

Repeated Sequences in Linear Genetic Programming Genomes

William B. Langdon

*Computer Science,
University College, London,
Gower Street, London, UK*

Wolfgang Banzhaf

*Computer Science,
Memorial University of Newfoundland,
St. John's, A1B 3X5, Canada*

Biological chromosomes are replete with repetitive sequences, microsatellites, SSR tracts, ALU, and so on, in their DNA base sequences. We started looking for similar phenomena in evolutionary computation. First studies find copious repeated sequences, which can be hierarchically decomposed into shorter sequences, in programs evolved using both homologous and two-point crossover but not with headless chicken crossover or other mutations. In bloated programs the small number of effective or expressed instructions appear in both repeated and nonrepeated code. Hinting that building-blocks or code reuse may evolve in unplanned ways.

Mackey–Glass chaotic time series prediction and eukaryotic protein localization (both previously used as artificial intelligence machine learning benchmarks) demonstrate the evolution of Shannon information (entropy) and lead to models capable of lossy Kolmogorov compression. Our findings with diverse benchmarks and genetic programming (GP) systems suggest this emergent phenomenon may be widespread in genetic systems.

“DNA whose sequence is not maintained by selection will develop periodicities as a result of random crossover.”
George P. Smith [17].

1. Introduction

It has been long noticed that there are emergent phenomena in genetic programming (GP) runs unintended by the human designer of the algorithm. Early on it was observed that code which does not change the output of the program (i.e., noneffective code) appears in many GP runs [1, 2, 3]. It was also noted that bloat affects many GP systems. Reasons for bloat and noneffective code have been examined in years past [4, 5, 6] and remedies have been developed that are more or less effective under particular circumstances (e.g., [7, 8, 9, 10]).

Here we would like to argue that noneffective code and bloat are only the tip of an iceberg and that there is more to be discovered about “emergent phenomena” in GP runs. Particularly, we would like to study the repetition of patterns in GP-evolved programs. These are instructions, or more interestingly, groups of instructions, that occur several times in a program. In fact long sequences of instructions which are repeated can sometimes be decomposed into shorter repeated sequences. This is interesting in itself and it parallels what has been found in natural genomes. Biologists have long noticed the curious existence of repeated sequences in genomic DNA.

Perhaps the reasons for emergence of repeated sequences is similar in biological and artificial evolutionary systems? What could we learn from biological explanations, and can we transfer understanding from evolutionary algorithms back into biology? What instruments are available for observing and examining repetitive sequences? Are there new representations of GP that might be more conducive to evolution once the reason for emergence of repeated sequences has been understood? Are we on the way to discover that evolution reuses code in a very interesting, yet hardly intelligible way? Are building blocks involved in the formation of repeated sequences? These and more questions are raised by our observations.

We first discuss the biological background to repeated sequences. Section 3 describes the two linear GP systems used for our experiments and the time series prediction and bioinformatics protein classification tasks they were applied to. Section 4 presents results of our experiments. Section 5 concludes.

2. Biological background

Biologists have discovered that there is a vast amount of repetition in the DNA of microbes, plants, and animals [11]. Less than 3% of a human genome consists of protein-coding genes but about 50% of it consists of repetitive sequences, many of viral origin [12, 13]. Initially biologists concentrated upon understanding the protein-coding part of genomes. However with whole genome analysis becoming more common, repetitive DNA is a lively subject of research [14, 15, 16].

There are many forms of repeated DNA. This multitude confirms that it is a complex phenomenon. There are satellites, minisatellites, and microsatellites. Repeats of different sizes are located next to each other along the genome. There are ALU repeats and interspersed repetitive sequences. Repeating sequences are found in coding, noncoding, and intergenic areas. Repeats are distributed over genomes and species and constitute a considerable fraction of all DNA in many organisms.

The search for causes began some time ago. Smith, in 1976, did numerical experiments in order to explain the evolution of repeated

DNA sequences [17]. His conclusion was that homologous crossover is a major factor in the emergence of repeated sequences. In more recent work, crossover and DNA duplication have been identified to play an important role. Driven by the inaccuracy of the DNA replication machinery, repeated sequences are both a consequence of misalignments and a cause for crossover [18]. Hsieh and Lee considered a model of bacterial genome growth working with a mechanism called “random segmental self-copying” [19]. This model was able to explain, at a statistical level, the distribution of patterns found in bacterial genomes. They concluded, since the statistical traces are still visible in the distribution of DNA patterns, that growth processes of genomes must have taken place.

In recent years, quantitative analysis tools have become available in molecular biology that allow a closer look at these phenomena [20, 21, 22]. These will provide the opportunity to observe even more closely how different repetitive patterns emerged during evolution. Applications of repetitive sequences are also starting to appear [23].

3. The linear genetic programming system

Genetic programming (GP) has been described many times [24, 25, 26]. Essentially GP applies the well-known genetic algorithm (GA) to the task of searching for a program which does what is needed. Unlike the typical GA, GP evolves both the numeric values and the form of the solution. This additional freedom allows trial solutions of changing complexity and opens up evolutionary computation to new, unexpected and emergent effects.

We decided to start with linear GP [27, 28, 29, 30] in which the evolved data or chromosome is a linear program. Such a linear structure is analogous to the DNA molecule in biological genomes. This makes linear GP a good place to start looking. A linear structure is also intrinsically easier to search for repeated substructures.

In order to show the wide-spread nature of repeated sequences, we used two radically different linear GP systems. GPengine, a standard academic C++ system, and a commercial machine code GP system, Discipulus. To allow repeatability, sections 3.1 to 3.7 describe GPengine and Discipulus in detail, while sections 3.8 and 3.9 describe the two benchmarks used. Sections 4.1 and 4.2 describe the results obtained, while the remainder of section 4 describes the repeated sequences evolved by crossover.

3.1 Tournament selection and steady state population

GPengine and Discipulus both use a steady state population and tournament selection. In a steady state population individuals are steadily

added and removed from the population [31]. In contrast, in a generational system, at the end of each generation, the whole population is replaced by the next. In steady state systems there are no distinct generations. In terms of measuring evolutionary time in a population of M individuals, a generation equivalent is the time taken to create and kill M individuals.

GPengine and Discipulus tournament selection are similar. Four distinct individuals are chosen at random¹ from the population. The fitness of the first two are compared, giving a winner and a loser. In the event that they have identical fitness (e.g., root mean squared [RMS] error) the tie is broken arbitrarily. The second pair are compared in the same way to give a second winner and loser.

The offspring produced from the two winners (by crossover, mutation, or copying [cloning]) replace the two losers. Note, that each tournament always produces exactly two children and the same method is used to produce both children.

