

The Case for Open Computer Programs

Darrel C Ince¹, Les Hatton² & John Graham-Cumming³

Increasingly more and more complex products are being generated by researchers, for example software, very large data sets and movies. Such products provide publication problems for research journals. We concentrate on software and argue that anything less than release of the actual source code used is inadequate for any scientific results that depend on computation. To this end, we review some of the copious literature of ambiguity and error in software development and describe some relevant case studies in climate and seismological software. We outline some simple steps to effect necessary change.

Introduction

This is an exciting time for scientific research: larger and larger computers are enabling theories to be investigated which were thought almost impossible a decade ago, robust hardware technologies are enabling data to be collected in the most inhospitable environments, larger quantities of data are being collected and an increasingly rich set of software tools have become available to analyse computer-generated data. However, there is a problem: that of reproducibility. This article attempts to examine this problem in the context of openly available program code. Our view is that anything less than release of the actual source code used is inadequate for any scientific results that depend on computation

T

here has been a recent debate concerning the release of program source code associated with scientific experiments^{1,2,3}, with views still ranging from total release to the release only of algorithmic descriptions.

A recent example of a journal asking for code release is that of a change in policy at *Science* which has now added code to the list of items that should be supplied by an author⁴. Other journals moving towards code availability include *Geoscientific Model Development*, which asks authors to provide the paper itself and a user manual with a strong recommendation for code release, and *Biostatistics*, which has appointed an editor to assess the reproducibility of the software and data associated with an article⁵.

¹ Dept of Computing Open University, Walton Hall, Milton Keynes, UK. MK7 6AA.

²(d.c.ince@open.ac.uk)

³ School of Computing and Information Systems, Kingston University, Kingston on Thames, UK. KT1 2EE.

⁴(lesh@oakcomp.co.uk)

⁵ 83 Victoria St., London SW1H 0HW. (jgc@jgc.org)

29In contrast, support for releasing only descriptions of relevant software in the body of an article or the
30methods section is exemplified by statements such as⁶:

31*Nature* does not require authors to make code available, but we do expect a description detailed enough to allow
32others to write their own code to do similar analysis.

33This article represents a response to this policy which we believe does not address problems that arise
34within e-science. Although our arguments are focussed on the implications of this statement it is
35worth pointing out that it is symptomatic of a wider problem: the scientific community places more
36faith in the results of computation than is justified by current understanding of the central role
37ambiguity in its many forms and numerical errors play in such computations.

38Much of the debate about data and code transparency has focused on issues that are associated with
39the philosophy of science, error validation and research ethics^{7,8}. We do not wish to rehearse them, but
40make some new points concerning the need for code release that draw on computing research.

41Of primary relevance to our discussion here are dissection and understanding. Code, even if not run
42by another researcher, provides an invaluable resource for these purposes. This will resonate with any
43practising scientist faced with trying to understand a poorly specified program and any contributor to
44the open source movement⁹. Code also provides building blocks for other researchers who want to
45extend the research work that it was used to support.

46Of equal importance is the need to address reproducibility. We note that two types of reproducibility
47are relevant: direct reproducibility and indirect reproducibility. The former involves the recompilation
48and rerunning of the code on, say, a different combination of hardware and systems software, in order
49to detect numerical computation problems^{10,11} and interpretation problems found in programming
50languages¹².

51An example of indirect reproducibility occurs where a researcher may want to independently validate
52something other than the code, for example a set of equations. Here the researcher may want to check
53that implementation aspects of the original work do not invalidate the original result before starting
54time-consuming reprogramming; to extract and check the detailed assumptions made by the
55researcher; and to run their own code against the original to check for statistical validity and explain
56any discrepancies.

57

58

59Any debate over reproducibility must of course be tempered by the undeniable benefits afforded by
60the explosion of Internet facilities, raw speed and data-handling capability that has occurred as a result
61of major advances in computer technology¹³. These have presented science with a great opportunity
62to address problems which would have been inconceivable even in the recent past. It is our view,

63however, the debate over code release is *essential* to resolve as soon as possible in order to garner the
64full benefits. On their own, finer computational grids, longer and more complex computations and
65larger datasets while highly attractive to scientific researchers, do not cast light on underlying
66computational uncertainties of proven intransigence and may even accentuate them.

