ::= pdiv <f2> ::= pow <a> ::= pow(<v>, <d>.<d>.<d>) (<v> + <d>.<d>.<d>) (<v> * <d>.<d>.<d>) <v> <d>.<d> <v> ::= x0 x1 x2 x3 x4 <d> ::= 0 1 2 3 4 5 6 7 8 9 </pre>	<pre> <s> ::= <e> <e> ::= <e> <o> <e> <a> (<e> <o> <e>) <a> <f1> (<e>, <e>) <a> <f2> (<e>, 2) <o> ::= + - * <f1> ::= pdiv <f2> ::= pow <a> ::= pow(<v>, <d>.<d>.<d>) (<v> + <d>.<d>.<d>) (<v> * <d>.<d>.<d>) <v> <d>.<d> <v> ::= x0 x1 x2 x3 x4 <d> ::= 0 1 2 3 4 5 6 7 8 9 </pre>	<pre> <s> ::= <e> <e> ::= <e> <o> <e> <a> (<e> <o> <e>) <a> <f1> (<e>, <e>) <a> <f2> (<e>, 2) <o> ::= + - * <f1> ::= pdiv <f2> ::= pow <a> ::= pow(<v>, <d>.<d>.<d>){10} (<v> + <d>.<d>.<d>){10} (<v> * <d>.<d>.<d>){10} <v> {10} <d>.<d> {10} <v> ::= x0 x1 x2 x3 x4{50} <d> ::= 0 1 2 3 4 5 6 7 8 9 </pre>
<pre> <e> ::= <e> + <e> {833} <e> - <e> {833} <e> * <e> {833} (<e> + <e>) {833} (<e> - <e>) {833} (<e> * <e>) {833} pdiv(<e>, <e>) {2500} pow(<e>, 2) {2500} pow(x0, 0.0) ... pow(x0, 9.9) {3} pow(x1, 0.0) ... pow(x1, 9.9) {3} pow(x2, 0.0) ... pow(x2, 9.9) {3} pow(x3, 0.0) ... pow(x3, 9.9) {3} pow(x4, 0.0) ... pow(x4, 9.9) {3} (x0 + 0.0) ... (x0 + 9.9) {3} (x1 + 0.0) ... (x1 + 9.9) {3} (x2 + 0.0) ... (x2 + 9.9) {3} (x3 + 0.0) ... (x3 + 9.9) {3} (x4 + 0.0) ... (x4 + 9.9) {3} (x0 * 0.0) ... (x0 * 9.9) {3} (x1 * 0.0) ... (x1 * 9.9) {3} (x2 * 0.0) ... (x2 * 9.9) {3} (x3 * 0.0) ... (x3 * 9.9) {3} (x4 * 0.0) ... (x4 * 9.9) {3} x0 x1 x2 x3 x4 {300} 0.0 ... 9.9 {15} </pre>	<pre> <e> ::= <e> + <e> {625} <e> - <e> {625} <e> * <e> {625} <e> pdiv <e> {625} (<e> + <e>) {625} (<e> - <e>) {625} (<e> * <e>) {625} (<e> pdiv <e>) {625} pow(<e>, 2) {2500} pow(x0, 0.0) ... pow(x0, 9.9) {2} pow(x1, 0.0) ... pow(x1, 9.9) {2} pow(x2, 0.0) ... pow(x2, 9.9) {2} pow(x3, 0.0) ... pow(x3, 9.9) {2} pow(x4, 0.0) ... pow(x4, 9.9) {2} (x0 + 0.0) ... (x0 + 9.9) {2} (x1 + 0.0) ... (x1 + 9.9) {2} (x2 + 0.0) ... (x2 + 9.9) {2} (x3 + 0.0) ... (x3 + 9.9) {2} (x4 + 0.0) ... (x4 + 9.9) {2} (x0 * 0.0) ... (x0 * 9.9) {2} (x1 * 0.0) ... (x1 * 9.9) {2} (x2 * 0.0) ... (x2 * 9.9) {2} (x3 * 0.0) ... (x3 * 9.9) {2} (x4 * 0.0) ... (x4 * 9.9) {2} x0 x1 x2 x3 x4 {200} 0.0 ... 9.9 {10} </pre>	<pre> <e> ::= + <e> <e> {625} - <e> <e> {625} * <e> <e> {625} pdiv <e> <e> {625} pow <e> 2 {625} pow x0 0.0 ... pow x0 9.9 pow x1 0.0 ... pow x1 9.9 pow x2 0.0 ... pow x2 9.9 pow x3 0.0 ... pow x3 9.9 pow x4 0.0 ... pow x4 9.9 + x0 0.0 ... + x0 9.9 + x1 0.0 ... + x1 9.9 + x2 0.0 ... + x2 9.9 + x3 0.0 ... + x3 9.9 + x4 0.0 ... + x4 9.9 * x0 0.0 ... * x0 9.9 * x1 0.0 ... * x1 9.9 * x2 0.0 ... * x2 9.9 * x3 0.0 ... * x3 9.9 * x4 0.0 ... * x4 9.9 x0 x1 x2 x3 x4 {100} 0.0 ... 9.9 {5} </pre>
<pre> G6 <e> ::= + <e> <e> - <e> <e> * <e> <e> pdiv <e> <e> pow <e> 2 pow <v> <GECodonValue{0.0 : 9.9 : 0.1}> + <v> <GECodonValue{0.0 : 9.9 : 0.1}> * <v> <GECodonValue{0.0 : 9.9 : 0.1}> <v> <v> <GECodonValue{0.0 : 9.9 : 0.1}> <GECodonValue{0.0 : 9.9 : 0.1}> <v> ::= x0 x1 x2 x3 x4 </pre>		

