

# GI Software with fewer Data Cache Misses\*

W.B. Langdon, Justyna Petke, Aymeric Blot, David Clark  
{W.Langdon,j.petke,david.clark}@ucl.ac.uk  
aymeric.blot@univ-littoral.fr

Department of Computer Science, University College London

## ABSTRACT

By their very name caches are often overlooked and yet play a vital role in the performance of modern and indeed future hardware. Using MAGPIE (Machine Automated General Performance Improvement via Evolution of software) we show genetic improvement GI can reduce the cache load of existing computer programs. Operating on lines of C and C++ source code using local search, Magpie can generate new functionally equivalent variants which generate fewer L1 data cache misses. Cache miss reduction is tested on two industrial open source programs (Google’s Open Location Code OLC and Uber’s Hexagonal Hierarchical Spatial Index H3) and two 2D photograph image processing tasks, counting pixels and OpenCV’s SEEDS segmentation algorithm.

Magpie’s patches functionally generalise. In one case they reduce data misses on the highest performance L1 cache dramatically by 47%.

## CCS CONCEPTS

• **Software and its engineering** → **Search-based software engineering**.

## KEYWORDS

genetic programming, genetic improvement, SBSE, linear representation, software resilience, automatic code optimisation, tabu, nonstationary noise, perf, world wide location, plus codes, zip code, OpenCV, image segmentation

## 1 INTRODUCTION

Jack Dongarra won the Turing Award in 2021. In celebration, in a recent article in the Communications of the ACM [10] he describes the current and foreseeable future limits of high performance computing considering that clock frequencies have not increased significantly for two decades but that Moore’s Law [33] continues to provide exponential increases in transistor count. Dongarra argues that this will continue to fuel the current trend to ever greater degrees of hardware parallelism. He also points out that already computing hardware is limited not by processor speed but by the time taken to get data to the compute engines. That is, computing is cheap, it is data movement that is expensive.

Although widespread, the trend to evermore hardware parallelism is exemplified by graphics cards (GPUs) and similar accelerator architectures (e.g. TPUs) which now form the back bone of both super computers and training deep neural networks. For a long time nVidia were reluctant to introduce general data caches

into their graphics cards, instead insisting the developers of computer games would know the data flows inside their programs and could therefore optimise their software to take best advantage of the huge data bandwidth available in GPUs. This was always unrealistic and, even for very experienced programmers, getting the best performance was very hard. Hence nowadays in an attempt to make performance programming easier, GPU manufacturers boast that their GPUs contain multiple levels of cache memory.

Dongarra tries to argue against the trend to use caches to make life easier for software developers, and says “Instead of just relying on hardware caches, new algorithms must be designed” [10, page 68]. However, he does not say how this will be done. Genetic programming has a long established tradition of inventing new solutions [17]. In contrast to GP, rather than starting evolution from scratch every time, Genetic Improvement (GI) [23, 38] finds updates on existing software. Although primarily used for fixing computer bugs [12] or speeding up code [24], we seek to strengthen the claim that it can optimise any *measurable* aspect of software by showing that, despite noise, the GI tool Magpie can sometimes increase the effectiveness of even the fastest level of data caches in a conventional Intel x86 desk top computer, even on compiler optimised code<sup>1</sup>. Thus GI and GP may form a two pronged attack on the problem of effectively using future data paths in highly parallel hardware, with perhaps GP inventing new algorithms and GI tailoring existing code to take better advantage of novel hardware.

Section 3 details the fitness test cases for Google’s OLC and Uber’s H3 digital mapping programs, the Blue image benchmark, and the OpenCV SEEDS resource intensive image segmentation algorithm. While Section 4 describes using a Coupon Collector argument to choose how much of the search space Magpie should sample. Following the results in Section 5 (see also Table 1), in Section 6 there is a brief description of the C and C++ source code changes made by Magpie, which in the image examples give a reduction in L1 data cache misses even on the existing compiler (GCC -O3) optimised code, and a discussion of non-stationary noise. We conclude in Section 7 that Magpie is ready to use and here it reduced L1 data cache misses by up to 47%. In all four examples the patches functionally generalise but only in the data hungry image processing examples do we see any sustained cache reduction (see Figures 1, 4, 7, 8, 9 and 10). But first we describe the background.

<sup>1</sup>CPU hardware provides both data and instruction caches. These are both fundamental to performance. However typically good locality and the small size of program machine code compared to data it acts upon, means instruction caches now and probably into the future, are less of an issue than data caches. In these experiments with an L1 instruction cache of 32KBytes there are only about a hundred instruction cache misses, while, for example, with SEEDS the number of L1 data misses exceeds 40 000 even on a small image such as Figure 3.

\*Long version of W.B. Langdon, Justyna Petke, Aymeric Blot, David Clark. 2023. Genetically Improved Software with fewer Data Cache Misses. In Genetic and Evolutionary Computation Conference Companion (GECCO ’23 Companion), July 15–19, 2023, Lisbon, ACM. <https://doi.org/10.1145/3583133.3590542>

## 2 BACKGROUND

Computers are universal. Everything relies on them. Even now many of the richest people are rich because they founded very successful software companies at about the turn of the century. IT technology in general and software in particular permeates and will continue to dominate the third millennium. Indeed the planet is addicted to software. Despite programming being more than 60 years old, software continues to be hand written. Automatic programming to a large extent remains a dream.

At GECCO-2009 Stephanie Forrest, ThanhVu Nguyen, Wes Weimer and Claire Le Goues [12] showed that genetic programming [5, 16, 39, 40, 42, 43] could automatically fix bugs in computer software [28, 44]. For the first time artificial intelligence (AI) is being applied to a major problem in software engineering on industry size programs [1]. Since then the field of automatic program repair (APR) has bloomed [32]. Inspired by Stephanie et al. we [23, 24] began applying genetic programming to improving human written software in many ways in addition to bug fixing [38].

Although genetic programming remains a common search technique in genetic improvement (GI), local search is increasingly popular [6]. In addition to ever more powerful computers, GI is able to scale because it does not start from scratch at every run but builds on existing software. Another enormous advantage is that GI can automatically double check with the existing painstakingly hand written code. As with regression testing [14], in effect, the existing program becomes its own specification. That is, it can be used as a test oracle for both automatic functional and non-functional improvements. This could be run time performance (be it elapse time, memory requirements, etc.) and also its functionality. We have mentioned above the success of automatic bug fixing, but functional improvements can include making the program give more accurate answers [26]. Evolution has also been used to adapt web software to colour blind users, and to help tune hearing aids for deaf patients [29]. The increasing success of automatically generated software tests [13] also naturally feeds into genetic improvement.

Genetic improvement research is often via bespoke one-off experiments. David White recognised this and proposed GIN [46] as a generic GI tool for Java programs. Gabin An [3, 4] proposed PyGGI for Python. Both tools have been extensively used and updated: [8, 30, 36, 37] and [2, 6, 15, 31, 41], and further GI tools have been proposed [18]. Nevertheless recently a user study said that GI lacked user friendly tools [47].

### 2.1 Magpie

In response to this Aymeric Blot wrote Magpie, which is not only user friendly but combines GI and parameter tuning. It was released last year as an open source project. Like PyGGI 2.0 [2] (from which it was developed) it is freely available from GitHub<sup>2</sup>. We use it to hopefully convince the reader that evolution can in principle improve any *comparable* measure of software quality. Although also written in Python, it aims to work with any computer programming language. It has been mostly tested on Apple Mac and Linux Laptops but aims to be generic enough to work under Microsoft windows. Here we test it on a Linux desktop and have not attempted to maintain compatibility with Microsoft.