Using this form of tournament selection in a steady state population means the best in the population cannot get worse. The best individual, however, is not immortal. If more than one individual has the smallest RMS error, the best individual may by chance be deleted (and replaced by the offspring of one of the other individuals which also had the smallest error).

Discipulus differs only in details. Discipulus saves the best program found so far as it runs. At the end of a run, we analyze the best of these. Note that this program may have been found hundreds of generations before the end of the run. While in GPengine runs, when we say “the best program,” we will mean the best program in the population at the end of the run.

■ 3.2 Linear genetic programming representation and program evaluation

In GPengine each individual consists of a linear sequence of instructions. After input to the program is provided by initializing certain registers, the sequence is executed and register values are changed accordingly. By convention, the output of a GP program resides in the first register (R0).

Each instruction takes two inputs, performs its (integer) calculation, and writes the output to a register. The first input is always a register. The second can either be a constant (0..127) or a register. Figure 1 describes a single instruction. We use eight 8-bit read-write registers. As mentioned, before the individual is executed, all the registers are initialized with data for the current fitness case. The sequence of instructions is obeyed from the start of the individual to its end. The final value in register R0 is the GP's output, that is, its prediction.

¹GPengine uses the C rand function.

Output R0..R7	Arg 1 R0..R7	Opcode + - * /	Arg 2 0...127 or R0..R7
------------------	-----------------	-------------------	----------------------------------

Figure 1. Format of a GPengine instruction.

Again, Discipulus is slightly different. It does not interpret the instructions, it executes them directly. To allow this they are Intel 486 machine code instructions, packed into 4-byte words (padded with nops where needed to fill four bytes). Inputs and constants are held in separate registers from the read/write floating point arithmetic registers.

■ 3.3 GPengine crossover (XOA and 2XO)

90% of tournaments are followed immediately by crossover of the two winners, yielding two children which overwrite the two losers. In the other 10% of cases, the losers are overwritten by copies of the winners. Two-point crossover is used (see Figure 2) however, GPengine appends to the end of the first parent if the code to be copied from the second does not overlap with the first. For this to happen the second parent must be longer than the first (see Figure 3). In a second set of experiments this append variation was disabled and insertion was used in all cases. If the chosen crossover points result in potential offspring that would exceed the maximum size (500 instructions) then the crossover is aborted and the loser is not overwritten. Note that the length checks are made independently, so the other crossover may proceed. Even if the loser is not replaced by crossover, it may still be changed by mutation.

Note that in GP we take it for granted that the parent programs are aligned at their starts. This provides a huge degree of both syntactic and semantic homology for free. This is similar to Nature, where chromosomes are crossed on a like-for-like basis. But at the detailed level where natural crossover occurs, Nature has to work to find matching DNA sequences to establish crossover points.

■ 3.4 Discipulus crossover (2XO and HXO)

Discipulus provides two crossover operators. The first, two-point, is essentially the same as GPengine's two-point crossover. The second is called "homologous" crossover (HXO) (see Figure 4). Here the same two crossover points are used in both parents. So while HXO exchanges code, the offspring are the same length as their parents [32].

■ 3.5 Variable length mutation: Headless chicken crossover

Initial programs are quite short. In order to study if crossover was uniquely responsible for repeating sequences we used a mutation opera-

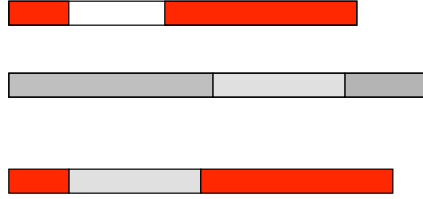


Figure 2. GPengine crossover. Two instructions are randomly chosen in each parent (top two genomes) as cut points. If the code to be inserted from the second parent at least partially overlaps the first, it is inserted in the normal way to give the child (lower chromosome). With headless chicken crossover, the inserted code is randomly generated.



Figure 3. XOA crossover. If there is no length overlap between code selected in the second parent and the first parent (top), the selected code fragment is appended to the whole of the code from the first parent. (That is, the middle portion of the first parent is not removed.)

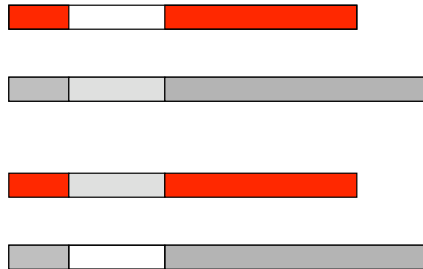


Figure 4. HXO crossover. As with the other crossover operations, two parents (top two programs) crossover to yield two child programs (bottom). In HXO the two crossover cut points are the same in both parents. Note that the code does not change its position relative to the start of the program (left edge) and the child programs are the same lengths as their parents.

tor that could change program lengths. We introduced headless chicken crossover (HCX) to linear GP in [33]. Although described as a crossover operation, only the length of the second parent has any influence on the child. Initially HCX works in the same way as two-point crossover (see Figure 2) except that, instead of inserting a code fragment taken from the second parent, a randomly generated sequence of code of the same length is inserted.

GPengine, unlike Discipulus, does not write-protect its inputs, this means that a long sequence of random instructions will eventually overwrite all the registers. Since the instructions are not reversible, each overwrite destroys information. If the random sequence is long enough it is virtually guaranteed to destroy all information in the registers. Once that happens a program's initial conditions cannot affect its subsequent behavior. Such programs are useless at predicting and so have large RMS errors. Assuming each overwrite is 100% destructive, a random sequence of about $8(\log 8 + \gamma) \approx 21.3$ instructions will render the offspring useless [34]. The expected size of the crossover fragment is $1/l \sum_{i=0}^{l-1} 1/2(l-i) = (l+1)/4$, where l is the number of instructions in the second parent. Hence we anticipate that runs using only HCX will not bloat much above 84 instructions.

When the second program is long enough, HCX becomes like a supersonic jet nozzle. Flow downstream of the nozzle is independent of that upstream. Similarly, program outputs (which are downstream of the random code) are independent of inputs. That is, they are disconnected from upstream perturbations.

■ 3.6 GPengine point mutation

After two children have been produced by crossover or by simply copying their parents (cloning), there is a 40% chance that they will be mutated. Point mutation consists of choosing uniformly at random exactly one instruction in the individual and changing it. Each of the four fields in the chosen instruction (cf. Figure 1) is equally likely to be changed. Apart from ensuring the new instruction is different, the mechanism is the same as that used to create the initial population. Note this means the second argument is approximately equally likely to be a constant (0..127) or a register (R0..R7). The other three fields are chosen uniformly from their legal values.