pen's state (position where it will draw next) onto and off the stack. Finally, `sqr` draws a square, and `crcl` draws a circle.

Table 6 shows the original grammar [33] as G0, with a call to a Python interpreter. It has no recursively defined binary operators: a single call to `<p> ::= <e>` stops recursion of symbol `<p>`, and likewise, a single call to `<e> ::= <v>` terminates the recursion of `<e>`. As such, G1 is identical to G0. There are, however, two sets of non-terminal symbols with linked productions: `<p>` with `<v>`, and `<e>` with `<o>`. G2 addresses this, through the explained approach of repetition of productions.

G3 reduces the number of non-terminal symbols to three. Although the first symbol (`<s>`) is of minor importance (it has a single associated production, so no codon is used, and it appears only once, at the start), the recursive nature

Table 6: Grammars used for the Shape Match (hard) problem, versions G0-G6; a set of productions followed by $\{n\}$ means they are repeated n times, whereas the notation $\alpha|\dots|\delta$ is shorthand for a production for each of the elements in the sequence $[\alpha..\delta]$; constants in G6 are created using `GECodeonValue` [25]

G0	G1 (same as G0)	G2
$\langle s \rangle ::= \text{python hard.py } \langle p \rangle$ $\langle p \rangle ::= \langle e \rangle$ $\langle e \rangle ::= \langle v \rangle$ $\langle o \rangle ::= s_0$ $\langle v \rangle ::= \text{sqr}$ $\langle crcl \rangle ::= \text{crcl}$	$\langle s \rangle ::= \text{python hard.py } \langle p \rangle$ $\langle p \rangle ::= \langle e \rangle$ $\langle e \rangle ::= \langle v \rangle$ $\langle o \rangle ::= s_0$ $\langle v \rangle ::= \text{sqr}$ $\langle crcl \rangle ::= \text{crcl}$	$\langle s \rangle ::= \text{python hard.py } \langle p \rangle$ $\langle p \rangle ::= \langle e \rangle$ $\langle e \rangle ::= \langle v \rangle$ $\langle o \rangle ::= s_0 s_0 s_0$ $\langle v \rangle ::= \text{sqr} \text{sqr}$ $\langle crcl \rangle ::= \text{crcl} \text{crcl}$
G3	G4	G5
$\langle s \rangle ::= \text{python hard.py } \langle p \rangle$ $\langle p \rangle ::= \langle e \rangle \{18\}$ $\langle e \rangle ::= \text{sqr} \{3\}$ $\langle crcl \rangle ::= \text{crcl} \{3\}$ $s_0 \langle e \rangle$ $s_1 \langle e \rangle$ $s_2 \langle e \rangle$ $s_3 \langle e \rangle$ $gro \langle e \rangle$ $shrnk \langle e \rangle$ $\langle e \rangle \{6\}$	$\langle s \rangle ::= \text{python hard.py } \langle p \rangle$ $\langle p \rangle ::= \langle e \rangle \{42\}$ $\langle e \rangle ::= \text{sqr} \{7\}$ $\langle crcl \rangle ::= \text{crcl} \{7\}$ $s_0 \langle e \rangle \{4\}$ $s_1 \langle e \rangle \{4\}$ $s_2 \langle e \rangle \{4\}$ $s_3 \langle e \rangle \{4\}$ $gro \langle e \rangle \{4\}$ $shrnk \langle e \rangle \{4\}$ $\langle e \rangle \{4\}$	$\langle s \rangle ::= \text{python hard.py } \langle p \rangle$ $\langle p \rangle ::= \text{sqr} \{7\}$ $\langle crcl \rangle ::= \text{crcl} \{7\}$ $s_0 \langle p \rangle \{2\}$ $s_1 \langle p \rangle \{2\}$ $s_2 \langle p \rangle \{2\}$ $s_3 \langle p \rangle \{2\}$ $gro \langle p \rangle \{2\}$ $shrnk \langle p \rangle \{2\}$ $pushed \langle p \rangle \{2\}$ $\text{sqr } \langle p \rangle \{7\}$ $\text{crcl } \langle p \rangle \{7\}$
G6		
$\langle s \rangle ::= \text{python hard.py } \langle p \rangle$ $\langle p \rangle ::= \langle v \rangle \langle t \rangle \langle p \rangle \langle v \rangle \langle p \rangle$ $\langle o \rangle ::= s_0 s_1 s_2 s_3 gro shrnk pushed$ $\langle v \rangle ::= \text{sqr} \text{crcl}$		