<sup>2</sup><https://github.com/bloa/magpie>

As of 27 November 2022, including examples and documentation, Magpie contains 4871 lines of code, mostly written in Python. It contains examples in Python, C, C++ and Ruby.

### 2.2 Updates to Magpie

In the course of previous work [22], we had enhanced Magpie to use Python’s ctypes to directly call the patched code. This allows us to reduce noise by collecting data directly on the individual C/C++ routines rather than the whole Magpie sub-process. Also to exclude the Python interpreter from our measurement, we clear the L1 data cache before invoking the patch code by setting and reading a fixed array of 32K bytes (the size of our CPU’s L1 data cache). To make the fitness robust, each patch is tested multiple times and since the mean is notoriously suspect to noisy outliers, we use instead the first quartile to summarise the measurements.

In the image segmentation example (Section 3.2) 17% of the lines consist of a single closing brace }. These are naturally interchangeable and so it was discovered that Magpie was wasting a lot of effort testing programs that were identical (apart from white space, etc). To prevent this, the compilation step keeps a “tabu” list [25]. Previously [25] we had a complicated tabu of both genotypic and phenotypic information (of up to 340 MBytes), here we simply keep a copy of each object file (on average 407 files per run, occupying about 30 MBytes). After compilation, semantically identical patches are rejected by simply comparing their object file with object files of the same size from previous patches. New unique patches are added to the tabu directory and identical ones are discarded without fitness testing. Although in the limit this could be slow as the tabu directory grows, in practice the time taken is negligible.

The Linux GNU perf utility allows access to many hardware performance counters. In particular we used the perf run time library `linux/perf_event.h` to collect L1 data and instruction cache misses, the count of instructions executed and elapsed time. (Only the L1 cache data misses are used by the fitness function.)

### 2.3 Datasets Background

We use four open source C/C++ examples. Two industrial geospatial programs, both written in C [21]. One from Google’s OLC and the other from Uber’s H3 (see Section 3.1 and Figure 1). And two C++ photographic image processing examples (Section 3.2). These are: our Blue benchmark from the 2018 Tarot summer school [19] (Figure 2) and OpenCV’s image segmentation code [27] (Figures 3 and 4).

## 3 FITNESS FUNCTION

Magpie attempts to run the patched program on all the test cases. If the patch passes them all, Magpie runs it again multiple times to try to get a good robust estimate of its performance.

In summary: Magpie uses multiple objectives to calculate a mutation’s fitness. In sequential (priority) order: 1) does the patch compile without error (warnings are ignored), 2) does the mutant software run without crashing or timing out on every test case, 3) are its outputs the same as those of the original code and 4) how many L1 data cache misses does perf record.

(1) The source code is compiled using the GCC compiler (version 10.2.1) with the same options and switches, e.g. `-O3`, as the

developers (Google, Uber, the OpenCV team) use. To avoid wasting time on reporting multiple errors, `-fmax-errors=1` is used to stop GCC on the first error. If the compiler succeeds in compiling the patch, it is linked with non-evolvable code outside the patch, and the C code that calls the `perf` run time library, to form a shared object library. (The GCC options `-shared` and `-fPIC` are used to create the shared library `prog.so`.) The Python interpreter uses `python ctypes` on the shared library to call the C interface routine which calls the `perf` runtime library and the patched code. After the patch has been run, the `perf` runtime library extracts hardware counters from within the CPU, which the interface code passes back to the Python interpreter, along with the outputs generated by the patch.

(2) Both Magpie (via the Python interpreter) and the mutant itself, can signal a problem via the Unix exit status. In either case, the main Magpie thread will discard the patch giving it poor fitness and then move onto generating and testing the next patch.

(3) For each test case the Python subprocess will check that output of the patch is as expected for each user supplied test case. For example, with OLC, Magpie checks that the patched code returned the same 16 characters as Google’s code for the test’s pair of latitude and longitude. If any characters are different or missing the test fails and fitness testing for that patch stops immediately. Before using Magpie, we ran the original OLC, H3, Blue and SEEDS programs on each test case, recorded their output and then this was automatically converted into a Python list data structure. For example, with OLC, the fitness training consists of ten latitude and longitude floating point numbers and ten 16 character strings, formatted as ten Python bracketed tuples.

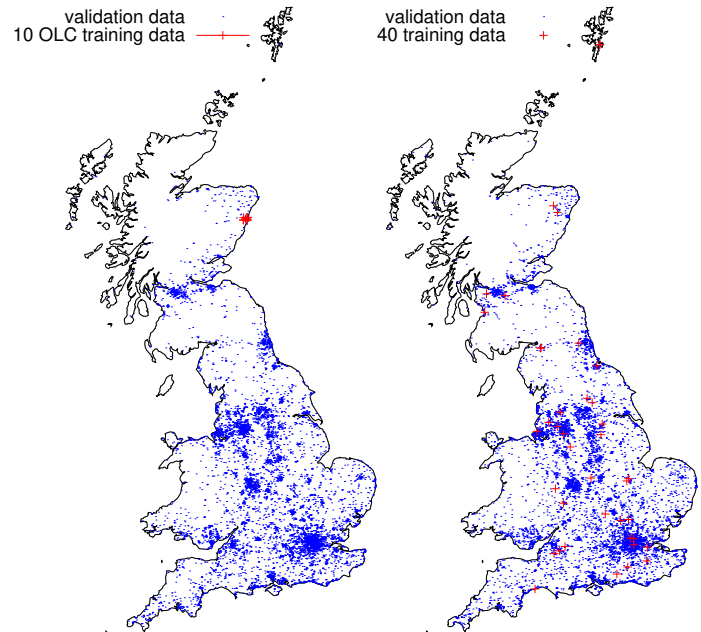
(4) Section 2.2 above describes how the `perf` C runtime library is integrated into Python. Magpie uses the first quartile (Q1) of all the patch’s repeated measurements to give its fitness measure (see also Section 2.2). Even in supposedly deterministic programs, the hardware counters for cache statistics, instructions run and elapsed time are noisy. Despite the use of robust statistics like Q1 on `perfs` L1 data cache misses, fitness remains noisy. But as we will see, in some cases Magpie is able to make progress.

### 3.1 Test Cases for Google’s OLC and Uber’s H3: GB Post Codes

We used the same test cases as before [21] when optimising OLC and H3.

Both Google’s Open Location Code (OLC) <https://github.com/google/open-location-code> (downloaded 4 August 2022) and Uber’s Hexagonal Hierarchical Geospatial Indexing System (H3) <https://github.com/uber/h3> (downloaded the previous day) are open industry standards (total sizes OLC 14 024 and H3 15 015 lines of source code). They include C programs which convert latitude and longitude into their own internal codes (see Table 1). For OLC we used Google’s 16 character coding and for H3 we used Uber’s highest resolution (`-r 15`) which uses 15 characters. Following our earlier work [21], we use as test cases the position of actual addresses.

For Google’s OLC we used the [21] dataset which was the location of the first ten thousand GB postcodes downloaded from [https://www.getthedata.com/downloads/open\\_postcode\\_geo.csv.zip](https://www.getthedata.com/downloads/open_postcode_geo.csv.zip) (dated



**Figure 1:** Left: Ten OLC training points randomly selected in the neighbourhood of Aberdeen (red). Holdout set (blue dots) GB post codes. Right: Forty training points randomly selected from ten H3 runtime classes. Holdout set (blue dots), locations of ten thousand random GB post codes (no overlap with H3 training or OLC (left) holdout data). OLC and H3 patches pass all 10 000 holdout tests. Dataset from [21].

16 March 2022). The addresses are alphabetically sorted starting with AB1 0AA, which is in Aberdeen. For training ten pairs of latitude and longitude were selected uniformly at randomly (see Figure 1). The unmutated code was run on each pair and its output saved (16 bytes). For each test case each mutant’s output is compared with the original output.