■ 3.7 Discipulus mutation

Discipulus mutation is similar to GPengine's point mutation. It, too, mutates offspring whether they have been created by crossover or copied (reproduced) directly from their parents. There are three types of Discipulus mutation. Each ensures that the offspring is still a valid machine code program. They do not change the program's size. They

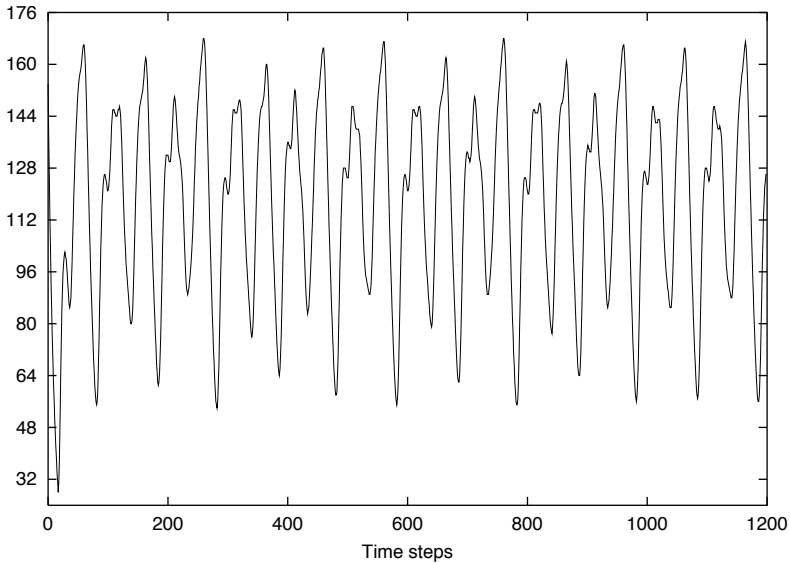


Figure 5. Discrete Mackey–Glass chaotic time series.

are (1) replace a 4-byte instruction block with another four bytes of randomly generated machine code, (2) replace an instruction (op-code) with another machine code instruction of the same length, (3) replace an instruction's operand (input, variable, or constant) with another input, variable, or constant.

■ 3.8 Mackey–Glass benchmark and GPengine

Since the goal is to study the long term behavior of an evolving population of programs we need a moderately difficult task. The population should continually improve and neither get stuck because the problem is too hard nor quickly find the optimal solution. We chose the problem of time series prediction as this is both hard and interesting. Indeed it has applications in scientific, medical, and financial modeling [35]. We used the IEEE benchmark Mackey–Glass chaotic time series.² Mackey–Glass is a continuous problem. The benchmark converts it to discrete time and we digitized the continuous data to give byte-sized integers (by multiplying by 128 and rounding to the nearest integer), see Figure 5.

The task for the GP is, given historical data, to predict the next value. Since the series is chaotic this cannot be done exactly. GP is given eight earlier values from the series. Arguably the most useful is that from the previous time step (which is loaded into R7) but values 2 time steps

²<http://neural.cs.nthu.edu.tw/jang/benchmark/>, $\tau = 17$, 1201 data points, sampled every 0.1).

ago, 4, 8, 16, 32, 64, and 128 time periods back are also available. As with the benchmark, values before the start of the sequence are set to zero. Note that the GP system only has 8-byte registers, and if it needs scratch registers, it may have to sacrifice one or more inputs to store intermediate results.

■ 3.9 Locating animal proteins with Discipulus

Once again we wanted a challenging problem but radically different from predicting tomorrow's share price from a time series. So we chose a binary classification (rather than continuous regression) problem from biology (rather than finance or engineering). It is possible to predict the function of a protein from the fractions of each amino acid from which it is made. This has the advantage that, in many cases, data is readily available from Swissprot or can also be inferred from gene DNA sequences. Naturally, with such a crude measure only limited predictions can be made. Nevertheless Reinhardt and Hubbard [36] showed amino acid composition can be used to predict the location of the protein from which it may be possible to infer something of its function. In [36] they used machine learning to differentiate between seven cellular locations spread across both animals (eukaryotic, they excluded plants) and microbes (prokaryotic).

We restrict ourselves to localizing animal proteins (one normally knows if a protein is animal or bacterial). Since we are seeking a binary classification problem, we evolve models which predict if an animal protein will be found in the cell nucleus or elsewhere (i.e., in the cell cytoplasm, in the mitochondria, or outside the cell [36]). We used the same Swissprot data for 2427 proteins as [36]. There are 1097 nuclear (and 1330 nonnuclear) sequences of amino acids (see Figure 6).

We counted the number of each amino acid in each protein and allocated nuclear proteins to class 1 and the others to class 0. This gave 2427 records, each containing 20 integers and a class label (0/1). Following [37] we split these evenly into a training and a test file.

■ 4. Experimental results

■ 4.1 Predicting the Mackey–Glass chaotic time series

Three pairs of two groups of 10 independent runs were made. In the first pair GPengine's default crossover (with append, XOA) was used. In the second pair two-point crossover (without append, 2XO) was used. Finally, the last pair used HCX. In the second of each pair, selection was turned off by deciding which individuals win or lose each tournament entirely at random. All runs use point mutation (cf. Table 1).

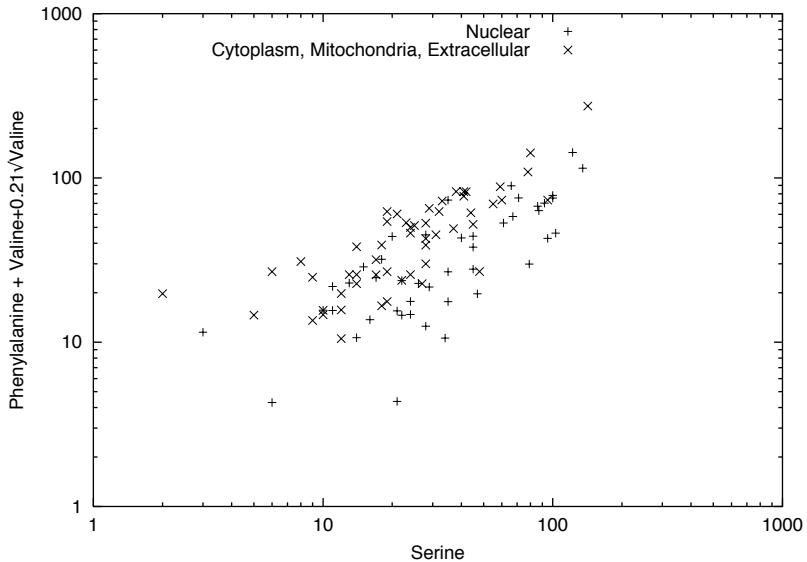


Figure 6. Number of amino acids in nuclear and nonnuclear proteins used with Discipulus. To reduce clutter only 5% of the data are presented. The three (of 20) amino acids plotted were selected by Discipulus as being the best discriminators. The nonlinear function of number of Valines was also suggested by an evolved model.