of symbols $\langle p \rangle$ and $\langle e \rangle$ in G0-G2 requires the presence of both. Also, after the incorporation of symbols $\langle o \rangle$ and $\langle v \rangle$ into $\langle e \rangle$, the resulting number of productions is even, which links them with those associated with $\langle p \rangle$, so the explained unlinking process was employed. G4 is similar, but reduces the bias of the \square operator to that of all other transformations.

G5 replaces \square with a prefix equivalent operator, *pushed*. This allows the definition of what is essentially a single non-terminal symbol grammar (not counting the $\langle s \rangle$ symbol, as explained above). The exact biases of G4 are impossible to keep, so a compromise was chosen. Finally, G6 uses three symbols, $\langle p \rangle$, $\langle o \rangle$ and $\langle v \rangle$, to create a compact and highly readable grammar.

5.3 Experimental Setup

The experimental setup used was the same as in Table 2, but using 50% non-coding tails at initialisation [27]. Also, as seen in Section 4.1, different grammars will generate different solution sizes at initialisation, with direct influence in their initial fitnesses. In order to properly analyse the effect of grammar design in the search capability of GE, experiments using different grammars were setup such that they generate similarly sized phenotype solutions at initialisation. As such, RND and PTC2 were setup as shown in Table 7. Finally, protected versions of some operators were used, such as division (1.0 if divisor $< 1e - 5$), inversion (1.0 if argument $< 1e - 5$), and square root ($\sqrt{|x|}$).

Table 7: Genome length (for RND) and min/max derivation steps (for PTC2) used during initialisation, for the three problems analysed

Problem	Parameter	g0	g1	g2	g3	g4	g5	g6
Keijzer6	Genome Length	30	30	30	15	15	15	23
	Min. Der. Steps	3	3	3	1	1	1	2
	Max. Der. Steps	31	31	31	15	15	15	23
Vladislavleva4	Genome Length	54	54	54	15	15	15	31
	Min. Der. Steps	4	4	4	1	1	1	2
	Max. Der. Steps	55	55	55	15	15	15	31
Shape Match (hard)	Genome Length	24	24	24	14	14	14	23
	Min. Der. Steps	4	4	4	3	3	2	3
	Max. Der. Steps	25	25	25	15	15	16	24

5.4 Results

Fig. 6 plots the mean best train results for the K6, V4 and Shape Match (hard) problem, using RND and PTC2 initialisation³. All grammar variants find similarly good solutions for the K6 experiment at the 50th generation, using PTC2, but much larger differences are found between the results obtained with grammars G0-G2 versus those with grammars G3-G6, when using RND. A similar result is observed for the V4 experiment, except that in this case, the difference between G0-G2 and G3-G6 is more evident when using PTC2.

Finally, the Shape Match (hard) experiments again show a large difference between both clusters of grammars, along with the positive effect in convergence of removing the high bias towards the \square operator in G4-G6 (applying it more than once has no practical effect). Although there is a very marked change in the speed of convergence towards an optimal solution, particularly with the reduction of non-terminal symbols, eventually all grammar solutions find similarly good solutions.

5.4.1 Significance and Test Results

In order to quantify the effect of each grammar design on the search effectiveness of GE, two-sample Mann-Whitney U-tests were calculated for final median best fit results, for all grammars. The results are shown in Table 8.

These are similar to what was observed in the K6, V4 and Shape Match problems. Overall, runs using G3-G5 are significantly better than all others. Small differences between G3-G5 are mostly problem domain specific. G0-G2 are often significantly better only amongst themselves, with the exception of a few noisy real-world datasets, such as Tower and Wine Quality. G6 seems to provide a compromise between performance and complexity of grammar. As before, G0 does particularly bad with RND initialisation, and differences between different grammars are less evident when using PTC2.

