

UC Berkeley

UC Berkeley Previously Published Works

Title

What Is the Price of Open-Source Software?

Permalink

<https://escholarship.org/uc/item/4572q9gs>

Journal

The Journal of Physical Chemistry Letters, 6(14)

ISSN

1948-7185

Authors

Krylov, Anna I
Herbert, John M
Furche, Philipp
[et al.](#)

Publication Date

2015-07-16

DOI

10.1021/acs.jpcllett.5b01258

Peer reviewed

This document is confidential and is proprietary to the American Chemical Society and its authors. Do not copy or disclose without written permission. If you have received this item in error, notify the sender and delete all copies.

What is the Price of Open-Source Software?

Journal:	<i>The Journal of Physical Chemistry Letters</i>
Manuscript ID:	jz-2015-01258h
Manuscript Type:	Viewpoint
Date Submitted by the Author:	13-Jun-2015
Complete List of Authors:	Krylov, Anna; University of Southern California, Dept. of Chemistry Herbert, John; The Ohio State University, Chemistry Furche, Filipp; University of California, Irvine, Chemistry Head-Gordon, Martin; University of California, Berkeley, Chemistry Knowles, Peter; Cardiff University, School of Chemistry Lindh, Roland; Chemistry - Ångström Laboratory, Theoretical Chemistry Manby, Frederick; University of Bristol, School of Chemistry Pulay, Peter; University of Arkansas, Chemistry and Biochemistry Skylaris, Chris-Kriton; University of Southampton, Werner, Hans-Joachim; Universitaet Stuttgart, Institut fuer Theoretische Chemie

SCHOLARONE™
Manuscripts

What is the price of open-source software?

Anna I. Krylov,^{1*} John M. Herbert,^{2†} Philipp Furche,³
Martin Head-Gordon,⁴ Peter J. Knowles,⁵ Roland Lindh,⁶ Frederick R. Manby,⁷
Peter Pulay,⁸ Chris-Kriton Skylaris,⁹ and Hans-Joachim Werner¹⁰

¹*Dept. of Chemistry, University of Southern California, Los Angeles, California, USA*

²*Dept. of Chemistry and Biochemistry, The Ohio State University, Columbus, Ohio, USA*

³*Dept. of Chemistry, University of California, Irvine, California, USA*

⁴*Dept. of Chemistry, University of California, Berkeley, California, USA*

⁵*School of Chemistry, Cardiff University, Cardiff, UK*

⁶*Dept. of Chemistry–Ångström Laboratory, Uppsala University, Uppsala, Sweden*

⁷*Center for Computational Chemistry, School of Chemistry, University of Bristol, Bristol, UK*

⁸*Dept. of Chemistry and Biochemistry, University of Arkansas, Fayetteville, Arkansas, USA*

⁹*School of Chemistry, University of Southampton, Highfield, Southampton, UK*

¹⁰*Institut für Theoretische Chemie, Universität Stuttgart, Stuttgart, Germany*

June 12, 2015

The notion that all scientific software should be open-source and free has been actively promoted in recent years, mostly from the top down via mandates from funding agencies^[1] but occasionally from the bottom up, as exemplified by a recent Viewpoint in this journal.^[2] A commonly articulated rationale is that the results of scientific research funded by government grants should be free for society and that the scientific community benefits from free access. The purpose of this Viewpoint is to examine the consequences of these opinions.

What is scientific software? Modern computational chemistry software is an extremely complex product based on advanced scientific ideas (models and theories) and sophisticated algorithms that transform these ideas from equations into useful tools. The development of practical software that can be used by non-experts to solve contemporary research problems requires considerable technical effort to produce and maintain robust, efficient, and validated code. Unlike the development of, *e.g.*, a smart-phone app, where the code base is small^[3] and a relatively large community can easily write extensions and add-ons, production of scientific software involves the curation of millions of lines of source code. The complexity of this code demands long-term user and developer support to maintain its integrity and performance while keeping up with new computer architectures, fixing bugs, and adding features. Recognizing the importance of these ideas, various funding agencies in the U.S. have made “sustainable software” a key priority in the distribution of research support.^[1] Sustainability is a critical goal, but one that can be realized in various ways.

Good software is important to science. Computational chemistry software is an essential scientific instrument that facilitates discovery and innovation far beyond the laboratories in which it is created, an achievement that was recognized by the 1998 and 2013 Nobel Prizes in Chemistry.^[4] Focusing on quantum chemistry software in this Viewpoint, we note that today any chemist can (with very little training) use numerous quantum chemistry programs as teaching and research tools that aid in the design and interpretation of experiments.