Objective:	Evolve a prediction for a chaotic time series.
Function set:	+ - × ÷ ^a (operating on unsigned bytes).
Terminal set:	Eight read-write registers, constants 0..127. Registers are initialized with historical values of time series. R0 128 time steps ago, R1 64, R2 32, R3 16, R4 8, R5 4, R6 2, and finally R7 with the previous value. Time points before the start of the series are set to zero.
Fitness:	Root mean error between GP prediction (final value in R0) and actual (averaged over 1201 time points).
Selection:	Steady state, tournament 2 by 2.
Initial population:	Random program's length is uniformly chosen from 1..14.
Parameters:	Population 500, max program size 500, 90% crossover, 40% mutation.
Termination:	125 500 individuals evaluated.

^aIf the second argument of ÷ is zero, it returns zero.

Table 1. GPengine parameters for Mackey–Glass time series prediction.

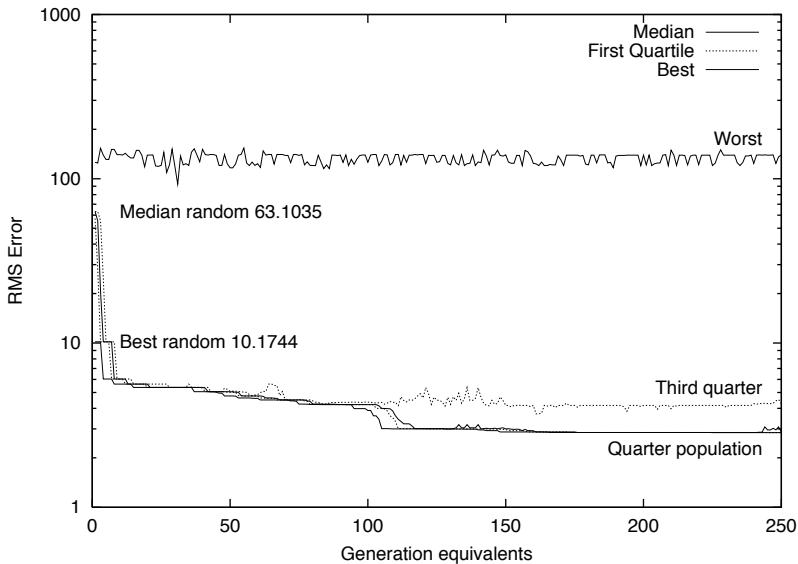


Figure 7. Evolution of Mackey–Glass prediction error (first of 10 runs). Note that the population chases after the best (lowest) fitness. For many generations at least 25% of the population has the same best fitness (sometimes more than half).

	RMS error										Mean
XOA	2.85	2.30	3.56	3.34	3.68	4.30	2.24	5.37	2.38	4.40	3.44
	29.40	6.26	30.20	30.17	30.18	8.03	30.07	30.17	19.17	30.22	24.39
2XO	3.53	3.47	1.60	4.27	5.37	2.43	3.81	5.37	5.37	2.72	3.79
	8.60	12.59	7.66	33.32	14.40	19.62	6.23	17.37	29.85	23.63	17.33
HCX	4.03	4.04	3.64	4.06	3.93	3.61	3.73	3.20	3.78	3.94	3.80
	9.95	6.32	9.95	11.71	16.59	15.83	7.92	7.37	10.71	8.60	10.49

Table 2. Best Mackey–Glass prediction error at the end of GPengine runs. Using two-point crossover with append (XOA), without append (2XO), and headless chicken crossover (HCX). The lower row of each pair refers to the runs without fitness selection.

In all 3×10 runs with selection, fitness improved and for many generations large parts of the population had the same fitness. Figure 7 shows the evolution of prediction error for the first of 10 runs (the others are similar, cf. Table 2).

Figure 8 shows the evolution of program size. Initially programs are between one and 15 instructions long, with a mean of seven. However, in runs with fitness selection and crossover (XOA and 2XO) length

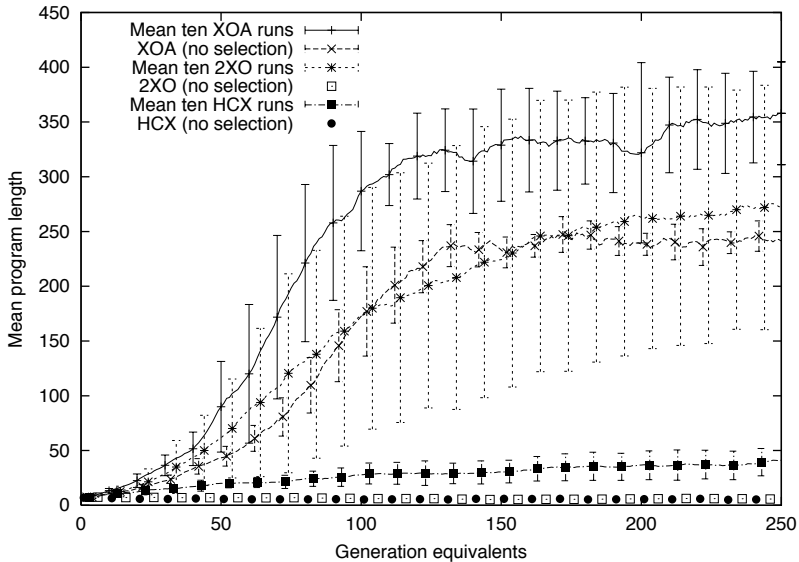


Figure 8. Evolution of mean program size and variation between Mackey–Glass prediction runs. (Variation between HCX and 2XO without selection, lowest plot, is too small to show, and so is omitted.) Except for the two-point crossover with append (XOA) runs, selection is required for bloat. With crossover (XOA or 2XO) the mean population size appears to increase exponentially, until constrained by the maximum size limit of 500.

quickly increases and the longest program is either 500 or very near to this limit. Such bloat was expected [4]. As predicted in section 3.5, in mutation-only runs (HCX) with selection, the increase in size is less dramatic. However it was a surprise to see bloat in runs without selection when using crossover with append (XOA) [38]. An initial thought was that this was due to the asymmetric append variation of the crossover operation. This appears to be correct, since when the variant is removed and normal two-point crossover linear GP is used instead, bloat does not appear without fitness selection. (See the lower lines in Figure 8.)