Regarding test performance, no validation or early-stopping approaches were employed in these experiments, so it is not unreasonable to expect some level of overfitting, particularly from approaches with very good training performance. The statistical significance results are shown in Table 9.

³ As the focus of this study is on grammar design, no regression performance improving techniques such as linear scaling [12] or cross-validation were used.

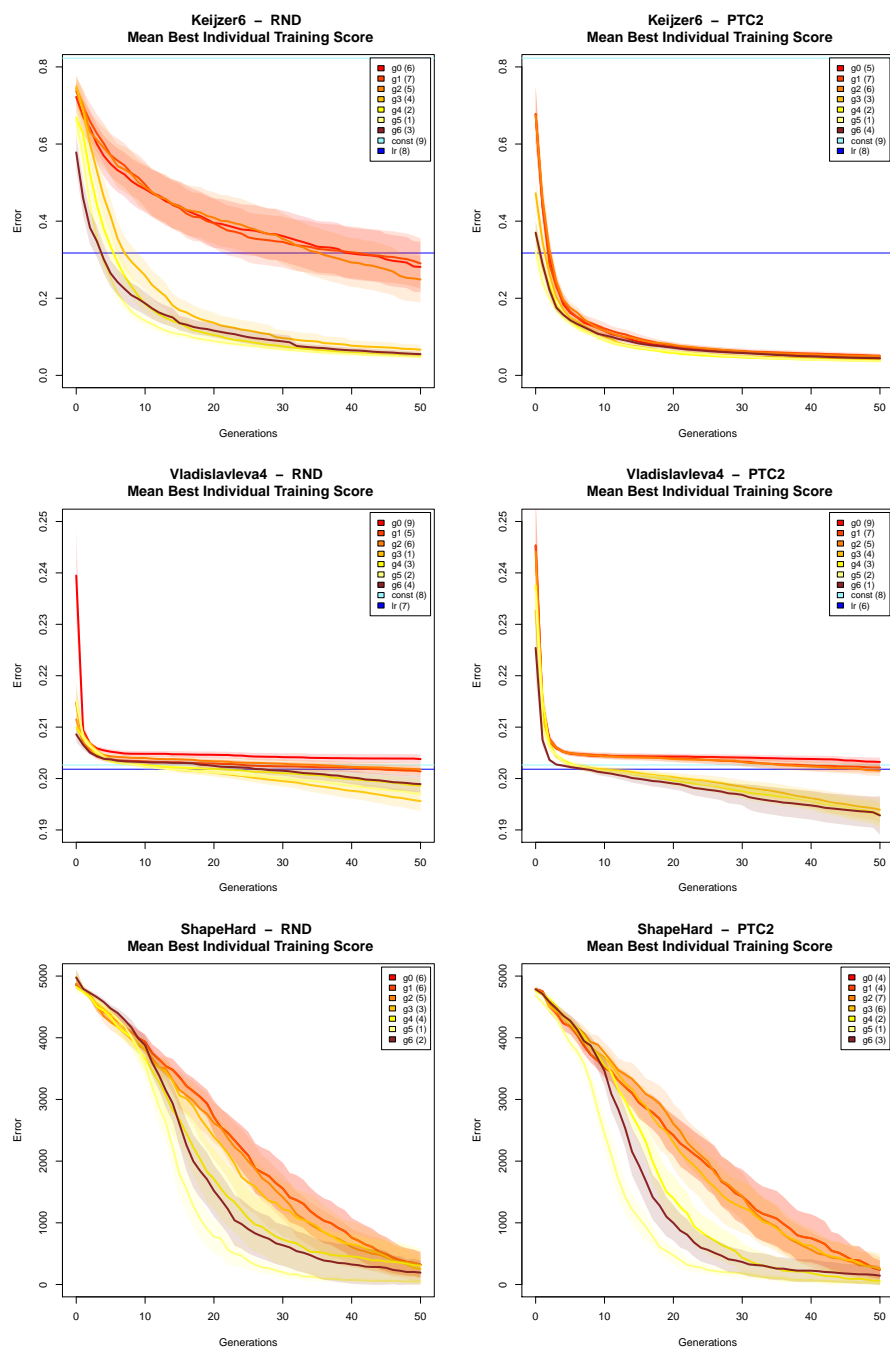


Fig. 6: Mean best train scores for the K6 (top), V4 (middle) and Shape match (hard) (bottom) experiments, using RND (left) or PTC2 (right), averaged across 100 independent runs; error is measured as the Root Mean Squared Error (RMSE) for K6 and V4, or the number of mismatched pixels for shape matching. Shades indicate 95% confidence intervals about the mean. Cyan and blue lines (where visible) indicate the RMSE of a constant predictor (mean of the train response variable) and a linear regression model of the training set, respectively.