*krylov@usc.edu

†herbert@chemistry.ohio-state.edu

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

A software package should be more than just a tool for end users, however; it should also be a platform to develop and test new models and algorithms. Maintaining a code base requires extensive validation, and given the complexity of modern computational methods, even testing of “pilot code” or a “proof-of-principle” implementation requires access to basic software infrastructure, *e.g.*, an integrals library, a self-consistent field procedure, efficient I/O and memory management, tools for manipulating tensors, etc. Modularity is a laudable goal, but in reality “interoperability” often comes at the expense of performance. In high-performance codes, the aforementioned components are tightly interwoven, to the extent that expert help is often required to modify key components or to develop non-standard interfaces to them. As such, the ability to innovate along either applied or theoretical lines depends crucially on the quality of the software and the availability of documentation and expert support.

As examples, consider two widely-used electronic structure programs, Q-CHEM^[5] and MOLPRO.^[6] These codes consist of ~ 5.5 and ~ 2.5 million lines of source code, respectively, written in multiple languages and each in continuous development over several decades. Q-CHEM incorporates scientific advances reported in more than 300 peer-reviewed scientific publications, while methods implemented for the first time in MOLPRO have led to 20 high-impact papers that have each been cited over 300 times. Neither code is static: more than 70 scientists are actively contributing to MOLPRO, and the Q-CHEM developer base numbers more than 100. Such agile innovation comes at a price, however. Significant effort is required to keep the code robust, efficient, and sound, and to provide the documentation that ensures the usability of new methods and the extensibility of older ones.

Software from academia is often developed with an emphasis on ideas rather than implementation, fed by the need for timely peer-reviewed journal publications that provide ongoing grant support and future jobs for graduate students. To bring new ideas to the production level, with software that is accessible to (and useful for) the broader scientific community, contributions from expert programmers are required. These technical tasks usually cannot—and generally should not—be conducted by graduate students or postdocs, who should instead be focused on science and innovation. To this end, Q-CHEM employs four scientific programmers. Other quantum chemistry codes (*e.g.*, MOLPRO,^[6] TURBOMOLE,^[7] JAGUAR,^[8] MOLCAS,^[9] PQS,^[10] and ONETEP^[11]) face the same challenges and adopt similar models to ensure sustainability.

It is important to distinguish these academically-led software ventures from purely commercial endeavors. The large majority of the code in a package like Q-CHEM is funded by the government, either through grants to academic groups or, in some cases, through technology grants to the company itself. The role of the company programmers is to enable sustainability through bug fixes, user support, release management, and the addition of features that academic developers either cannot or will not add themselves. Programmers employed by the company place emphasis on functionality, robustness, and performance, more so than scientific innovation. They are directly addressing the “reproducibility problem”.^[2] Sales revenue cannot support the entire development cost of an academic code, but it contributes critically to its sustainability. The cost that the customer pays for a code like Q-CHEM reflects this funding model: it is vastly lower than the development cost, particularly for academic customers but also for industry. It primarily reflects the sustainability cost.

Software is not data. In his Viewpoint,^[2] Gezelter argues that both software and data should be open, yet it is important not to conflate the two. Software is not data, and simply because it is *feasible* to put software on the internet does not imply that it *should* be posted. Software is a product that contains an intellectual component (models and algorithms) but owes its existence to additional technical efforts. Such efforts include implementation of minor but useful (or requested) features that are not publishable in the peer-reviewed literature. This is not to say that details should be withheld as proprietary information. Just as models and algorithms should be described in full detail in scientific publications, so too should implementation details be specified, along with performance metrics (timings and scaling data) and benchmarks (energies and other computed properties). Nevertheless, the software itself is a product, not a scientific finding, more akin to, say, an NMR spectrometer—a sophisticated instrument—than to the spectra produced by that instrument.

Consider an analogy from the field of photovoltaics. Scientific findings concerning the mechanistic details of charge generation and exciton propagation in a given material are results that merit discussion in the

1
2
3 peer-reviewed literature. However, creating a new solar panel based on this research requires significant
4 additional engineering effort, which is most commonly conducted in an industrial setting. This is a common
5 mechanism for technology transfer, by means of which society benefits from academic research. Likewise,
6 new telecommunication technologies, information storage media, computer chips, etc., are products that
7 build upon—but are not equivalent to—scientific findings. Going from a journal article to a product in one’s
8 home or office requires a significant investment of resources that is often impossible to achieve in the absence
9 of a commercial platform. Software is not different.
10

11 ***There is no free software.*** The creation of scientific software is a labor-intensive process, and its
12 support and curation even more so. The question is, how do we pay for these labor costs? The answer is
13 clear in the case of commercial software, where license fees are used to defray the costs of development and
14 support. In this model, users buy the software that fits their research needs and affords them the highest
15 productivity. The decision mechanism and price-versus-deliverables choices are similar to those faced when
16 purchasing a computer or other lab equipment. Just like a new laser system may enable the pursuit of new
17 science, software that offers a competitive advantage is a sensible investment of research funds.