Uber’s H3 was treated similarly (see right of Figure 1)

### 3.2 Test Cases for Blue and OpenCV SEEDS

We had previously produced a simple example of the GISMO GI system for students attending the 2018 TAROT summer school on Software Testing, Verification and Validation [19] [http://www.cs.ucl.ac.uk/staff/W.Langdon/ftp/gp-code/opencv\\_gp.tar.gz](http://www.cs.ucl.ac.uk/staff/W.Langdon/ftp/gp-code/opencv_gp.tar.gz). We generated ten random images (see Figure 2) and calculated the number of “blue” pixels in each. These were used by Magpie as training data, with a goal of optimising their code to minimise L1 cache misses. Notice the training images contain  $96 \times 128 = 12\,288$  coloured pixels, occupying 49 152 bytes, and so exceed the L1 data cache. We removed the comments, leaving 100 lines of C++ code.

In contrast OpenCV is an enormous suite of C++ image processing tools. At the beginning of 2023 OpenCV’s open source repository on GitHub comprised more than two million lines of code (mostly C, C++ and XML). Therefore, we selected an important routine: the state-of-the-art OpenCV SEEDS superPixels image segmentation algorithm. This figured in the \$50K OpenCV Challenge,



Figure 2: One of ten “Blue” random  $96 \times 128$  images. Example contains 2135 blue pixels [19].



Figure 3: OpenCV SEEDS image segmentation training  $204 \times 153$  image. Downloaded from [http://www.cs.ucl.ac.uk/staff/W.Langdon/ftp/gp-code/opencv\\_gp.tar.gz](http://www.cs.ucl.ac.uk/staff/W.Langdon/ftp/gp-code/opencv_gp.tar.gz) [27].

and we had previously used it in GI experiments to reduce run time whilst respecting its API [27].

The OpenCV code of the SEEDS SuperPixel algorithm is 1500 lines of C++ code, but the important routines are held in one file `updatePixels.cpp`. After removing comments and empty lines there are 319 lines. Unfortunately the SEEDS algorithm is compute intensive and so instead of using full images obtained from [27], we reduced the training data to  $1/16$  (see Figure 3). Notice the  $204 \times 153$  training image contains 31 212 coloured pixels (124 848 bytes) and so the major data structure used by the SEEDS algorithm exceeds the L1 data cache.

## 4 MAGPIE SEARCH

Magpie allows a clean separation of user supplied parameters for each experiment from its generic code.

It is often the case that there is a “settling in” period when programs take longer to run than usual. Magpie solves this by starting with a “warm up”. Before starting any search Magpie generates empty patches and tests them. (Originally Magpie created 4 empty patches, we now use 11.) The fitness of the warm up patches are discarded, except that Magpie reports the performance of all later patches as a proportion of the average warm-up fitness. Note therefore that lower Magpie fitness are better.

In addition to the 11 warm-up patches, Magpie generated 700 OLC, 19 077 H3, 904 Blue and 3151 SEEDS patches (see Table 1). The



Figure 4: Example holdout  $2448 \times 3264$  image (tested at  $100\% \frac{1}{2}$ ,  $\frac{1}{4}$ ,  $\frac{1}{8}$  and  $\frac{1}{16}$  sizes). One of 30 SEEDS examples randomly selected from <http://www.cs.ucl.ac.uk/staff/W.Langdon/egp2014/granada/> The automatically evolved SEEDS C++ code was validate multiple times at five resolutions on each.

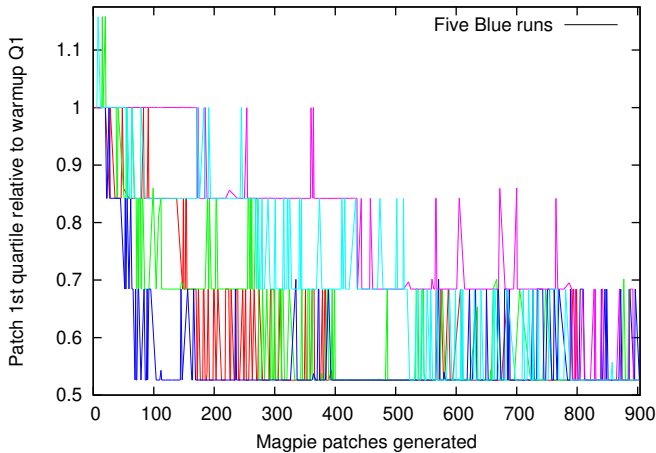
Table 1: Left: size of C/C++ sources to be optimised (comments and blank lines removed). Middle: averages for up to five Magpie runs. Columns 3–4 size of patch. Column 5 best fitness (average reduction in L1 data cache misses). Right: Column 6 size of search space explored. Column 7 fraction of mutants which compile, run ok and give correct answer. Column 8 average Magpie run times for 1 core on a 3.6 GHz Intel i7-4790 desktop with GCC 10.2.1 and Python 3.10.1

Example	LOC	Mutant		Magpie		
		size	L1D	steps	% run ok	duration
OLC	134	1–5	5%	700	23%	66 secs
H3	1615	6–18	6%	19077	33%	3.9 hours
blue	100	7–10	47%	904	30%	1.7 hours
SEEDS	319	2–7	7%	3151	2%	5.2 hours

OLC and H3 (700 and 19077) values are taken from our previous work [21]. As before [21], we use a coupon collector [11] argument to calculate how many random samples would be needed to be almost certain of visiting every line of the C/C++ source code at least once. (The H3 source code to be optimised is much bigger than the others, see Table 1, hence the larger search effort.) In all four cases we used Magpie’s run time reduction option: `python3 -m bin.magpie_runtime`.

Magpie used a single main thread on an otherwise mostly idle 32 GB 3.60 GHz Intel i7-4790 desktop CPU running networked Unix CentOS 7, using Python 3 version 3.10.1 and version 10.2.1 of the GNU C compiler. In all four cases there was a lot of variation in values recorded by perf. (For example, see Figures 7 and 8.)





**Figure 5: Blue patches. perf L1D counts for five Magpie runs relative to warm-up. The fitness (lower is better) of only patches which past all ten training cases is plotted.**

## 5 RESULTS

The results are summarised in Table 1, and Figures 5 and 6 show the fitness of Blue and SEEDS patches found as Magpie was running. Note in particular that the results in columns 3 and 4 (Mutant size and L1D) are for the best fitness found by Magpie during its runs. That is, in Table 1, the L1D improvement is as measured during training. In the cases of the two geographic programs (OLC and H3), while the patches retain their functional ability to pass up to 10 000 holdout tests, the desired improvement in cache performance did not generalise. Indeed it appears after taking care of the noise, there is no real difference in L1 data cache misses between the original and patch code. This is in great contrast to the two data rich image programs, where the patch does give reduced data cache misses on images of the same size as the training data (See Figures 9 and 10). Again both Blue and SEEDS patches also generalise in terms of still giving the correct answer on unseen images. However, the right-hand side of Figure 10 (blue  $\times$ ) shows the SEEDS patch does not give a reduction in L1 data cache misses in images more than four times larger than the training image (Figure 3).