67

68**A Case-study of Difficult Reproducibility**

69One of us (JGC) detected two errors in software in 2009 that had been developed by the United
70Kingdom's Meteorological Office. This organization has a high standard of software development¹⁴.

71The Meteorological Office produces (in conjunction with the University of East Anglia's Climatic
72Research Unit) downloadable gridded temperature anomaly datasets known as HadCRUT and
73CRUTEM3.

74On examining the available datasets and the description of the algorithm¹⁵, a number of errors were
75identified. One set of errors were procedural, but two indicate defects in the software implementation
76of one of the algorithms in the paper.

77These latter errors highlight the difficulty in translating an English-language description (albeit with
78some formulae expressed mathematically) into software. The first error involved incorrect
79computation of 'normal values' (historical average temperatures) in a number of records in Oceania.
80The Meteorological Office confirmed the errors and issued an update to CRUTEM3.

81The other two errors relate to the calculation of station errors (an estimate of the error in any average
82temperature reading). The station error was incorrectly calculated because of two different
83programming errors.

84The forensic effort in detecting these errors was large. The only way to verify the implementation
85was to write appropriate software from scratch and then try to understand the differences which arose.

86This does not in any way reflect badly on the original authors. The rewriting simply plays the part of
87peer review and it is normal to find such errors. Indeed, the discovery of such errors in 'working'
88software is exceedingly common in all computing, even when the software has been in use for some
89very considerable time. This was emphatically demonstrated in a seminal IBM study¹⁶, demonstrating
90that fully a third of all the software failures in the study took longer than 5,000 execution *years* to fail
91for the first time.

92

93**The Failure of Code Descriptions**

94The curse of ambiguity

95Ambiguity in program descriptions has haunted software development from its earliest days. Such
96ambiguity can occur at the lexical, syntactic or semantic level¹⁷; this problem is regarded as so
97axiomatic that its avoidance or minimization is routinely taught at the undergraduate level in
98computing degrees.

99It is still a major research area in computing in which, for example, the use of tools for the detection
100of ambiguity are constructed¹⁸ the avoidance of ambiguity in major projects studied¹⁹ and analyses that
101aid the requirements specification process established¹⁷. Ambiguity in some written work is not a
102result of incompetence or bad practice, but is a natural consequence of using natural language²⁰; it is
103unavoidable.

104It should also be noted that computer scientists regard a computer program as a mathematical object: a
105whole branch of computing is devoted to the formal semantics of software. In spite of this, the
106execution of a program that manipulates the type of floating point numbers used by scientists is
107dependent on a whole host of factors outside the consideration of a program as a mathematical
108object²¹.

109

110One proposed solution to the problem in ambiguity is to pay a large amount of attention to the
111description of a computer program; perhaps expressing it mathematically or in natural language
112augmented by mathematics. However, one objection to scientists employing mathematical
113specifications is that it would require researchers to acquire skills that are only peripheral to their
114work (set theory, predicate calculus and proof methods). Perhaps worse, such a solution offers no
115guarantees of the absence of defect²². A recent study²³ that examined software diversity issues
116discovered, as a by-product, problems with the interpretation of a very small, semi-mathematical
117specification of a simple computer program.

118Even a mathematical specification or a specification consisting of mathematics and natural language
119converted into software produces results which are highly susceptible to *numerical* error; this is still
120the subject of considerable research^{24,25,26,27}.

121Errors exist even with ‘perfect’ descriptions

122Errors continue to intrude even if the description is ‘perfect’. First, there are programming errors.
123Over the years, researchers have quantified the occurrence rate of such defects in the approximate
124range 1-10 per thousand lines of source code²⁸.

125Second, there are also errors associated with the numerical properties of scientific software, where,
126for example, rounding errors can occur when a large number of computations are repeatedly executed,

127for example in weather forecasting²⁹. While there is considerable research in this area, for example in
128the area of algorithms³⁰, in the area of verification³¹ and in fundamental practice³², much of it is
129published in outlets not routinely accessed by scientists in generic journals such as *Computers and*
130*Mathematics with Applications*, *Mathematics in Computer Science* and *SIAM Journal on Scientific*
131*Computing*.