Table 8: Mann-Whitney U-tests of final median best training fit results, for all problems attempted; each number is a count of all grammars against which a significantly differently better performance was measured (results across 100 independent runs)

	G0	G1	G2	G3	G4	G5	G6	
RND	BreastCancerW	0	0	0	5	3	0	0
	Dow	0	0	0	3	3	3	3
	EvenParity5	0	0	0	4	4	4	0
	Forest	0	1	1	4	4	5	1
	Housing	0	1	1	3	3	3	3
	Keijzer6	0	0	0	3	4	4	3
	Korns12	0	0	0	3	4	6	4
	Multiplexer11	0	1	1	3	3	1	4
	Pagie1	0	0	0	4	4	4	0
	Tower	0	2	2	4	5	2	1
	Vladislavleva4	0	0	0	4	0	4	3
	WineQualityRed	0	1	1	5	5	2	2
	WineQualityWhite	0	1	1	3	3	3	3
	ShapeEasy	0	0	0	0	0	0	0
	ShapeMedium	0	0	0	0	0	0	0
	ShapeHard	0	0	0	0	0	1	0
	Sum	0	7	7	48	45	42	27
Mean	0	0.438	0.438	3	2.812	2.625	1.688	
Median	0	0	0	3	3	3	1.5	
Std. Dev.	0	0.6292	0.6292	1.633	1.797	1.857	1.580	
PTC2	BreastCancerW	0	0	0	4	4	1	0
	Dow	0	0	0	2	0	2	3
	EvenParity5	0	0	0	1	1	4	1
	Forest	0	0	0	2	4	1	0
	Housing	0	0	0	3	2	2	2
	Keijzer6	0	0	0	3	4	4	0
	Korns12	0	0	0	3	4	5	6
	Multiplexer11	0	1	0	0	0	0	4
	Pagie1	0	0	0	4	4	4	2
	Tower	0	0	1	1	3	1	1
	Vladislavleva4	0	0	0	3	3	3	3
	WineQualityRed	1	1	0	2	2	0	0
	WineQualityWhite	1	1	3	2	1	0	1
	ShapeEasy	0	0	0	0	0	0	0
	ShapeMedium	0	0	0	0	0	0	0
	ShapeHard	0	0	0	0	0	0	0
	Sum	2	3	4	30	32	27	23
Mean	0.125	0.188	0.25	1.875	2	1.688	1.438	
Median	0	0	0	2	2	1	1	
Std. Dev.	0.342	0.403	0.774	1.408	1.713	1.778	1.788	

As expected, the test results are far less clear cut. Results are still better when using G3-G5, but are very problem-dependent. This mostly results from problems such as K12: it is such a hard problem to solve using the original function and terminal set [14] (which uses no trigonometric functions), that any small improvement in training performance invariably led to a degradation in test performance.

Table 9: Mann-Whitney U-tests of final median best test fit results (of best training solutions), for all problems with a test set; each number is a count of all grammars against which a significantly differently better performance was measured (results across 100 independent runs)

		G0	G1	G2	G3	G4	G5	G6
RND	BreastCancerW	0	0	0	0	0	1	0
	Dow	0	0	0	4	2	2	0
	Forest	6	2	1	0	0	0	3
	Housing	0	1	1	3	3	3	3
	Keijzer6	0	0	0	2	4	4	3
	Korns12	3	3	6	3	0	0	1
	Pagiel	0	0	0	4	4	4	0
	Tower	1	0	0	0	0	0	1
	Vladislavleva4	0	1	1	0	0	0	1
	WineQualityRed	0	1	1	5	5	2	2
	WineQualityWhite	5	0	0	0	0	0	0
	Sum	15	8	10	21	18	16	14
	Mean	1.364	0.727	0.909	1.909	1.636	1.455	1.273
	Median	0	0	0	2	0	1	1
Std. Dev.	2.248	1.009	1.758	1.973	2.014	1.635	1.272	
PTC2	BreastCancerW	0	0	0	0	4	4	0
	Dow	0	0	0	1	2	0	0
	Forest	1	0	0	0	0	0	1
	Housing	0	0	0	3	2	2	2
	Keijzer6	1	0	0	1	2	2	0
	Korns12	4	3	3	3	2	0	0
	Pagiel	0	0	0	4	3	4	2
	Tower	1	0	0	0	0	0	1
	Vladislavleva4	1	1	0	0	0	0	1
	WineQualityRed	0	1	0	2	1	0	0
	WineQualityWhite	0	0	0	0	0	0	0
	Sum	8	5	3	14	16	12	7
	Mean	0.727	0.455	0.273	1.273	1.455	1.091	0.636
	Median	0	0	0	1	2	0	0
Std. Dev.	1.191	0.934	0.905	1.489	1.368	1.640	0.809	