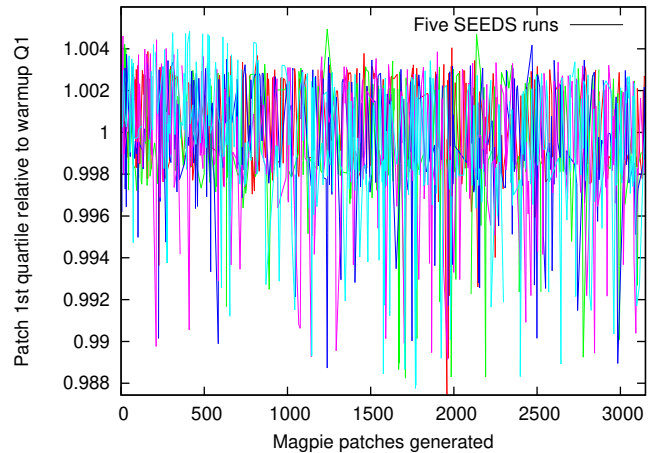
On average 72% of OLC patches fail to compile<sup>3</sup>, about 5% fail one or more test cases, while the remaining 23% pass all ten fitness tests. Table 2 summarises the statistics for all four experiments. The pattern for H3 and Blue is similar. However the tabu list used with SEEDS means whenever a patch fails a tabu check, it is marked by Magpie as if it had failed to compile. Hence the low figure for SEEDS’ ok column<sup>4</sup>. Note mostly patches which compile, run ok and pass all the tests.

Although Magpie has a nice tool for minifying patches, we did not use it due to the noisy nature of our perf based fitness measure.

In all four experiments, the final patch generated the same results as the original program. In the case of the two geographic tools

<sup>3</sup>Previously we used specialist mutation operators with LLVM IR which ensured all mutants compiled successfully to machine code [21].

<sup>4</sup>In [20] we used a similar idea to test if mutated code is identical by inspecting the X86 assembler generated by the GNU gcc compiler. Also Mike Papadakis et al. [35] compared compiler output to look for equivalent mutants.



**Figure 6: SEEDS patches. perf L1D counts for five Magpie runs relative to warm-up. To suppress clutter data poor fitness, i.e. above 1.005, are not plotted. (Lower fitness, y-axis, is better.)**

**Table 2: Summary of Magpie patches. Average of five runs on each L1 data cache experiment.**

	Compile error	Test failed	time out	too big	ok
olc	72%	5%	0%	0%	23%
h3	61%	4%	0%	2%	33%
blue	56%	13%	1%	0%	30%
seeds	87%	10%	0%	0%	2%

(OLC and H3), the holdout set contains “missing data”, i.e. postal addresses without a latitude, longitude location. In these 85 cases OLC produces a default output, while H3 aborts (with a non-zero Unix error code) and an error message. The H3 patches similarly detected and reported the error. On the other 9915 holdout locations the patch similarly returns the same output as the original H3. That is both OLC and H3 patches pass all 10 000 hold out tests. However in neither case, were we able to show the fitness seen in the Magpie runs, translated to re-running them. (See also Sections 6.1.1 and 6.1.2.)

Our results are summarised in Table 1 column 5 “L1D” which gives the percentage reduction in L1 data cache misses during training. Unfortunately, as will be discussed in the next section, the improvements in cache use reported during training with OLC and H3, did not generalise and out of sample, there is no reduction in L1 data cache misses. In contrast on both image processing examples we find significant (non-parametric Mann-Whitney test) reduction in L1 data cache misses. For the Blue benchmark it is 47%. And 1% for the patch to the OpenCV image segmentation, SEEDS, C++ code on 30 unseen unrelated images of the same size as the training image.

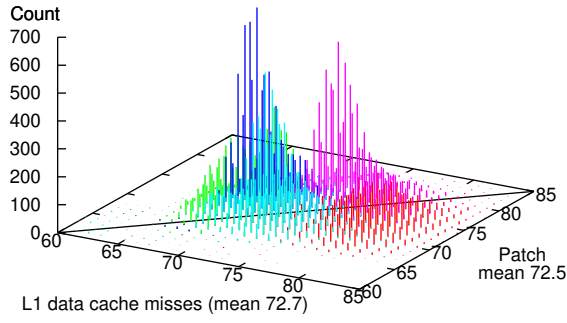


Figure 7: Magpie OLC patch sample performance on 10 000 hold out postal addresses. The patch passes all 10 000 examples. Unlike H3, OLC is quite happy to process without error the 85 invalid locations. The 10 000 tests have been repeated five times on the same computer. The presence of five distinct distributions is shown by five colours. In each repeat, the original and patch have significantly different distributions, but this is an *artifact*, see Section 6.2. Left hand side of Figure 1 shows the 10 training data location and these 10 000 holdout post codes.

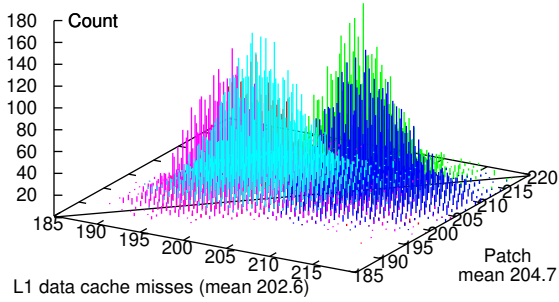


Figure 8: Magpie H3 patch out of sample performance on 9915 hold out postal addresses. The patch passes all 10 000 examples. There are 85 invalid locations (for which the patch correctly generates an error message) which are excluded from this comparison. The patch incurs on average 2.1 extra L1 data cache misses. Note spread of data  $\sigma=5$ . As with Figure 7 we use five colours to show five distinct distributions cause by repeating the 9915 tests five times. The right hand side of Figure 1 shows the 40 training data location and these 9915 holdout post codes.

## 6 DISCUSSION

### 6.1 Types of Improvement Found

**6.1.1 OLC.** In one run Magpie found a single patch which on the ten training locations gave an 8% reduction in L1 data cache misses. It simply deletes one line. The removed line is part of OLC’s command line verification. Since all the tests use well formed command lines, checks for errors in parsing the command line are never triggered. Therefore the deleted line is never used. The 8% improvement reported by Magpie appears to be just noise in the perf measurements. That is (as with H3, next section), the patch does not give sustained reductions in L1 data cache misses. However, as would be

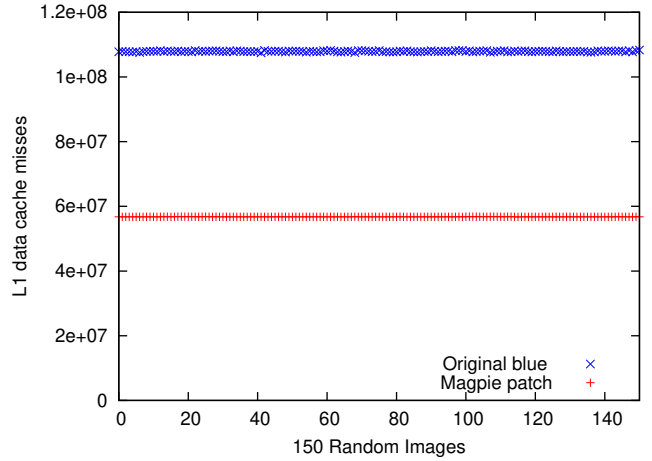


Figure 9: Magpie “Blue” patch out of sample performance on 150 hold out images. (Figure 2 contains one of the training images.) In all cases the patch returns the same count of blue pixels as the original code, but incurs only about half as many L1 data cache misses.