■ 4.2 Predicting protein location

For each experiment on the Swissprot animal proteins (cf. section 3.9) Discipulus was run 10 times, always in classification mode and with identical parameters but with different initial random number seeds. As far as possible we used Discipulus' default configurations (actual values are given in Table 3). In both experiments the maximum program size was 2048 bytes. In the first experiment, “homologous crossover” was

Objective:	Evolve a prediction of nuclear or nonnuclear location for animal proteins based on their amino acid composition.
Terminal set:	Two read-write FPU registers, 43 randomly chosen constants. Number (integer) of each of the 20 amino acids in the protein. (Codes B and Z are ambiguous. Counts for B were split evenly between aspartic acid D and asparagine N. Those for Z between glutamic acid E and glutamine Q.)
Fitness:	DSS [37, 39]. Parsimony not used.
Selection:	Steady state, tournament 2 by 2.
Initial population:	Each random program's length is chosen uniformly at random from the range 4 to 80 bytes.
Parameters:	Population 500 (10 × 50 demes), max program size 2048 (bytes), 95% crossover (either all 2XO or 95% HCX and 5% 2XO), and 95% mutation (three types 30%, 30%, 40%).
Termination:	500 000 individuals evaluated.

Table 3. Discipulus parameters used in animal protein location prediction experiments. Only the maximum program size and HCX were changed from factory defaults.

	Percent correct nuclear vs. nonnuclear prediction										Mean
2XO	82	80	80	80	80	81	82	78	79	81	80.3
HXO	81	81	79	80	80	79	78	81	78	81	80.0

Table 4. Accuracy of best animal protein location predictors evolved using Discipulus two point (2XO) and HXO (homologous) crossover. Performance on training and test datasets is always within a few percent of each other (indicating little over fitting) therefore we report performance averaged across both.

disabled, whilst in the second it was left at its default rate (95%). Table 4 gives the accuracy of the best models found by each run.

As an aside, when Discipulus was allowed to run in its default multiple-run mode, it evolved a team solution with an accuracy of 87.1%. Reinhardt and Hubbard [36] did not do a nuclear *versus* non-nuclear classification. Instead they report results for the easier problem of classifying animal nuclear proteins against each of their other classes one at a time. Table 2 in [36] gives mean accuracies of 84.9%, 84.8%, and 86.1% when classifying nuclear proteins against each of the three locations individually. Given the variance of the neural network runs, one cannot say that Discipulus did significantly better, but it is doing at least as well on a harder version of the problem.

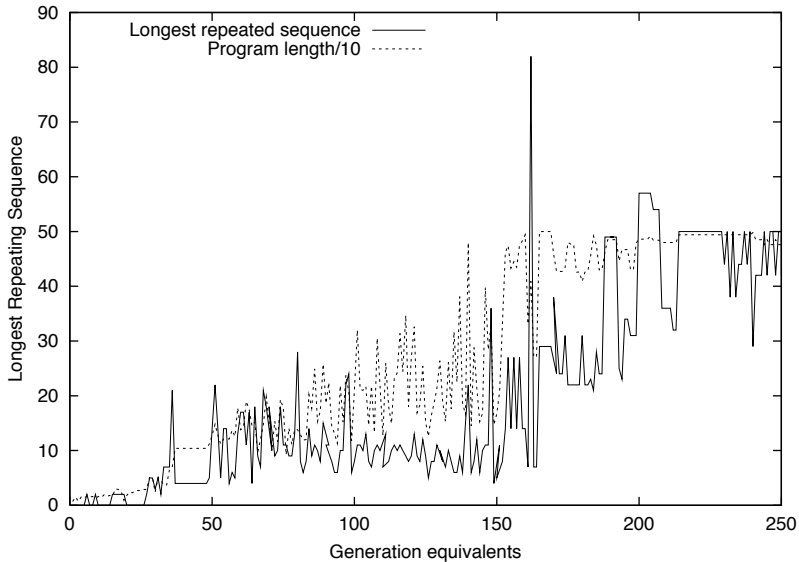


Figure 9. Evolution of length of longest repeated sequence of instructions in the best Mackey–Glass prediction program produced by the first run with two-point crossover (2XO) and fitness selection. The length of the programs is also shown.

4.3 Repeated program instruction sequences

In the random initial GPengine programs there are no repeated sequences. They are overwhelmingly unlikely to arrive by chance. However, as crossover, mutation, and selection get to work and programs grow, instructions start to become repeated. Initially just single instructions are repeated but the length and number of repeats increases (see Figure 9).

All the best programs of the 10 runs with the append crossover variant (XOA) contained repeated sequences (see Figure 10). The longest sequences contained from 12 to 62 instructions. All of these occurred twice, however the programs also contained other, distinct, shorter sequences which occurred multiple times. Again the XOA runs without selection throw up a surprise: eight of the 10 best programs³ contain sequences of instructions which are repeated. All 10 runs with two-point crossover without the append variation (2XO), produce repeated sequences, however none of the 10 2XO runs without fitness selection produced repeated instructions.

Both the 10 2XO and 10 HXO Discipulus runs evolved binary classification programs with repeated instructions like those found in

³Even in the absence of selection one can observe the quality of programs by evaluating the fitness function.

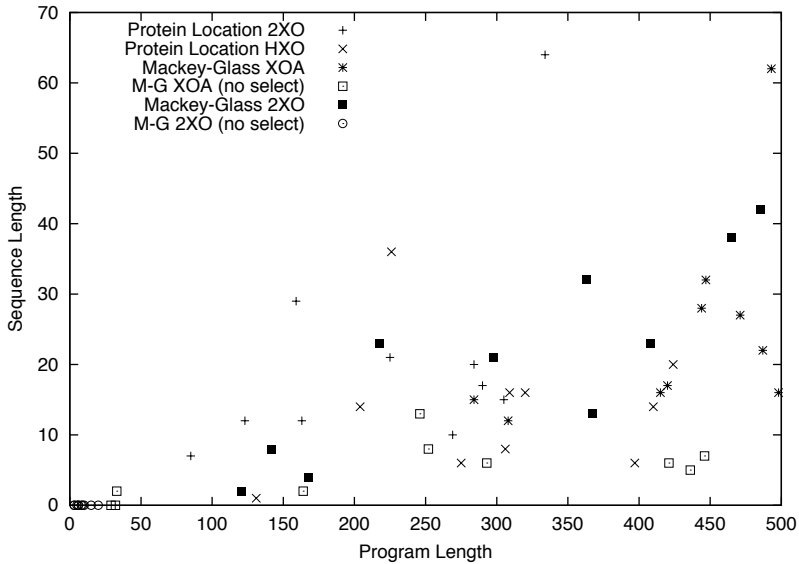


Figure 10. Length of longest repeated sequence of instructions in the best prediction programs. 2×10 Discipulus protien location runs and 4×10 Mackey–Glass prediction runs. With fitness selection, all three types of crossover evolved repeating sequences, as do eight of the 10 XOA runs without fitness selection. No 2XO runs evolve repeated sequences when tournaments are random.

GPengine 2XO runs. In 19 cases the longest sequence was between six and 64 instructions. (The remaining run found a short program containing a few individual instructions repeated a small number of times, but no repeated sequences.)

