

# Evals is not enough: why we should report wall-clock time.

John R. Woodward  
University of Stirling  
Stirling  
Scotland, United Kingdom  
jrw@cs.stir.ac.uk

Alexander E.I. Brownlee  
University of Stirling  
Stirling  
Scotland, United Kingdom  
sbr@cs.stir.ac.uk

Colin G. Johnson  
University of Kent  
Kent  
England, United Kingdom  
C.G.Johnson@kent.ac.uk

## ABSTRACT

Have you ever noticed that your car never achieves anywhere near the fuel economy claimed by the manufacturer. Does this seem unfair? Is it unscientific? Would you like the same situation to occur in Genetic Improvement (GI)? Comparison will always be a difficult business [6]. However, guidelines are discussed in [2] [4] [3]. This paper asks if reported number of evaluations, or if wall-clock time is also important, and argues that reporting time is even more important when doing GI when compared to traditional GP.

How do we fairly compare two GP systems, written in different programming languages? Counting the number of evaluations of the cost function is a fair approach. This also means you are comparing the GP systems, and not how efficiently they are implemented by humans. However, with GI we will typically compare systems which are applied to the same language (i.e. a GI systems targeted at Java, may not even be applied to a C program). GI is the process of program improvement by meta-heuristics such as Genetic Algorithms and Genetic Programming (GP), but not exclusively genetic-based operators. Importantly, unlike GP, GI starts with an existing program.

## Keywords

Genetic Improvement, Genetic Programming

## 1. POSITION

The halting problem states that we cannot in general determine if a program will halt. This poses a deep issue for GI, but is particularly important if we only count the number of evaluation. The easiest solution is to set a time-out parameter, after which, if a program has not halted, we force termination on it. However, if this parameter is set too low, it will prevent quality programs from being produced. Conversely, if this parameter is set too high, valuable wall-clock time will be wasted executing programs which do not halt. We may be able to gather some valuable information as to how to set this parameter for different test cases, given runs

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

WOODSTOCK '97 El Paso, Texas USA

© 2016 ACM. ISBN 123-4567-24-567/08/06...\$15.00

DOI: 10.475/123\_4

of the existing program. An open question is, does setting the time out parameter a little higher than the time needed for the existing program to run on a test case help GP discover new and better programs, even though they may not ultimately need the extra time. In other words, as a program is manipulated, its descendants may walk along a path through programs which require a longer runtime, but eventually lead to programs with better properties.

With the growth of online repositories, the code for a GI system can be made available to the research community, along with the test cases. This will help to independently verify published results and allow for comparisons of different GI algorithms on different hardware.

Typically we compare two metaheuristics with a fixed number of evaluations on a given set of problem instances. This is the easiest (and reasonably fair) way to compare algorithms. However, it does not take account of the amount of time to generate the next program. This could be a relatively simple process which is cheap to compute (e.g. randomly exchanging lines of code in a program). However, it could also be more computationally expensive (e.g. instrumenting the existing program to gather information about the execution of the program). This difference is ignored if we are only counting evaluations of the target program.

A possible fair comparison would be to limit the number of clock cycles (a low level measure of time). But when transferred to a different machine, GI may take a different number of clock cycles, but may be a better indication of how good a GI system is than physical time or number of evaluations of the cost function. The number of clock ticks could then be used to estimate wall-clock time.

By measuring wall-clock time, in addition to the number of evaluations, we are presenting the whole situation (i.e. the total time for GI to operate and time for the sampled programs to execute on the chosen test cases). Transferring the GI to a faster machine will result in speed up, but it is for the end user to take final responsibility for which machine they run the algorithms on. e.g. they may run GI on one hardware architecture, but later execute the genetically improved programs on a different machine with a different configuration.

We can classify programs into 4 types, depending on how they are executed.

- “1” where all nodes in the syntax tree are executed once (e.g. programs constructed with a function set  $f1$  of arithmetic operators  $\{+, -, *, \%\}$ ).
- “0-1” where nodes are either executed once or not (e.g. programs constructed with a function set  $f2$  contain-

ing logical operators {AND, OR, NOT}, where short circuiting is used.)

- “bounded” programs containing for loops with a determined number of iterations (bounded execution time).
- “unbounded” programs with while loops with an unknown termination condition (unbounded execution time).

Generally, most of GP is with the former two types of program (“1” and “0-1”), while most of GI is with the latter two types (“bounded” and “unbounded”). Adopting a GI approach, which deals with software, forces us to confront the fact that different programs can take vastly different amounts of time to execute. (of course GP work exists using Turing Complete instruction sets [5, 1], and there is no reason why GI could not be applied to programs consisting of instruction sets such as  $f_1$  or  $f_2$ ). With the first two types, programs will execute in a comparatively short amount of time (bound by the size of the program). While with the last two types, programs make take an awfully long time to terminate (possibly not halting). Therefore, reporting just the number of evaluations can be more misleading with GI when compared to traditional GP.