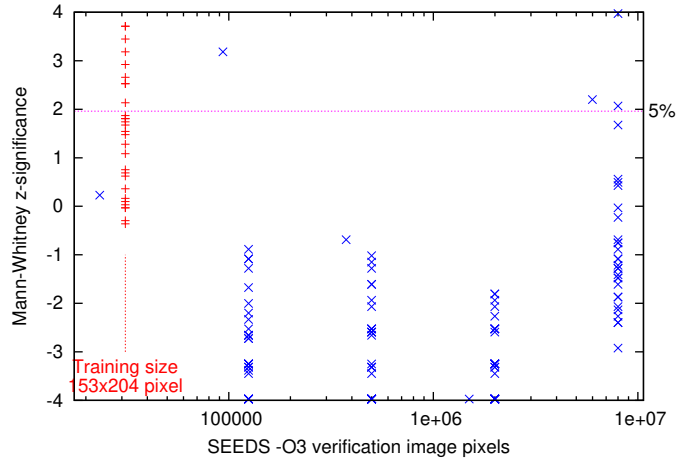


Figure 10: Magpie SEEDS patch out of sample performance on 150 hold out images of various sizes (e.g. Figure 4). In all cases the patch segmented the image identically to OpenCV. Measurement of L1 cache still very noisy even when repeated 11 times. Nevertheless the patch tends to reduce L1 data misses on 153×204 images (red +) of the same size as the training image (Figure 3) and worse on larger images (blue ×). Note log x-axis.

expected, it does generalise functionally and passes all 10 000 hold out postal addresses (shown with blue + on the left of Figure 1).

**6.1.2 H3.** In one H3 run the best patch found contains 12 changes. Two are LineReplacement, nine are LineDeletion and there is one LineInsertion, giving 14 line changes, all of them non-functional. Two delete lines of code that are never used. One includes a header file, which due to conditional compilation has no effect. Two add

lines to const arrays which are never used. Nine delete data in const arrays, in which only in two cases is the array referenced and then in all cases, only parts of the array which are unchanged are used. As would be expected, the patched code passes all 10 000 validation cases.

**6.1.3 SEEDS.** In one run Magpie found a patch with seven changes which together reduced reported L1 data cache misses by 7.4%:

Deletion 254, Deletion 204, Insertion before 102 of 105, Deletion 247, Insertion before 87 of 66, Deletion 63, Insertion before 137 of 109.

In three places GI exploited information available at run time but not to the GCC -O3 compiler by using the fact that `int seeds_prior` is always 2. Thus line `if( seeds_prior )` is always true and can be deleted and the C++ `switch( seeds_prior )` statement always takes case 2:. Therefore case 1: and case 5: can be deleted without changing the programs semantics. In fact so too can the case 3: and case 4: statements. Indeed in other runs Magpie removed the case 3: and/or case 4:. (These redundant case statements are on lines 247–249 and 254.)

These can be seen as traditional GI speed-up changes, in which unnecessary operations are removed. They may also have a potentially beneficial impact on the data cache, although only removing `if( seeds_prior )` on line 63 has a direct impact on data usage. Reducing the code size may have a subtle second order impact on the data cache by changing the machine code generated by the optimising compiler.

Deletion of line 204 does the opposite, since it removes `if( labelA != labelB )` which is only true infrequently. Thus the patch causes `update()` on the next line to be called more often. But each time it calculates the same label for pixel  $(w-1,y)$  and so does not change the image’s segmentation. In the fitness function there is no trade-off: the patch can waste as many CPU instructions as it likes, only reducing L1 data cache misses are important. The `update()` following line 204 is called in a loop of 153 iterations in which data adjacent to  $(w-1,y)$ , i.e.  $(0,y+1)$ , is needed on the next loop iteration. Therefore the unnecessary call of `update()` may have the beneficial impact of keeping data in the L1 cache<sup>5</sup>. This argument ought to hold with images of any size, but Figure 3 suggests the overall benefit to the L1 data cache does not hold each time the image size is increased by a factor of four. This may be because with larger images the other arrays `update()` uses, nullify the benefit of trying to keep the label pixel data in the L1 cache.

The other three changes insert copies of lines of code which again do not change values but cause them to be recalculated:

```
priorB = threebyfour(x, y, labelB);    (before line 87)
int a11 = Labels((y - 1) * width + (x - 1)); (line 102)
int a23 = Labels((y) * width + (x + 1));  (line 137)
```

The compiler recognises that `int a11` and `int a23` are unused, however this does not ensure it ignores them and so it may generate code to access labels for pixels  $(x-1,y-1)$  and  $(x+1,y)$ . Both occur in a doubly nested loop which scans the image in the wrong order: i.e. the inner loop samples the data a long way apart, rather than accessing neighbouring cache lines. Hence the available 512 L1 cache lines may be under great stress. For the patch adding `int a11`, it seems label  $(x-1,y-1)$  is often needed almost immediately, so

<sup>5</sup>Cache lines are 64 bytes and so can hold labels for 16 horizontally adjacent pixels.

duplicating the calculation of `int a11` perhaps costs little in terms of cache usage. Whereas the `int a23` patch reads the label for pixel  $(x+1,y)$  which has been recently read, and so again recalculating `int a23` may impose little on the cache. The nested loops makes analysis hard and there may be important data locality effects, similar those suggested for `update()` in the earlier paragraphs. These effects may totally overwhelm the computer’s 512 L1 data cache lines in bigger images.

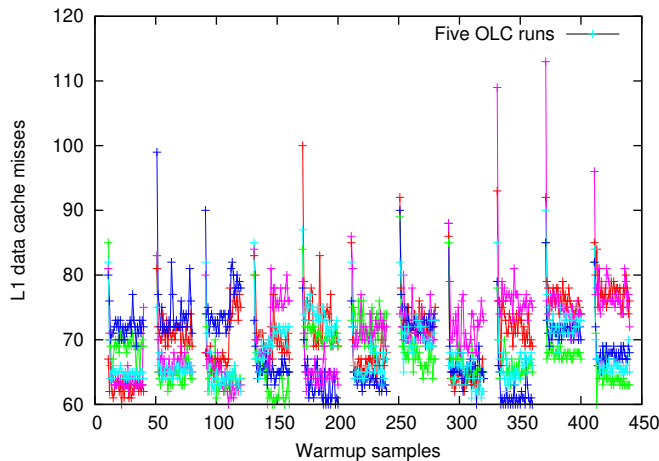
In contrast the additional call to `threebyfour(x, y, labelB);` to set `priorB` inserted before line 87, is in a pair of nested loops which accesses pixel data in the “right” order, i.e along the x-axis. If called, `threebyfour` reads four adjacent pixel labels in three adjacent rows. These will occupy between three and six cache lines (depending how the x-dimension straddles 64 byte boundaries). Like the line 204 patch discussed above, which causes `update()` to be called many more times, given the nested loops and that `threebyfour` has been recently called, the extra call to `threebyfour` may be beneficial in keeping data (which is already in the L1 cache) in it. Data near these  $3 \times 4$  pixels, and hence in the same cache lines, is almost certain to be needed immediately by the next loop iteration.

## 6.2 Non-stationary noise, Profiling, Mutation and other Search Operators

As mentioned in Section 2.2, motivated by large positive outliers often seen in run time measurement, we have used the first quartile in the fitness function, as it is a robust statistic which is relatively immune to positive outliers. Although Figures 7 and 8 do not plot the small number of large positive outliers, they show that the distribution of cache misses is very noisy and also that it is more symmetric and Gaussian like than expected. Given the apparent symmetry it may be that the median would give a more consistent fitness measure than the first quartile (see also Figure 11).

However, the multi-modal distributions shown in Figures 7 and 8 with different colours, show another little discussed problem: The L1D noise is *not* stationary [34], but subject to some unknown drift. Classical arguments, which assume multiple measurements are independent and identically distributed (IID), suggest increasing the number of measurements  $n$  will reduce the impact of noise in proportion to  $\sqrt{n}$ . However this misses the fact if the noise is non-stationary, then measurements taken during a Magpie run (or indeed during any evolutionary computing EC run) will drift. Not only during the EC run itself, but also when the evolved artifact is used. It may be we need, not only to take multiple L1D measurements per fitness evaluation, but also, during the EC run, to make estimates of the drift. Perhaps this might be done by running some known fixed example code. If online drift estimation turns out to be effective, it is likely that effort spent on combating noise during EC runs would be well worth while, as it should lead to better more robust solutions.