Figure 10 plots the variation of maximum length of repeated instructions in each of the 6×10 best of run programs against their size. As an alternative to saying that repetitive sequences are due to the crossover operator, Figure 10 suggests that the length of the programs (i.e., bloat) is more important. To some extent this is borne out by the runs with HCX. The 10 runs with selection produced the best predictors of between 18 and 76 instructions (those without selection contained between 2 and 21 instructions). None contained repeated instructions.

Figure 11 shows the location of repeated sequences in a single evolved program stressing the sheer number of repeated instructions. It also shows some repeated instructions are part of long sequences (up to 10% of program length) and that repeated instructions appear throughout the program. (The tartan pattern in Figure 11 suggests there may be other structures which have yet to be investigated.)

Our results strongly suggest that crossover is responsible for repeated sequences. However we cannot rule out the possibility that bloated long

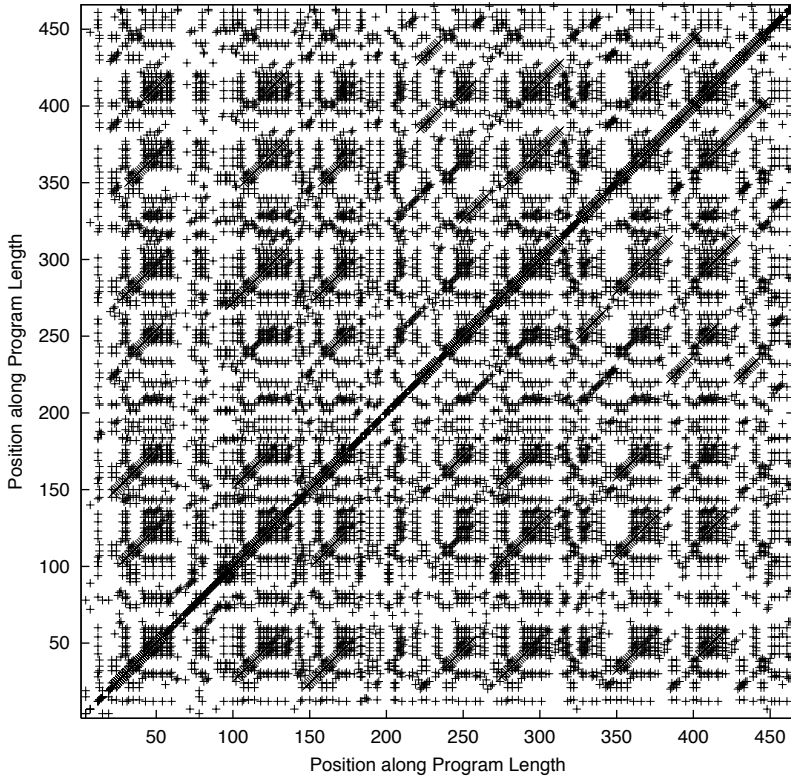


Figure 11. Location of repeated instructions in the best Mackey–Glass prediction program produced by the third run with two-point crossover (2XO) and fitness selection. Instructions that are part of repeated sequences longer than 10 are plotted with \times , those less, with $+$. This run was chosen as it best highlights the structure. Notice that almost every instruction is repeated, so the diagonal is almost solid. (The figure is symmetric about the diagonal.) Each cross in a vertical line, say $x = 400$, indicates an instruction which is identical (to that at position 400).

programs produced by mutation might also contain repeated sequences, although this does seem unlikely. As mentioned before, mutation is unable to produce long programs, so we cannot test this.

■ 4.4 Effective code

Rapid increase in length is a characteristic of bloat [4]. We used algorithm 2 from [40] to analyze the best predictors of the GPengine 2XO and XOA runs. We also used Discipulus' own intron removal tool. These showed that the majority of instructions have no impact on

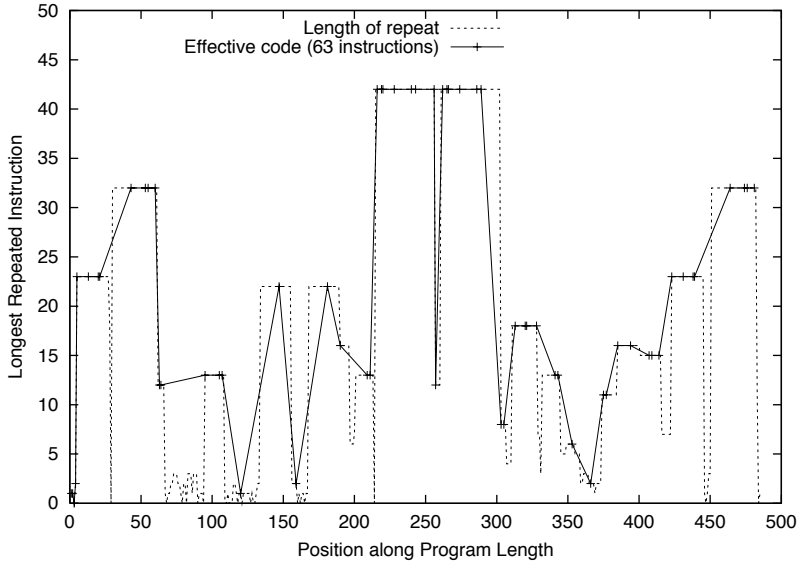


Figure 12. Distribution of repeated sequences along the length of the best Mackey–Glass predictor at end of first 2XO run. It achieves an RMS error of 3.5 using only data from one and eight previous time steps. The solid line highlights the location of its 63 effective instructions. An animation can be found via <http://www.cs.ucl.ac.uk/staff/W.Langdon/gecco20041b/>.

the output of the programs, that is, they are ineffective code (introns). Figure 12 shows the distribution of instructions which could affect the prediction along the length of one program. The other bloated best of 2XO GPengine runs are similar to that shown in Figure 12. However in three runs there was less bloat. Their best predictors are much shorter, containing only one effective instruction, which is near the end. There is no obvious correlation between whether an instruction is effective and how many times it is repeated.

As the fraction of ineffective code increases so mutation is more and more unlikely to change a program’s performance. This is (part of) the reason why Mackey–Glass populations converge towards a single fitness, cf. Figure 7. Similarly as evolution increases the number of repeated instructions, mutation increasingly often duplicates an existing instruction. However, on average, 88% of point mutations to the best Mackey–Glass programs evolved by 2XO produce a new instruction.