132Third, there are well-known ambiguities in some of the internationally standardised versions of
133commonly used programming languages in scientific computation¹². An alarming example of a
134particularly subtle form problem is described by Monniaux²¹:

135“More subtly, on some platforms, the exact same expression, with the same values in the same variables, and the
136same compiler, can be evaluated to different results, depending on seemingly irrelevant statements (printing
137debugging information or other constructs that do not openly change the values of variables).”

138This is known as an order of evaluation problem and many programming languages are subject to its
139wilful ways. This executional ambiguity is quite deliberate and is present to allow a programming
140language compiler more flexibility in its optimisation strategy.

141Even if you interpret the Nature policy on code release as covering algorithmic descriptions using
142pseudo code, mathematical specifications employing equations or algorithmic descriptions there are
143still the problems detailed in this section associated with numerical accuracy; moreover, there is no
144guarantee that tools such as pseudo code and other algorithmic descriptions do not give rise to
145ambiguities²³. Because of this we would regard the non-availability of code as a serious impediment to
146reproducibility.

147An exemplifying case study

148One particular case study makes the above points forcefully. The study compared 9 different
149commercial implementations of the same seismic data processing algorithms, developed
150independently³³. In this particular study, several of the above sources of ambiguity were successfully
151excluded; the same dataset was used; the signal processing algorithms used were unambiguously
152specified in mathematics; and the same programming language was used (Fortran 77). The individual
153companies used software processes which would be considered of a “high standard” in the sense
154described earlier in this article.

155The results were alarming to say the least. Approximately 200,000 lines of code were exercised in
156each of the packages in a 14 stage pipeline where the output of each stage was the input to the next.
157The signal processing algorithms used would be familiar to many scientists – Wiener deconvolution,
158acoustic wave equation solutions, Fast Fourier Transforms and numerous common statistical
159procedures.

160The initial stage involved reading 32-bit pressure data from tapes recorded in a marine environment.
161Starting with the six significant figures of single precision floating point arithmetic in the input data,
162four significant figures of agreement were excised from the final output leading to variations between
163the package results in the second and sometimes first significant figure. However these data are used
164by geologists to site extremely expensive marine drilling rigs and the particular geological features of
165interest (unconformity gas traps) required around 3 significant figures to delineate accurately.
166Perhaps even worse, the variations between the nine different outputs were shown by a process of
167feedback to be unequivocally related to programming errors *which had remained undiscovered in the*
168*respective packages for several years* and the differences between the outputs were non-random.
169Even porting exactly the same software between different architectures using the same input data lost
1702 out of 6 significant places¹¹.

171Although conducted some years ago, the study is just as relevant today. The language is still in use in
172one dialect or another in scientific research, the same software assurance procedures are still widely
173used and scientific programmers are still scientific programmers and subject to human fallibility.

174Similar latent defects also surfaced in an air-traffic control system, again of the order of 200,000 lines
175of code and developed using mathematical specifications²².

176Even when programs are very much simpler such errors remain surprisingly common. As recently as
1772010 (updated Sept 2011) Microsoft issued a warning that because of the implementation of floating
178point numbers in the popular Excel spread-sheet this ‘may affect the results of some numbers or
179formulas due to rounding and/or data truncation’ (<http://support.microsoft.com/kb/78113>,
180accessed Oct 17 2011). When document and design ambiguities are thrown into the mix, the problem
181of validation to the level we associate with conventional experiments in the natural sciences, is
182challenging to say the least.

183Perfection is no guarantee of reproducibility

184Finally, the problems reported here are not just confined to software issues, but can also arise in
185machine deployment where the results from code can diverge when hardware and software
186configurations are changed²¹. So even perfection in one's own software environment does not
187guarantee reproducibility. As a result, for true reproducibility and consistency not only would we
188urge code release, but also a description of the hardware and software environment that the program
189was executed and developed.