Magpie has targeted whole functions that could possibly be called. Particularly in our largest example, H3, this is not sufficient. Since there are both lines of code and pre-set data values that are never used, but Magpie is wasting effort on trying to optimise them. It is common in GI to profile the program to be evolved, and then to target only code that is indeed executed [24]. In these examples,



**Figure 11: L1 data cache misses for Magpie’s OLC warm-up (when the same program is repeatedly measured). Notice: the presence of outliers, that often the first measurement of a group is much bigger than the rest, and the systematic, as well as random, variation between runs of the same software on the same computer.**

profiling was not used. Indeed the presence of a huge volume of unused data in H3 hints at a further problem. Earlier profiling has concentrated upon code execution. It appears not to have considered that data might be created (and so need maintaining) which is not used by the code during every-day mundane operations.

Of course Magpie has more sophisticated syntax aware approaches and they also might benefit from profile data. Apparently more recent version of Magpie, already rule out simple patches that move `#include` files or simply swap lines containing just a single curly bracket “}”. These may help, but we fear that as usual, fitness driven evolution will find a way to exploit code changes by making other, as yet unthought of, apparently “useless” changes.

The trick of enforcing that each patch make a new semantic change by keeping a tabu list (as was used with SEEDS) will not deal with the problem of wasted effort being spent on generating patches that mutate either unused code or unused data. The tabu list it uses to prevent duplicates is based on the object file created by the compiler. Changes to either unused code or unused pre-set data would change the object file and so although useless would appear to be plausible semantic changes. The problem is exacerbated here with Magpie’s local search, as apparently “useless” changes to un-executed code or unused data can, due to the noisy fitness function, appear to be beneficial and so drive search in unproductive directions. However we need to be cautious, as the behaviour of caches is often proprietary and the exact implications of even “obviously useless” changes is in practice unknowable. For example, even when the GCC compiler (with `-O3`) issues a warning saying the patch introduces an “unused variable”, it may change the machine code it generates. So potentially changing the behaviour of the caches.

So far with Magpie we only used a few types of mutation: Line Deletion<sup>6</sup>, LineInsertion and LineReplacement. Many others are feasible. Similarly we have only used Magpie’s local search, and other strategies could be explored. Indeed Magpie already supports genetic programming. In future Magpie’s parameter search might also be applied to aspects outside the source code, such as the compilation and linking.

## 7 CONCLUSIONS

We have taken a new open source genetic improvement tool written in Python (Magpie) and applied it to industrial C source codes from Google and Uber, the “Blue” image processing problem and OpenCV’s state-of-the-art image segmentation C++ code. In particular to the never before attempted optimisation of the deepest level of the hardware data cache hierarchy, which is essential to modern computers. For OpenCV’s SEEDS a 1% reduction on compiler `-O3` optimised code was found, whilst “Blue” L1 data cache misses were almost halved (47%).

In retrospect perhaps it was optimistic to expect search to find ways to make a large impact on the data cache on our two industrial geo-positioning examples (OLC and H3). After all they have a large amount of code for only a very small data input (two floats). Perhaps more importantly a typical run only incurs a few hundred L1 data cache misses, making it a small target to improve. Nonetheless it is heartening that in a minute or a few hours<sup>7</sup> on compiler optimised code improvements were found and that these at least generalise functionally.

## ACKNOWLEDGMENTS

I am grateful to H.Wierstorff for help with gnuplot. Supported by EP/P023991/1 and the Meta OOPS project.

## REFERENCES

- [1] Nadia Alshahwan. 2019. Industrial experience of Genetic Improvement in Facebook. In *GI-2019, ICSE workshops proceedings*, Justyna Petke, Shin Hwei Tan, William B. Langdon, and Westley Weimer (Eds.). IEEE, Montreal, 1. <http://dx.doi.org/10.1109/GI.2019.00010> Invited Keynote.
- [2] Gabin An, Aymeric Blot, Justyna Petke, and Shin Yoo. 2019. PyGGI 2.0: Language Independent Genetic Improvement Framework. In *Proceedings of the 27th Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering ESEC/FSE 2019*, Sven Apel and Alessandra Russo (Eds.). ACM, Tallinn, Estonia, 1100–1104. <http://dx.doi.org/10.1145/3338906.3341184>
- [3] Gabin An, Jinhan Kim, Seongmin Lee, and Shin Yoo. 2017. PyGGI: Python General framework for Genetic Improvement. In *Proceedings of Korea Software Congress (KSC 2017)*. Busan, South Korea, 536–538. <https://coinse.kaist.ac.kr/publications/pdfs/An2017aa.pdf>
- [4] Gabin An, Jinhan Kim, and Shin Yoo. 2018. Comparing Line and AST Granularity Level for Program Repair using PyGGI. In *GI-2018, ICSE workshops proceedings*, Justyna Petke, Kathryn Stolee, William B. Langdon, and Westley Weimer (Eds.). ACM, Gothenburg, Sweden, 19–26. <http://dx.doi.org/10.1145/3194810.3194814>
- [5] Wolfgang Banzhaf, Peter Nordin, Robert E. Keller, and Frank D. Francone. 1998. *Genetic Programming – An Introduction; On the Automatic Evolution of Computer Programs and its Applications*. Morgan Kaufmann, San Francisco, CA, USA. <https://www.amazon.co.uk/Genetic-Programming-Introduction-Artificial-Intelligence/dp/155860510X>
- [6] Aymeric Blot and Justyna Petke. 2021. Empirical Comparison of Search Heuristics for Genetic Improvement of Software. *IEEE Transactions on Evolutionary Computation* 25, 5 (Oct. 2021), 1001–1011. <http://dx.doi.org/10.1109/TEVC.2021.3070271>

<sup>6</sup>Delete is commonest way programmers speed up code [9].

<sup>7</sup>Magpie with the GNU C/C++ compiler (`-O3`) processed between 0.7 and 10 patches per second depending on OLC, H3, Blue or SEEDS. Whereas with Clang 14.0.0 in our earlier work [21] on OLC and H3, we processed between 0.25 and 1.5 LLVM IR patches per second, depending on experiment and if using `-O3` or not.