The number of evaluations in a GP system, and therefore a GI system, could be counted on at least 4 different levels: These being the number of:

- programs evaluated (with fitness evaluated by executing a fixed set of test cases),
- test cases evaluated (the number of test cases can vary over during the training),
- the number of nodes when a program is executed (e.g. some programs take longer to execute than others), or
- the number of nodes (weighted) when a program is executed (e.g. some instructions take longer to execute than others),

during a run. By choosing one of these methods to count evaluations over a different method, we may be able to demonstrate one GI technique is superior to another. However, choosing a different method of evaluation, our claim could be invalidated.

When a program is executed in a GI framework, there are at least 4 possible outcomes; a program

- crashes
- fails the test cases
- passes the test case
- is terminated as being possibly non-halting

One GI system may produce syntactically incorrect programs, or avoid runtime errors, while another GI system may avoid incorrect programs, and have intelligent genetic operators which make use of white-box information obtained by instrumenting the program. A GI which uses clever test case prioritization will detect earlier when to bailout and stop testing a program. Just counting the number of evaluations will not distinguish between these scenarios and is therefore a crude measure of performance. For example, a GP evolving programs that computes polynomials will always succeed and therefore, in this case, it make more sense

to count evaluations of the fitness function. Whereas, with a GI system, as the software being evolved is more complicated (taking different amounts of time to execute, failing, or not even halting), time becomes more of a pressing issue.

It may be the case that algorithm  $A_1$  runs faster than algorithm  $A_2$  on machine  $M_1$ , but on machine  $M_2$  the opposite is true. If Computing Science were treated as a natural science, we would report behaviour over a number of machines (e.g. the most popular machines). As much of GI research targets non-functional properties, GI will be aimed at multi-objective optimization. However, as we compare GI systems themselves, we are making multi-objective comparisons, making comparisons even more difficult. Just as we can cherry pick which benchmark instances we select to showcase our the performance of an algorithm, (or equivalently over-tune our algorithm on those benchmark instances), we can cherry pick architectures for GI. We should be as open as possible when making comparisons.

In conclusion, making comparisons will always be problematic. However, the situation is more difficult with GI than with GP, as the scope as to what constitutes an evaluation is broader (involving a possible program crash), and can use more time (as there is possible non-termination of programs). This paper is not claiming we should abandon comparing programs based on the number of evaluations of a fitness function, but to promote the debate and raise awareness. Counting evaluations of the cost function is sensible with a black-box setting, when the evaluations of the cost function are similar. However, we can use white-box approach with GI, calling on expensive but useful instrumentation for example, and therefore the picture is more complex. Therefore we should be as open as we can be when reporting results.

## 2. ADDITIONAL AUTHORS

## 3. REFERENCES

- [1] B. Harvey, J. Foster, and D. Frincke. Towards byte code genetic programming. In W. Banzhaf, J. Daida, A. E. Eiben, M. H. Garzon, V. Honavar, M. Jakiela, and R. E. Smith, editors, *Proceedings of the Genetic and Evolutionary Computation Conference*, volume 2, page 1234, Orlando, Florida, USA, 13-17 July 1999. Morgan Kaufmann.
- [2] J. N. Hooker. Testing heuristics: We have it all wrong. *Journal of Heuristics*, 1(1):33–42, 1995.
- [3] D. S. Johnson. A Theoretician’s Guide to the Experimental Analysis of Algorithms. In *5th and 6th DIMACS Implementation Challenges*. American Mathematical Society, 2002.
- [4] G. Kendall, R. Bai, J. Blazewicz, D. C. P., M. Gendreau, R. John, J. Li, B. McCollum, E. Pesch, R. Qu, N. Sabar, G. V. Berghe, and A. Yee. Good laboratory practice for optimization research. *J Oper Res Soc*, 67(4):676–689, Apr 2016.
- [5] P. Nordin and W. Banzhaf. Evolving Turing-complete programs for a register machine with self-modifying code. In L. J. Eshelman, editor, *Genetic Algorithms: Proceedings of the Sixth International Conference (ICGA95)*, pages 318–325, Pittsburgh, PA, USA, 15-19 July 1995. Morgan Kaufmann.

- [6] G. J. E. Rawlins. *Compared to What?: An Introduction to the Analysis of Algorithms*. Computer Science Press, Inc., New York, NY, USA, 1992.