190

191

192Potential barriers and proposed solutions

193 There are a number of barriers to the release of code. These include a shortage of tools that package
194 up code, data and research articles; a shortage of central scientific repositories or indexes for program
195 code; an understandable lack of perception of the computational problems with scientific code leading
196 to an assumption that program descriptions are adequate (something we address in this article); and
197 the fact that the development of program code is a subsidiary activity in the scientific effort.

198 A modest proposal

199 An effective step forward would be if journals adopted a standard for declaring the degree of source
200 accessibility associated with a scientific paper. A number of simple categories illustrate the idea:

201 • **SourceCode-Full:** Full release of *all* source code used to produce the published results along with
202 the make recipes to recreate the executables and a series of regression tests to build confidence.
203 (Anybody who has used Perl modules from CPAN (<http://cpan.org>, accessed 27-09-2011) will
204 understand the value of this).

205 • **SourceCode-Partial:** Full release of source code written by the researcher accompanied by
206 associated documentation of ancillary packages used, for example commercial scientific
207 subroutine libraries.

208 • **SourceCode-Marginal:** Release of executable code and an application programming interface to
209 allow other researchers to write test cases.

210 • **SourceCode-None.** No code provided.

211 This would alert both the readers and authors of a journal article to the fact that the issue is
212 important and it would also serve to highlight the degree to which results might be reproduced
213 independently.

214 There remain however some potential stumbling blocks a number of which can easily be resolved
215 using existing facilities.

216 Intellectual Property Rights

217 Clearly if there is evidence of commercial potential such as a patent or some copyright then there is a
218 problem. It is difficult to see how to deal with this within the financial limitations of journals.
219 Perhaps the simplest solution is for a journal to flag the software as “SourceCode-None” until such
220 time as they can be reproduced, either because the software goes into the public domain or is released
221 under some free licence. No slight is meant here. It simply says that for the moment, the results are
222 not currently reproducible.

223 Limited access

224 Researchers may not have access to at least some of the software packages that are used for
225 development. This would, we suggest, not be a problem for most researchers: their institutions would
226 normally provide such software. If it was a problem then a journal could mark a publication as
227 “SourceCode-Partial”. The release would still be valuable in that it would address issues such as
228 dissection detailed in the first part of our article and would enable rewriting using other programming
229 languages.

230 Procedure

231 Adopting the simple disclosure of the availability of source code will help make it clear to the
232 readership of a journal that this is an important issue whilst also giving them an idea of the degree of
233 code release. However, we would further suggest that journals adopt a standard that specifies that
234 supplementary material supporting a research article describe each of the released modular
235 components of any software used and that editors and reviewers be empowered to include an
236 appraisal of this in their judgment about the publication potential of the article. A good example of
237 this is the way that the journal *Geoscientific Model Development* asks authors to describe their
238 program code.

239 Logistics

240 The logistics of releasing and storing code whilst maintaining a cooperative development environment
241 have been solved by the open source community in the last two decades. SourceForge
242 (<http://www.sourceforge.net/>, accessed 26-09-2011) is an excellent exemplar and already home to
243 numerous scientific packages. We would urge funding agencies to investigate and adopt similar
244 solutions.

245 Packaging

246 There are a number of tools that enable code, data and the text of the article that depends on them to
247 be packaged up. Two examples here are *Sweave* associated with the programming language R and the
248 text processing system *Latex*, and *GenePattern-Word RRS* a system specific to genomic research³⁴.

249 It is still early days however. There are some projects that are exemplars, for example Donoho and co-
250 workers at Stanford University have developed software packages that allow anyone with access to
251 the *Matlab* system to reproduce figures from their harmonic analysis articles, inspect source code,
252 change parameters and access datasets³⁶.

253

254 **A Pathway to Implementation**

255 Our thesis is that previous journal and funding body strictures relating to software implementations of
256 scientific ideas that have held in the 20th century have become obsolete.

257 We have suggested one modest solution to code availability in this article. There are a number of
258 further steps that journals, academies and educational organisations might consider taking:

259• Research funding bodies should commission research and development on tools that enable code
260 to be integrated with other elements of some scientific research such as data, graphical displays
261 and the text of an article.