- [7] Aymeric Blot and Justyna Petke. 2022. MAGPIE: Machine Automated General Performance Improvement via Evolution of Software. arXiv. <http://dx.doi.org/10.48550/ARXIV.2208.02811>
- [8] Alexander E. I. Brownlee, Justyna Petke, Brad Alexander, Earl T. Barr, Markus Wagner, and David R. White. 2019. Gin: genetic improvement research made easy. In *GECCO '19*, Manuel Lopez-Ibanez, Thomas Stuetzle, Anne Auger, Petr Posik, Leslie Pezre, Caceres, Andrew M. Sutton, Nadarajen Veerapen, Christine Solnon, Andries Engelbrecht, Stephane Doncieux, Sebastian Risi, Penousal Machado, Vanessa Volz, Christian Blum, Francisco Chicano, Bing Xue, Jean-Baptiste Mouret, Arnaud Liefooghe, Jonathan Fieldsend, Jose Antonio Lozano, Dirk Arnold, Gabriela Ochoa, Tian-Li Yu, Holger Hoos, Yaochu Jin, Ting Hu, Miguel Nicolau, Robin Purshouse, Thomas Baeck, Justyna Petke, Giuliano Antoniol, Johannes Lengler, and Per Kristian Lehre (Eds.). ACM, Prague, Czech Republic, 985–993. <http://dx.doi.org/10.1145/3321707.3321841>
- [9] James Callan, Oliver Krauss, Justyna Petke, and Federica Sarro. 2022. How Do Android Developers Improve Non-Functional Properties of Software? *Empirical Software Engineering* 27 (2022), Article 113. <http://dx.doi.org/10.1007/s10664-022-10137-2> Topical Collection: Software Performance.
- [10] Jack J. Dongarra. 2022. The Evolution of Mathematical Software. *Commun. ACM* 65, 12 (Dec 2022), 66–72. <http://dx.doi.org/10.1145/3554977>
- [11] William Feller. 1957. *An Introduction to Probability Theory and Its Applications* (2 ed.). Vol. 1. John Wiley and Sons, New York.
- [12] Stephanie Forrest, ThanhVu Nguyen, Westley Weimer, and Claire Le Goues. 2009. A genetic programming approach to automated software repair. In *GECCO '09*, Guenther Raidl, Franz Rothlauf, Giovanni Squillero, Rolf Drechsler, Thomas Stuetzle, Mauro Birattari, Clare Bates Congdon, Martin Middendorf, Christian Blum, Carlos Cotta, Peter Bosman, Joern Grahl, Joshua Knowles, David Corne, Hans-Georg Beyer, Ken Stanley, Julian F. Miller, Jano van Hemert, Tom Lenaerts, Marc Ebner, Jaume Bacardit, Michael O'Neill, Massimiliano Di Penta, Benjamin Doerr, Thomas Jansen, Riccardo Poli, and Enrique Alba (Eds.). ACM, Montreal, 947–954. <http://dx.doi.org/10.1145/1569901.1570031> GECCO 2019 10-Year Most Influential Paper Award, Best paper.
- [13] Gordon Fraser and Andrea Arcuri. 2011. EvoSuite: automatic test suite generation for object-oriented software. In *8<sup>th</sup> European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering (ESEC/FSE '11)*. ACM, Szeged, Hungary, 416–419. <http://dx.doi.org/10.1145/2025113.2025179>
- [14] Mark Harman, William B. Langdon, and Westley Weimer. 2013. Genetic Programming for Reverse Engineering. In *20th Working Conference on Reverse Engineering (WCRE 2013)*, Rocco Oliveto and Romain Robbes (Eds.). IEEE, Koblenz, Germany, 1–10. <http://dx.doi.org/10.1109/WCRE.2013.6671274> Invited Keynote.
- [15] Linsey Kitt and Myra B. Cohen. 2021. Partial Specifications for Program Repair. In *GI @ ICSE 2021*, Justyna Petke, Bobby R. Bruce, Yu Huang, Aymeric Blot, Westley Weimer, and W. B. Langdon (Eds.). IEEE, internet, 19–20. <http://dx.doi.org/10.1109/GI52543.2021.00012>
- [16] John R. Koza. 1992. *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. MIT Press, Cambridge, MA, USA. <http://mitpress.mit.edu/books/genetic-programming>
- [17] John R. Koza, Forrest H Bennett III, and Oscar Stiffelman. 1999. Genetic Programming as a Darwinian Invention Machine. In *Genetic Programming, Proceedings of EuroGP '99 (LNCS, Vol. 1598)*, Riccardo Poli, Peter Nordin, William B. Langdon, and Terence C. Fogarty (Eds.). Springer-Verlag, Goteborg, Sweden, 93–108. [http://dx.doi.org/10.1007/3-540-48885-5\\_8](http://dx.doi.org/10.1007/3-540-48885-5_8)
- [18] Oliver Krauss. 2022. Amaru - A Framework for combining Genetic Improvement with Pattern Mining. In *GECCO 2022*, Bobby R. Bruce, Vesna Nowack, Aymeric Blot, Emily Winter, W. B. Langdon, and Justyna Petke (Eds.). Association for Computing Machinery, Boston, USA, 1930–1937. <http://dx.doi.org/10.1145/3520304.3534016>
- [19] W. B. Langdon. 2018. *Genetic Improvement GISMoe Blue Software Tool Demo*. Technical Report RN/18/06. University College, London, London, UK. [http://www.cs.ucl.ac.uk/fileadmin/user\\_upload/blue.pdf](http://www.cs.ucl.ac.uk/fileadmin/user_upload/blue.pdf)
- [20] W. B. Langdon. 2020. Genetic Improvement of Genetic Programming. In *GI @ CEC 2020 Special Session*, Alexander (Sandy) Brownlee, Saemundur O. Haraldsson, Justyna Petke, and John R. Woodward (Eds.). IEEE Computational Intelligence Society, IEEE Press, internet, paper id24061. <http://dx.doi.org/10.1109/CEC48606.2020.9185771>
- [21] William B. Langdon, Afnan Al-Subaih, Aymeric Blot, and David Clark. 2023. Genetic Improvement of LLVM Intermediate Representation. In *EuroGP 2023: Proceedings of the 26th European Conference on Genetic Programming (LNCS, Vol. 13986)*, Gisele Pappa, Mario Giacobini, and Zdenek Vasicek (Eds.). Springer Verlag, Brno, Czech Republic, 244–259. [http://dx.doi.org/10.1007/978-3-031-29573-7\\_16](http://dx.doi.org/10.1007/978-3-031-29573-7_16) forthcoming.
- [22] William B. Langdon and Bradley J. Alexander. 2023. Genetic Improvement of OLC and H3 with Magpie. In *12th International Workshop on Genetic Improvement @ ICSE 2023*, Vesna Nowack, Markus Wagner, Gabin An, Aymeric Blot, and Justyna Petke (Eds.). IEEE, Melbourne, Australia. [http://www.cs.ucl.ac.uk/staff/W.Langdon/ftp/papers/langdon\\_2023\\_GI.pdf](http://www.cs.ucl.ac.uk/staff/W.Langdon/ftp/papers/langdon_2023_GI.pdf) Accepted.
- [23] W. B. Langdon and M. Harman. 2010. Evolving a CUDA Kernel from an nVidia Template. In *2010 IEEE World Congress on Computational Intelligence*, Pilar Sobrevilla (Ed.). IEEE, Barcelona, 2376–2383. <http://dx.doi.org/10.1109/CEC.2010.5585922>
- [24] William B. Langdon and Mark Harman. 2015. Optimising Existing Software with Genetic Programming. *IEEE Transactions on Evolutionary Computation* 19, 1 (Feb. 2015), 118–135. <http://dx.doi.org/10.1109/TEVC.2013.2281544>
- [25] William B. Langdon, Brian Yee Hong Lam, Justyna Petke, and Mark Harman. 2015. Improving CUDA DNA Analysis Software with Genetic Programming. In *GECCO '15*, Sara Silva, Anna I Esparcia-Alcazar, Manuel Lopez-Ibanez, Sanaz Mostaghim, Jon Timmis, Christine Zarges, Luis Correia, Terence Soule, Mario Giacobini, Ryan Urbanowicz, Youhei Akimoto, Tobias Glasmachers, Francisco Fernandez de Vega, Amy Hoover, Pedro Larranaga, Marta Soto, Carlos Cotta, Francisco B. Pereira, Julia Handl, Jan Koutnik, Antonio Gaspar-Cunha, Heike Trautmann, Jean-Baptiste Mouret, Sebastian Risi, Ernesto Costa, Oliver Schuetz, Krzysztof Krawiec, Alberto Moraglio, Julian F. Miller, Pawel Widera, Stefano Cagnoni, JJ Merelo, Emma Hart, Leonardo Trujillo, Marouane Kessentini, Gabriela Ochoa, Francisco Chicano, and Carola Doerr (Eds.). ACM, Madrid, 1063–1070. <http://dx.doi.org/10.1145/2739480.2754652>
- [26] William B. Langdon, Justyna Petke, and Ronny Lorenz. 2018. Evolving better RNAfold structure prediction. In *EuroGP 2018: Proceedings of the 21st European Conference on Genetic Programming (LNCS, Vol. 10781)*, Mauro Castelli, Lukas Sekanina, and Mengjie Zhang (Eds.). Springer Verlag, Parma, Italy, 220–236. [http://dx.doi.org/10.1007/978-3-319-77553-1\\_14](http://dx.doi.org/10.1007/978-3-319-77553-1_14)
- [27] William B. Langdon, David R. White, Mark Harman, Yue Jia, and Justyna Petke. 2016. API-Constrained Genetic Improvement. In *Proceedings of the 8th International Symposium on Search Based Software Engineering, SSBSE 2016 (LNCS, Vol. 9962)*, Federica Sarro and Kalyanmoy Deb (Eds.). Springer, Raleigh, North Carolina, USA, 224–230. [http://dx.doi.org/10.1007/978-3-319-47106-8\\_16](http://dx.doi.org/10.1007/978-3-319-47106-8_16)
- [28] Claire Le Goues, Michael Pradel, and Abhik Roychoudhury. 2019. Automated Program Repair. *Commun. ACM* 62, 12 (Dec. 2019), 56–65. <http://dx.doi.org/10.1145/3318162>
- [29] Pierrick Legend, Claire Bourgeois-Republique, Vincent Pean, Esther Harboun-Cohen, Jacques Levy-Vehel, Bruno Frachet, Evelyne Lutton, and Pierre Collet. 2007. Interactive evolution for cochlear implants fitting. *Genetic Programming and Evolvable Machines* 8, 4 (Dec. 2007), 319–354. <http://dx.doi.org/10.1007/s10710-007-9048-4> special issue on medical applications of Genetic and Evolutionary Computation.
- [30] Sherlock A. Licorish and Markus Wagner. 2022. On the Utility of Marrying GIN and PMD for Improving Stack Overflow Code Snippets. ArXiv. <https://dblp.org/rec/journals/corr/abs-2202-01490.bib>
- [31] Ibrahim Mesecan, Michael C. Gerten, James I. Lathrop, Myra B. Cohen, and Tomas Haddad Caldas. 2021. CRNRepair: Automated Program Repair of Chemical Reaction Networks. In *GI @ ICSE 2021*, Justyna Petke, Bobby R. Bruce, Yu Huang, Aymeric Blot, Westley Weimer, and W. B. Langdon (Eds.). IEEE, internet, 23–30. <http://dx.doi.org/10.1109/GI52543.2021.00014> Winner Best Paper.
- [32] Martin Monperrus. 2018. Automatic Software Repair: A Bibliography. *Comput. Surveys* 51, 1 (Jan. 2018), article no 17. <http://dx.doi.org/10.1145/3105906>
- [33] Gordon E. Moore. 1965. Cramming more components onto integrated circuits. *Electronics* 38, 8 (April 19 1965), 114–117.
- [34] David Moskowit. 2016. *Automatically Defined Templates for Improved Prediction of Non-stationary, Nonlinear Time Series in Genetic Programming*. Ph.D. Dissertation. College of Engineering and Computing, Nova Southeastern University, USA. [http://nsuworks.nova.edu/gscis\\_etd/953/](http://nsuworks.nova.edu/gscis_etd/953/)
- [35] Mike Papadakis, Yue Jia, Mark Harman, and Yves Le Traon. 2015. Trivial Compiler Equivalence: A Large Scale Empirical Study of a Simple Fast and Effective Equivalent Mutant Detection Technique. In *37th International Conference on Software Engineering (ICSE 2015)*. Florence. <http://pages.cs.aueb.gr/~mpapad/papers/ICSE15B.pdf> To appear.
- [36] Justyna Petke and Aymeric Blot. 2020. Refining Fitness Functions in Test-Based Program Repair. In *The First International Workshop on Automated Program Repair (APR@ICSE 2020)*, Shin Hwei Tan, Sergey Mechtchev, Martin Monperrus, and Mukul Prasad (Eds.). Association for Computing Machinery, internet, 13–14. <http://dx.doi.org/10.1145/3387940.3392180>
- [37] Justyna Petke and Alexander Brownlee. 2019. Software Improvement with Gin: a Case Study. In *SSBSE 2019 (LNCS, Vol. 11664)*, Shiva Nejati and Gregory Gay (Eds.). Springer, Tallinn, Estonia, 183–189. [http://dx.doi.org/10.1007/978-3-030-27455-9\\_14](http://dx.doi.org/10.1007/978-3-030-27455-9_14)
- [38] Justyna Petke, Saemundur O. Haraldsson, Mark Harman, William B. Langdon, David R. White, and John R. Woodward. 2018. Genetic Improvement of Software: a Comprehensive Survey. *IEEE Transactions on Evolutionary Computation* 22, 3 (June 2018), 415–432. <http://dx.doi.org/doi:10.1109/TEVC.2017.2693219>
- [39] Riccardo Poli, William B. Langdon, and Nicholas Freitag McPhee. 2008. *A field guide to genetic programming*. Published via <http://lulu.com> and freely available at <http://www.gp-field-guide.org.uk>. <http://www.gp-field-guide.org.uk> (With contributions by J. R. Koza).