■ 4.5 Entropy and information content

There are $8 \times 8 \times 4 \times (128 + 8) = 34816$ GPengine legal instructions (cf. Figure 1) Since $\log_2 34816 = 15.087463$, a randomly chosen in-

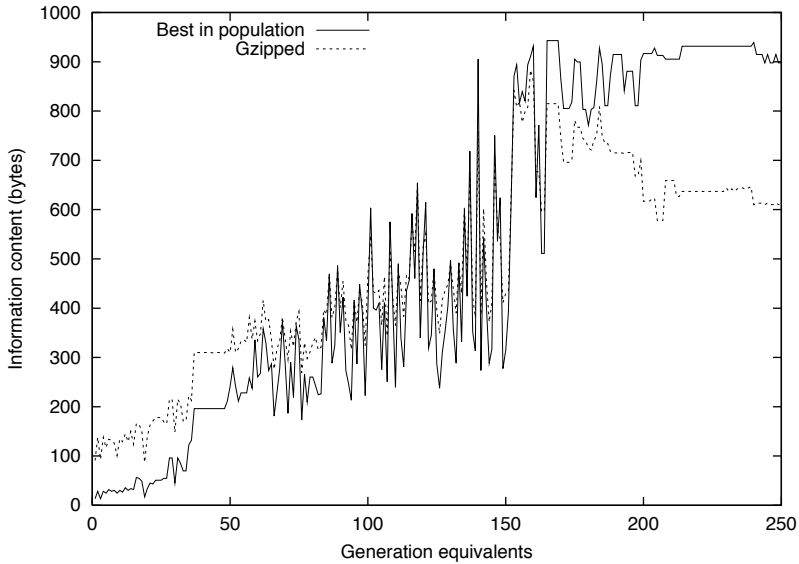


Figure 13. Evolution of information content in the best Mackey–Glass prediction program produced by the first run with two-point crossover (2XO) and fitness selection (as measured by gzip). For comparison the length of the programs are also plotted (but normalized so as to give their precompressed information content).

struction contains slightly more than 15 bits of information. Using this measure suggests that as the population bloats each predictor contains more information. A crude way of estimating actual information content is to compress the programs using gzip [26]. Figure 13 shows that information content increases over time but as programs contain more repeated sequences, gzip’s Lempel–Ziv algorithm is able to compress them. Smaller size gives a lower estimate of information content of the programs. Figure 13 shows that gzip (with default parameters and a simple ASCII text format) initially imposes an overhead of about 100 bytes. After about generation 150, gzip is able to recognize patterns and use them to compress the programs. For comparison our Mackey–Glass benchmark, without compression, contains 8576 bits (1072 bytes) of information ($1201 \times \log_2 141 = 8576$). Lossy (programmatic) Kolmogorov compression is possible and achieves 7.6 times more compression than gzip’s lossless compression.

5. Conclusion

Approximately half of human DNA is composed of repeated sequences. Initial experiments using two very different linear genetic programming

(GP) systems (one evolving Intel floating point machine code, one using a C++ interpreter of integer arithmetic) on two different benchmarks (one time series prediction, the other a bioinformatics binary classification) suggest that where artificial genetic systems have the space to evolve repeated patterns, they will emerge. In our runs, both the number of repeated sequences and their length increase with time but so, too, does the length of the programs. These bloated programs are not random. Evidence suggests that crossover is responsible.

We have observed the evolution of long repeated sequences of instructions. The chances of them being found purely at random are infinitesimal. However, while we anticipate these sequences occur widely, so far we have only observed them in linear GP. Of course it is interesting to see whether the same, or other emergent phenomena, occur in tree GP. Most importantly, can these observations be used to help us build better systems in the future? Finally, could experiments of this type in artificial evolution give insight for biologists? For example, are statistical distributions of repeated sequences comparable to what happens in real genomes?

Acknowledgments

The linear genetic programming system, GPengine, was given by Peter Nordin. We would like to thank Paul Gillard, Marian Wissink, Nolan White, Dick Furnstahl, Frank D. Francone, Markus Conrads, and David H. Jones. Support was provided by a grant from the visitor program of the Department of Computer Science, Memorial University of Newfoundland.

References

- [1] A. Singleton, Walter Tackett [2] citing private communication from Andy Singleton as proposing the “intron” explanation for bloat in GP trees (1994).
- [2] W. A. Tackett, “Recombination, Selection, and the Genetic Construction of Computer Programs,” Ph.D. Thesis, University of Southern California, Department of Electrical Engineering Systems, USA, 1994.
- [3] L. Altenberg, “Emergent Phenomena in Genetic Programming,” in *Evolutionary Programming: Proceedings of the Third Annual Conference*, San Diego, CA, USA, 24–26 February 1994, edited by A. V. Sebald and L. J. Fogel (World Scientific Publishing, pp. 233–241).
- [4] W. B. Langdon, T. Soule, R. Poli, and J. A. Foster, “The Evolution of Size and Shape,” in *Advances in Genetic Programming 3* edited by L. Spector, W. B. Langdon, U.-M. O’Reilly, and P. J. Angeline (MIT Press, Cambridge, MA, USA, June 1999, ch. 8, pp. 163–190).

- [5] P. J. Angeline, "Subtree Crossover Causes Bloat," in *Genetic Programming 1998: Proceedings of the Third Annual Conference*, University of Wisconsin, Madison, WI, USA, 22–25 July 1998, edited by J. R. Koza, W. Banzhaf, K. Chellapilla, K. Deb, M. Dorigo, D. B. Fogel, M. H. Garzon, D. E. Goldberg, H. Iba, and R. Riolo (Morgan Kaufmann, pp. 745–752).
- [6] W. Banzhaf and W. B. Langdon, "Some Considerations on the Reason for Bloat," *Genetic Programming and Evolvable Machines*, 3(1) (2002) 81–91.
- [7] P. Nordin, F. Francone, and W. Banzhaf, "Explicitly Defined Introns and Destructive Crossover in Genetic Programming," in *Advances in Genetic Programming 2*, edited by P. J. Angeline and K. E. Kinneer, Jr. (MIT Press, Cambridge, MA, USA, 1996, ch. 6, pp. 111–134).
- [8] F. D. Francone, M. Conrads, W. Banzhaf, and P. Nordin, "Homologous Crossover in Genetic Programming," in *Proceedings of the Genetic and Evolutionary Computation Conference*, Orlando, FL, USA, 13–17 July 1999, edited by W. Banzhaf, J. Daida, A. E. Eiben, M. H. Garzon, V. Honavar, M. Jakiela, and R. E. Smith, volume 2 (Morgan Kaufmann, pp. 1021–1026).
- [9] W. B. Langdon, "Size Fair and Homologous Tree Genetic Programming Crossovers," *Genetic Programming and Evolvable Machines*, 1(1/2) (2000) 95–119.
- [10] J. V. Hansen, "Genetic Programming Experiments with Standard and Homologous Crossover Methods," *Genetic Programming and Evolvable Machines*, 4(1) (2003) 53–66.
- [11] R. J. Britten and D. E. Kohnen, "Repeated Sequences in DNA," *Science*, 161 (1968) 529–540.
- [12] A. F. A. Smit, "The Origin of Interspersed Repeats in the Human Genome," *Current Opinions in Genetics and Development*, 6 (1996) 743–748.
- [13] C. Patience, D. A. Wilkinson, and R. A. Weiss, "Our Retroviral Heritage," *Trends in Genetics*, 13 (1997) 116–120.
- [14] J. R. Lupski and G. M. Weinstock, "Short, Interspersed Repetitive DNA Sequences in Prokaryotic Genomes," *Journal of Bacteriology*, 174 (1992) 4525–4529.
- [15] G. Toth, Z. Gaspari, and J. Jurka, "Microsatellites in Different Eukaryotic Genomes: Survey and Analysis," *Genome Research*, 10 (2000) 967–981.
- [16] G. Achaz, E. P. C. Rocha, P. Netter, and E. Coissac, "Origin and Fate of Repeats in Bacteria," *Nucleic Acids Research*, 30 (2002) 2987–2994.
- [17] G. P. Smith, "Evolution of Repeated DNA Sequences by Unequal Crossover," *Science*, 191(4227) (1976) 528–535.