5.5 Analysis

The results obtained show just how much the design of a grammar can affect the search capability of GE. The most obvious performance improvements come from the construction of balanced grammars, and the reduction of the number of non-terminal symbols. These come at a price, however, particularly the latter: the complexity of the generated grammars can make them particularly hard to read and/or modify by hand. However, the application of these transformations is a deterministic process, meaning that they can be applied automatically: this allows manual grammar modifications to be applied to the simpler versions of the grammars, with subsequent application of automatic transformations. In any case, even a partial (manual) application of some of these modifications can be of use, as seen with the results for G6 variants.

The usefulness of some of these applications is problem-dependent. This is particularly the case of bias-related transformations, expressed in the results obtained with G2 and G4. For the problems attempted, the unlinking process employed with G2 grammars does not seem to confer any performance advantage when compared to G1 grammars. Likewise, the removal of biases towards certain operators em-

ployed in G4 grammars does not seem to confer any advantage, in the problems attempted. Finally, the switch from infix to prefix notation only makes a small difference for a few problems (for better or worse).

It is worth pointing out again that these results relate only to GE using a linear genome representation. Derivation tree based approaches make use of different search operators, and as such, the effect of grammar design is markedly different.

Another observation is the need to understand grammar design, when setting up any initialisation procedure (even something as simple as RND initialisation). The size of genotype structures required to generate certain solution (phenotype) size ranges needs to be adapted, depending on the grammar used.

6 Conclusions

Grammar design is one of the main tasks in GE: in order to write a good grammar, deep knowledge in both GE and the application domain are required. The experiments examined in this chapter, however, show that this is not necessarily always the case: there are specific grammar design principles that can be applied when attempting to solve any problem using GE, which do not require domain knowledge. Of the transformations analysed, the creation of recursion-balanced grammars and the reduction of the number of non-terminal symbols are particularly useful for improving the performance of GE, when using a linear genome representation.

Fixing symbol biases does not necessarily lead to better performance: this is completely problem dependent, and it can both degrade or improve performance. But it does affect search space exploration, as shown in the experiments conducted. This leads to two recommendations:

- When designing a GE grammar to compare its performance with systems such as GP, respect the symbol biases of the original system;
- When applying GE to a real-world problem, where symbol biases are known, use this knowledge to bias the exploration of the search space.

GE is like any other search algorithm, and in fact like any tool: it has its advantages and disadvantages, and will only provide its best performance if correctly used. There has been a recent surge of publications criticising GE's performance [42,18,21,17], some even deeming that its performance “resembles that of random search” [42]. But most of the results provided were the result of using badly designed grammars, and poor experimental setup. The analysis and results presented in this chapter aim to provide another step towards achieving the goal of *Understanding Grammatical Evolution*.

References

1. J. Byrne, M. O'Neill, J. McDermott, and A. Brabazon. An analysis of the behaviour of mutation in grammatical evolution. In A. I. E.-A. et. al, editor, *European Conference on Genetic Programming, EuroGP 2010*, volume 6021 of *LNCS*, pages 14–25. Springer, 2010.
2. P. Cortez, A. Cerdeira, F. Almeida, T. Matos, and J. Reis. Modeling wine preferences by data mining from physicochemical properties. *Decision Support Systems*, 47(4):547–553, November 2009.