- [40] Conor Ryan. 1999. *Automatic Re-engineering of Software Using Genetic Programming*. Genetic Programming, Vol. 2. Kluwer Academic Publishers. <http://www.amazon.com/exec/obidos/ASIN/0792386531/qid%3D943291341/102-9266197-5591202>
- [41] Marta Smigielska, Aymeric Blot, and Justyna Petke. 2021. Uniform Edit Selection for Genetic Improvement: Empirical Analysis of Mutation Operator Efficacy. In *GI @ ICSE 2021*, Justyna Petke, Bobby R. Bruce, Yu Huang, Aymeric Blot, Westley Weimer, and W. B. Langdon (Eds.). IEEE, internet, 1–8. <http://dx.doi.org/10.1109/GI52543.2021.00009>
- [42] Lee Spector, W. B. Langdon, Una-May O'Reilly, and Peter J. Angeline (Eds.). 1999. *Advances in Genetic Programming 3*. MIT Press, Cambridge, MA, USA. <http://dx.doi.org/10.7551/mitpress/1110.001.0001>
- [43] Leonardo Vanneschi and Sara Silva. 2023. *Lectures on Intelligent Systems*. Springer. <http://dx.doi.org/10.1007/978-3-031-17922-8>
- [44] Westley Weimer, Stephanie Forrest, Claire Le Goues, and ThanhVu Nguyen. 2010. Automatic program repair with evolutionary computation. *Commun. ACM* 53, 5 (June 2010), 109–116. <http://dx.doi.org/10.1145/1735223.1735249>
- [45] Westley Weimer, ThanhVu Nguyen, Claire Le Goues, and Stephanie Forrest. 2009. Automatically Finding Patches Using Genetic Programming. In *International Conference on Software Engineering (ICSE) 2009*, Stephen Fickas (Ed.). Vancouver, 364–374. <http://dx.doi.org/10.1109/ICSE.2009.5070536>
- [46] David R. White. 2017. GI in No Time. In *GI-2017*, Justyna Petke, David R. White, W. B. Langdon, and Westley Weimer (Eds.). ACM, Berlin, 1549–1550. <http://dx.doi.org/doi:10.1145/3067695.3082515>
- [47] Shengjie Zuo, Aymeric Blot, and Justyna Petke. 2022. Evaluation of Genetic Improvement Tools for Improvement of Non-functional Properties of Software. In *GECCO 2022*, Bobby R. Bruce, Vesna Nowack, Aymeric Blot, Emily Winter, W. B. Langdon, and Justyna Petke (Eds.). Association for Computing Machinery, Boston, USA, 1956–1965. <http://dx.doi.org/10.1145/3520304.3534004> Winner Best Paper.