- [18] O. Elemento, O. Gascuel, and M. P. Lefranc, “Reconstructing the Duplication History of Tandemly Repeated Genes,” *Molecular Biology and Evolution*, **19** (2002) 278–288.
- [19] L. C. Hsieh and H. C. Lee, “Model for the Growth of Bacterial Genomes,” *Modern Physics Letters*, **16** (2002) 821–827.
- [20] G. Benson, “Tandem Repeats Finder: A Program to Analyze DNA Sequences,” *Nucleic Acids Research*, **27** (1999) 573–580.
- [21] J. W. Bizzaro and K. A. Marx, “Poly: A Quantitative Analysis Tool for Simple Sequence Repeat (SSR) Tracts in DNA,” *BMC Bioinformatics*, **4** (2003) 22–28.
- [22] A. Taneda, “Adplot: Detection and Visualization of Repetitive Patterns in Complete Genomes,” *Bioinformatics*, **5** (2004) 701–708.
- [23] Z. Izsvak, Z. Ivics, and P. B. Hackett, “Repetitive Elements and their Applications in Zebrafish,” *Biochemical Cell Biology*, **75** (1997) 507–523.
- [24] J. R. Koza, *Genetic Programming: On the Programming of Computers by Means of Natural Selection* (MIT Press, Cambridge, MA, USA, 1992).
- [25] W. Banzhaf, P. Nordin, R. E. Keller, and F. D. Francone, *Genetic Programming—An Introduction; On the Automatic Evolution of Computer Programs and Its Applications* (Morgan Kaufmann, dpunkt.verlag, 1998).
- [26] W. B. Langdon, *Genetic Programming and Data Structures: Genetic Programming + Data Structures = Automatic Programming!* volume 1 of *Genetic Programming* (Kluwer, Boston, 1998).
- [27] P. Nordin, “Evolutionary Program Induction of Binary Machine Code and Its Applications,” Ph.D. Thesis, der Universitat Dortmund am Fachereich Informatik, 1997.
- [28] M. S. Atkin and P. R. Cohen, “Learning Monitoring Strategies: A Difficult Genetic Programming Application,” in *Proceedings of the 1994 IEEE World Congress on Computational Intelligence*, Orlando, FL, USA, 27–29 June 1994 (IEEE Press, pp. 328–332a).
- [29] W. Banzhaf, “Genetic Programming for Pedestrians,” in *Proceedings of the Fifth International Conference on Genetic Algorithms, ICGA-93*, University of Illinois at Urbana-Champaign, 17–21 July 1993, edited by S. Forrest (Morgan Kaufmann, p. 628).
- [30] T. Perkis, “Stack-based Genetic Programming,” in *Proceedings of the 1994 IEEE World Congress on Computational Intelligence*, Orlando, FL, USA, 27–29 June 1994, volume 1 (IEEE Press, pp. 148–153).

- [31] G. Syswerda, “Uniform Crossover in Genetic Algorithms,” in *Proceedings of the Third International Conference on Genetic Algorithms*, George Mason University, 4–7 June 1989, edited by J. D. Schaffer (Morgan Kaufmann, pp. 2–9).
- [32] P. Nordin, W. Banzhaf, and F. D. Francone, “Efficient Evolution of Machine Code for CISC Architectures Using Instruction Blocks and Homologous Crossover,” in *Advances in Genetic Programming 3* edited by L. Spector, W. B. Langdon, U.-M. O’Reilly, and P. J. Angeline (MIT Press, Cambridge, MA, USA, June 1999, ch. 12, pp. 275–299).
- [33] P. J. Angeline, “Subtree Crossover: Building Block Engine or Macromutation?” in *Genetic Programming 1997: Proceedings of the Second Annual Conference*, Stanford University, CA, USA, 13–16 July 1997, edited by J. R. Koza, K. Deb, M. Dorigo, D. B. Fogel, M. Garzon, H. Iba, and R. L. Riolo (Morgan Kaufmann, pp. 9–17).
- [34] W. B. Langdon, “Convergence Rates for the Distribution of Program Outputs,” in *GECCO 2002: Proceedings of the Genetic and Evolutionary Computation Conference*, New York, 9–13 July 2002, edited by W. B. Langdon, E. Cantú-Paz, K. Mathias, R. Roy, D. Davis, R. Poli, K. Balakrishnan, V. Honavar, G. Rudolph, J. Wegener, L. Bull, M. A. Potter, A. C. Schultz, J. F. Miller, E. Burke, and N. Jonoska (Morgan Kaufmann Publishers, pp. 812–819).
- [35] H. Oakley, “Two Scientific Applications of Genetic Programming: Stack Filters and Nonlinear Equation Fitting to Chaotic Data,” in *Advances in Genetic Programming*, edited by K. E. Kinnear, Jr. (MIT Press, 1994, ch. 17, pp. 369–389).
- [36] A. Reinhardt and T. Hubbard, “Using Neural Networks for Prediction of the Subcellular Location of Proteins,” *Nucleic Acids Research*, **26**(9) (1998) 2230–2236.
- [37] F. D. Francone, *Discipulus Owner’s Manual*, version 3.0 draft edition, 11757 W. Ken Caryl Avenue F, PBM 512, Littleton, Colorado, 80127-3719, USA, 2001.
- [38] W. B. Langdon and R. Poli, *Foundations of Genetic Programming* (Springer-Verlag, 2002).
- [39] C. Gathercole, “An Investigation of Supervised Learning in Genetic Programming,” Ph.D. Thesis, University of Edinburgh, 1998.
- [40] M. Brameier and W. Banzhaf, “A Comparison of Linear Genetic Programming and Neural Networks in Medical Data Mining,” *IEEE Transactions on Evolutionary Computation*, **5**(1) (2001) 17–26.