3. P. Cortez and A. Morais. A data mining approach to predict forest fires using meteorological data. In J. N. et. al, editor, *Portuguese Conference on Artificial Intelligence, EPIA 2007*, pages 512–523. APPIA, December 2007.
4. D. Fagan, M. Nicolau, E. Hemberg, M. O’Neill, A. Brabazon, and S. McGarraghy. Investigation of the performance of different mapping orders for ge on the max problem. In S. S. et. al, editor, *European Conference on Genetic Programming, EuroGP 2011, Torino, Italy, April 27-29, 2011, Proceedings*, volume 6621 of *Lecture Notes in Computer Science*, pages 286–297. Springer, 2011.
5. D. Fagan, M. O’Neill, E. Galván-López, A. Brabazon, and S. McGarraghy. An analysis of genotype-phenotype maps in grammatical evolution. In A. I. E.-A. et. al, editor, *European Conference on Genetic Programming, EuroGP 2010*, volume 6021 of *LNCS*, pages 62–73. Springer, 2010.
6. S. Forstenlechner, D. Fagan, M. Nicolau, and M. O’Neill. A grammar design pattern for arbitrary program synthesis problems in genetic programming. In J. M. et. al, editor, *European Conference on Genetic Programming, EuroGP 2017*, volume 10196 of *LNCS*, pages 262–277. Springer, 2017.
7. S. Forstenlechner, M. Nicolau, D. Fagan, and M. O’Neill. Grammar design for derivation tree based genetic programming systems. In M. I. H. et. al, editor, *European Conference on Genetic Programming, EuroGP 2016*, volume 9594 of *LNCS*, pages 199–214. Springer, 2016.
8. C. Gathercole and P. Ross. An adverse interaction between crossover and restricted tree depth in genetic programming. In J. R. K. et. al, editor, *Genetic Programming 1996: First Annual Conference*, pages 291–296. MIT Press, 1996.
9. R. Harper. GE, explosive grammars and the lasting legacy of bad initialisation. In *IEEE Congress on Evolutionary Computation, CEC 2010*, pages 2602–2609, 2010.
10. E. Hemberg. *An Exploration of Grammars in Grammatical Evolution*. PhD thesis, University College Dublin, University College Dublin, Ireland, 2010.
11. E. Hemberg, N. McPhee, M. O’Neill, and A. Brabazon. Pre-, in- and postfix grammars for symbolic regression in grammatical evolution. In *IEEE Workshop and Summer School on Evolutionary Computing, 2008*, pages 18–22, 2008.
12. M. Keijzer. Improving symbolic regression with interval arithmetic and linear scaling. In C. R. et. al, editor, *European Conference on Genetic Programming, EuroGP 2003*, volume 2610 of *LNCS*, pages 70–82. Springer, 2003.
13. M. Keijzer, M. O’Neill, C. Ryan, and M. Cattolico. Grammatical evolution rules: The mod and the bucket rule. In J. A. F. et. al, editor, *European Conference on Genetic Programming, EuroGP 2002*, volume 2278 of *LNCS*, pages 123–130. Springer, 2002.
14. M. F. Korns. Accuracy in symbolic regression. In R. R. et. al, editor, *Genetic Programming Theory and Practice IX*, Genetic and Evolutionary Computation, pages 129–151. Springer, New York, 2011.
15. J. R. Koza. *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. MIT Press, Cambridge, MA, 1992.
16. M. Lichman. UCI machine learning repository. <http://archive.ics.uci.edu/ml>, 2013.
17. N. Lourenço, J. Ferrer, F. B. Pereira, and E. Costa. A comparative study of different grammar-based genetic programming approaches. In J. M. et. al, editor, *European Conference on Genetic Programming, EuroGP 2017*, volume 10196 of *LNCS*, pages 311–325. Springer, 2017.
18. N. Lourenço, F. B. Pereira, and E. Costa. Unveiling the properties of structured grammatical evolution. *Genetic Programming and Evolvable Machines*, 17(3):251–289, 2017.
19. S. Luke. Two fast tree-creation algorithms for genetic programming. *IEEE Transactions on Evolutionary Computation*, 4(3):274–283, 2000.
20. S. Luke and L. Panait. A comparison of bloat control methods for genetic programming. *Evolutionary Computation*, 14(3):309–344, September 2006.
21. E. Medvet. A comparative analysis of dynamic locality and redundancy in grammatical evolution. In J. M. et. al, editor, *European Conference on Genetic Programming, EuroGP 2017*, volume 10196 of *LNCS*, pages 326–342. Springer, 2017.
22. M. Nicolau. Automatic grammar complexity reduction in grammatical evolution. In R. P. et. al, editor, *Genetic and Evolutionary Computation Conference, GECCO 2004*, 2004.
23. M. Nicolau. Understanding grammatical evolution: initialisation. *Genetic Programming and Evolvable Machines*, 18(4):1–41, Dec 2017.
24. M. Nicolau, A. Agapitos, M. O’Neill, and A. Brabazon. Guidelines for defining benchmark problems in genetic programming. In *IEEE Congress on Evolutionary Computation, CEC 2015*, 2015.