262• Research funding bodies should provide meta-data repositories that describes both program and
263 data produced by researchers. The Australian National Data Service which acts as an index to data
264 held by Australian research organisations is one example of this

265• Journals should expect researchers to provide some modular description of the components of the
266 software that supports a research result; referees should be empowered to include an appraisal of
267 this as part of their reviewing task. An example of this can be seen in a recent article published in
268 *Geoscientific Model Development*³⁶.

269• Science departments should expand their educational activities into reproducibility. Clearly such
270 departments should, in the main, be teaching material that relevant to the science; however,
271 courses on statistics, programming and experimental method could be easily expanded and
272 combined to form offerings that have reproducibility at their heart.

273

274 **Acknowledgements**

275 We would like to thank David Hales of the Department of Design at the Open University and Paul
276 Piwak of the Department of Computing for pointing out some reproducibility work outside
277 computing. JGC is grateful to Ilya Goz for pointing out that there appeared to be calculation errors in
278 the CRUTEM data from the Meteorological Office

279

280 **Author Contribution**

281 DCI, LH and J G-C contributed to all aspects of this article.

282 **References**

- 2831 Barnes, S. Publish your computer code: it is good enough. *Nature* **467**, 753-753 (2010).
- 2842 McCafferty, D. Should code be released? *Comm. ACM* **53**, 16-17 (2010).
- 2853 Merali, Z. Computational science:...error. *Nature* **467**, 775-777 (2010).
- 2864 Hanson, B., Sugden, A. & Alberts, B. Making data maximally available. *Science* **331**, 649-649
287 (2011).
- 2885 Peng, R.D. Reproducible research and biostatistics. *Biostatistics* **10**, 405-408 (2009).
- 289 **This work provides a succinct and convincing argument for reproducibility. *Biostatistics***
290 **are at the forefront of ensuring that code and data are provided for other researchers.**
- 2916 Editorial. Devil in the details. *Nature* **470**, 305-306 (2011).
- 292
- 2937 Pederson, T. Empiricism is not a matter of faith. *Computational Linguistics* **34**, 465-470
294 (2008).
- 2958 Donoho, D.L. An invitation to reproducible computational research. *Biostatistics* **11**, 385-388
296 (2010).
- 2979 Raymond, E.S. *The Cathedral and the Bazaar*, Sebastopol, CA: O'Reilly (2001).
- 29810 He, Y. & Ding, C.H.Q. Using accurate arithmetics to improve numerical reproducibility and
299 stability in parallel applications. *Journl. of Supercomputing* **18**, 259-277 (2001).
- 30011 Hatton, L., Wright, A., Smith, S., Parkes, G., Bennett, P. & Laws, R. The seismic kernel
301 system-a large scale exercise in fortran 77 portability. *Software Practice and Experience* **18**,
302 301-329 (1988).
- 30312 Hatton, L. The T experiments: errors in scientific software. *IEEE Computational Science and*
304 *Engineering* **4**, 27-38 (1997).
- 30513 Hey, T., Tansley, S. & Tolle, K. (Eds) *The Fourth Paradigm: Data-Intensive Scientific*
306 *Discovery*. Seattle, WA: Microsoft Press (2009).
- 30714 Easterbrook, S.M. & Johns, T.C. Engineering the software for understanding climate change.
308 *Computing in Science and Engineering* **11**, 65-74 (2009).
- 30915 Brohan, P., Kennedy, J.J., Harris, I., Tett, S.F.B. & Jones, P.D. Uncertainty estimates in
310 regional and global observed temperature changes: a new dataset from 1850. *Journl.*
311 *Geophysical Research* **111**, doi:10.1029 (2006).
- 31216 Adams, E.N.: Optimizing preventive service of software products. *IBM Journal of Research*
313 *and Development* **28**, 2-14 (1984)

- 31417 Gervasi, V. & Zowghi, D. On the role of ambiguity in RE. In *Requirements Engineering: Foundation for Software Quality* (Eds. Wieringa, R. & Persson, A) 248-254. Berlin: Springer-Verlag (2010).
- 315
316
- 31718 Yang, H., Willis, A., De Roeck, A. & Nuseibeh, B. Automatic detection of nocuous coordination ambiguities in natural language requirements. In *Proc. of the IEEE/ACM International Conference on Automated Software Engineering*, 53-62 (2010).
- 318
319
- 32019 de Bruijn, F. and Dekkers, H. Ambiguity in natural language software requirements: a case study. In *Requirements Engineering: Foundation for Software Quality* (Eds. Wieringa, R & Persson, A.) 233-247. Berlin: Springer-Verlag (2010).
- 321
322
- 32320 van Deemter, K. *Not Exactly*. Oxford: Oxford University Press (2010).
- 32421 Monniaux, D. The pitfalls of verifying floating-point computations. *ACM Trans. on Programming Languages and Systems* **30**, 12:1-12:41 (2008).
- 325
- 32622 Pfleeger, S.L. & Hatton, L. Investigating the influence of formal methods. *IEEE Computer* **30**, 33-43 (1997).
- 327
- 32823 van der Meulen, M. J. P. and Revilla, M.A. The effectiveness of diversity in a large population of programs. *IEEE Trans. Software Engineering* **34**, 753-764 (2008).
- 329
- 33024 Revol, N. Standardized interval arithmetic and interval arithmetic used in libraries. *Proc. 3rd Intl. Congress on Mathematical Software*, 337-341 (2010).
- 331
- 33225 Rump, S.M., Ogita, T. & Oishi, S. Accurate floating point summation part 1: faithful rounding. *SIAM Jnl. of Scientific Computing* **31**, 189-224 (2008).
- 333
- 33426 Rump, S. M. Accurate and reliable computing in floating-point arithmetic. *Proc. 3rd International Congress on Mathematical Software*, 105-108 (2010).
- 335
- 33627 Badin, M., Bic, L., Dillencourt, M. & Nicolau, A. Improving accuracy for matrix multiplications on GPUs. *Scientific Programming* **19**, 3-11 (2011).
- 337
- 33828 Boehm, B, Rombach, H.D and Zelkowitz M.V. (Eds) *Foundations of Empirical Software Engineering: the Legacy of Victor R. Basili*, Berlin: Springer (2005).
- 339
- 34029 Thomas, S.J., Hacker, J.P., Desagne, M. & Stull, R.B. An ensemble analysis of forecast errors related to floating point performance. *Weather and Forecasting* **17**, 898-906 (2002).
- 341
- 34230 Pan, V. Y., Murphy, B., Qian, G., & Rosholt, R.E. A new error-free floating-point summation algorithm. *Computers and Mathematics with Applications* **57**, 560-564 (2009).
- 343
- 34431 Boldo, S & Muller, J-M. Exact and approximated error of the FMA. *IEEE Trans. Computers*. **60**, 157-164 (2011).
- 345

34632 Kahan, W. Desperately needed remedies for the undebuggability of large floating-point
347 computations in science and engineering, IFIP/SIAM /NIST Working Conference on
348 Uncertainty Quantification in Scientific Computing.

34933 Hatton, L. & Roberts, A. How accurate is scientific software? *IEEE Trans. Software*
350 *Engineering* **20**, 785-797 (1994).

35134 Mesirov, J.P. Accessible reproducible research. *Science* **327**, 415-416 (2010).

35235 Donoho, D.L., Maleki, A., Rahman, I.U., Shahram, M. & Stodden, V. Reproducible research
353 in computational harmonic analysis. *Computing in Science and Engineering* **11**, 8-18 (2009).

354 **Donoho and fellow researchers have been at the forefront of reproducibility for many years.**
355 **This article describes their work and their experiences, particularly with respect to tools for**
356 **code presentation. Note that this journal rennumbers pages for each issue, the issue number is 1.**

357

35836 Yool, A. , Popova, E.E. & Anderson, T.R. MEDUSA-1.0: a new intermediate complexity
359 plankton ecosystem model for the global domain. *Geoscientific Model Development* **4**. 381-417
360 (2011).

361 **An example of an article from a journal that proactively encourages code release. Appendix B**
362 **of the article provides a description of the software.**

363 **Methods**

364 Supplementary information including the software that was used to check the meteorological database
365 can be found linked to the online version of this article.

366