25. M. Nicolau and I. Dempsey. Introducing grammar based extensions for grammatical evolution. In *IEEE Congress on Evolutionary Computation, CEC 2006*, pages 2663–2670, 2006.
26. M. Nicolau and M. Fenton. Managing repetition in grammar-based genetic programming. In T. Friedrich, editor, *Genetic and Evolutionary Computation Conference - GECCO 2016, Denver, CO, USA, July 20-24, 2016, Proceedings*, pages 765–772. ACM, 2016.
27. M. Nicolau, M. O’Neill, and A. Brabazon. Termination in grammatical evolution: Grammar design, wrapping, and tails. In *IEEE Congress on Evolutionary Computation, CEC 2012*, pages 1–8, 2012.
28. M. O’Neill and A. Brabazon. mpga: The meta-grammar genetic algorithm. In M. K. et. al, editor, *European Conference on Genetic Programming, EuroGP 2005*, volume 3447 of *LNCS*, pages 311–320. Springer, 2005.
29. M. O’Neill and C. Ryan. Evolving multi-line compilable c programs. In R. P. et. al, editor, *European Workshop on Genetic Programming, EuroGP 99*, volume 1598 of *LNCS*, pages 83–92. Springer Berlin Heidelberg, 1999.
30. M. O’Neill and C. Ryan. *Grammatical Evolution - Evolutionary Automatic Programming in an Arbitrary Language*, volume 4 of *Genetic Programming*. Kluwer Academic, 2003.
31. M. O’Neill, C. Ryan, M. Keijzer, and M. Cattolico. Crossover in grammatical evolution. *Genetic Programming and Evolvable Machines*, 4(1):67–93, 2003.
32. M. O’Neill, C. Ryan, and M. Nicolau. Grammar defined introns: An investigation into grammars, introns, and bias in grammatical evolution. In L. S. et. al, editor, *Genetic and Evolutionary Computation Conference, GECCO 2001*, pages 97–103. Morgan Kaufmann, 2001.
33. M. O’Neill, J. M. Swafford, J. McDermott, J. Byrne, A. Brabazon, E. Shotton, C. McNally, and M. Hemberg. Shape grammars and grammatical evolution for evolutionary design. In G. R. et al., editor, *Genetic and Evolutionary Computation Conference, GECCO 2009*, pages 1035–1042. ACM, 2009.
34. L. Pagie and P. Hogeweg. Evolutionary consequences of coevolving targets. *Evolutionary Computation*, 5(4):401–418, 1997.
35. D. Robilliard, S. Mahler, D. Verhaghe, and C. Fonlupt. Santa fe trail hazards. In E.-G. T. et. al, editor, *International Conference on Evolution Artificielle, EA 2005*, volume 3871 of *LNCS*, pages 1–12. Springer, 2005.
36. F. Rothlauf and M. Oetzel. On the locality of grammatical evolution. In P. C. et. al, editor, *European Conference on Genetic Programming, EuroGP 2006*, volume 3905 of *LNCS*, pages 320–330. Springer, 2006.
37. C. Ryan and A. Azad. Sensible initialisation in grammatical evolution. In E. C.-P. et. al, editor, *Genetic and Evolutionary Computation Conference, GECCO 2003*. AAAI, 2003.
38. C. Ryan, J. Collins, and M. O’Neill. Grammatical evolution: Evolving programs for an arbitrary language. In W. B. et. al, editor, *European Workshop on Genetic Programming, EuroGP 1998*, volume 1391 of *LNCS*, pages 83–95. Springer, 1998.
39. C. Ryan, M. Keijzer, and M. Nicolau. On the avoidance of fruitless wraps in grammatical evolution. In E. C.-P. et. al, editor, *Genetic and Evolutionary Computation Conference, GECCO 2003*, volume 2724 of *LNCS*, pages 1752–1763. Springer, 2003.
40. J. Tavares and F. B. Pereira. Automatic design of ant algorithms with grammatical evolution. In A. M. et. al, editor, *European Conference on Genetic Programming, EuroGP 2012*, volume 7244 of *LNCS*, pages 206–217. Springer, 2012.
41. E. J. Vladislavleva, G. F. Smits, and D. den Hertog. Order of nonlinearity as a complexity measure for models generated by symbolic regression via pareto genetic programming. *IEEE Transactions on Evolutionary Computation*, 13(2):333–349, 2009.
42. P. A. Whigham, G. Dick, J. Maclaurin, and C. A. Owen. Examining the “best of both worlds” of grammatical evolution. In S. Silva, editor, *Genetic and Evolutionary Computation Conference, GECCO 2015*, pages 1111–1118. ACM, 2015.
43. D. R. White, J. McDermott, M. Castelli, L. Manzoni, B. W. Goldman, G. Kronberger, W. Jaskowski, U.-M. O’Reilly, and S. Luke. Better GP benchmarks: community survey results and proposals. *Genetic Programming and Evolvable Machines*, 14(1):3–29, 2013.