18 Interestingly, Gezelter^[2] specifically mentions the Quantum Chemistry Program Exchange (QCPE), an
19 early repository for open-source software, as having contributed greatly to the growth of the field. It is
20 therefore telling to note that QCPE was supported initially by the Air Force Office of Scientific Research
21 and then by the National Science Foundation, before subsequently becoming a fee-for-software service with
22 paid employees to do the time-consuming work of testing, documenting, and distributing the contributed
23 programs.^[12] Gezelter acknowledges the cost of maintaining scientific software and suggests alternative
24 models to defray these costs including selling support, consulting, or an interface, all the while making
25 the source code available for free.^[2] These suggestions strike us as naïve, something akin to giving away
26 automobiles but charging for the mechanic who services them. Such a model creates a financial incentive to
27 release a less-than-stellar product into the public domain, then charge to make it useful and usable. It is
28 better to release a top-of-the-line product, for a nominal fee.

29 Is “free” software genuinely free of charge to individual researchers? Consider software developed in the
30 U.S. national labs. These ventures are supported by full-time scientific programmers employed specifically
31 for the task, and the cost to support and develop these products is subtracted from the pool of research
32 funding available to the rest of the community. The individual researcher pays for these codes, in a sense,
33 with his rejected grant proposals in times of lean funding. In contrast to using one’s own performance metrics
34 to guide software purchases, within this system one has no choice in what one pays for. In other words, “free
35 software” is not free for you; the only sense in which it is “free” is that you are freed from making a choice
36 about how to spend your research money.

37 Computational chemistry software must balance the needs of two audiences: users, who gauge their pro-
38 ductivity based on the speed, functionality, and user-friendliness of a given program; and developers, who
39 may be more concerned with whether the structure “under the hood” provides an environment that fosters
40 innovation and ease of implementation. As a quantitative example, consider that the cost of supporting
41 a postdoctoral associate (salary plus benefits) is perhaps \$4,800/month. If the use of well-supported com-
42 mercial software can save two weeks of a postdoc’s time, then this would justify an expense of \gtrsim \$2,000 to
43 purchase a software license. This amount exceeds the cost of an academic license for many computational
44 chemistry programs. Given the choice between a free product and a commercial one, a scientist should make
45 a decision based on her own needs and her own criteria for doing innovative research.

46 ***What is “open source”?*** The term “open source” is ubiquitous but its meaning is ambiguous. Some
47 codes are “free” but are not open,^[13] while others make the source code available albeit without binary
48 executables so that responsibility for compilation and installation is left to the user. Insofar as the use
49 of commercial quantum chemistry software is a mainstay of modern chemical research and teaching, there
50 exists a broad consensus that the commercial model offers the stability and user support that the community
51 desires. Strict coding guidelines can be enforced within a model where source code access is limited to
52 qualified developers, and this kind of stability offers one counterbalance to the “reproducibility crisis”.^[2]
53 To the extent that such a crisis exists, it has occurred *in spite of* the existence of open-source electronic
54 structure codes such as GAMESS,^[14] NWChem,^[15] and CP2K.^[16]
55
56
57
58
59
60

1
2
3 Occasionally the open-source model is touted on the grounds that one can use the source code to learn
4 about the underlying algorithms, but this hardly seems relevant if the methods and algorithms are published
5 in the scientific literature. Source code itself rarely constitutes enjoyable reading, and using source code to
6 learn about an algorithm is a last resort forced by poorly written scientific papers. Better peer review is a
7 more desirable solution.
8

9 A more practical use of openly available source code is to reuse parts of it in other programs, provided
10 that the terms of the software license allow this. Often they do not. Some ostensibly “open” chemistry
11 codes forbid reuse, or even redistribution.^[13,17] Others, such as CP2K,^[16] use the restrictive General Public
12 License^[18] that requires any derivative built on the original code to be open-source itself. Variation in design
13 structure from one program to the next also severely hampers transferability, even if the license terms are
14 amenable.

15 Access to source code allows developers to introduce their own innovations, but this is distinct from
16 reuse. To facilitate innovation by developers, source code needs only to be available to people who intend
17 to build upon it. This is commonly accomplished in the framework of “closed-source” software projects by
18 granting academic groups access to the source code for development purposes. Given the large number of
19 developers for many quantum chemistry codes, this developer community should not be envisaged as some
20 small, secluded cabal but rather as a rich, diverse community of academic scientists.

21 Let us analyze accessibility using Q-CHEM as a specific example. This is commercial software whose source
22 code is not in the public domain, because that would eliminate the product that Q-Chem, Inc., is selling.
23 Consider, however, how many developers benefit from Q-CHEM as a platform for their own innovation: this
24 community exceeds 100 scientists from at least twelve countries,^[19] for whom the code is open. We call this
25 model *open teamware*.^[5] Moreover, any licensed user can obtain access to the source code upon signing a
26 non-disclosure agreement. Given the size of the user base, this is likely a significantly larger group than
27 the number of people who care to look at many existing open-source packages. Does an “open-source” code
28 that serves just a few people offer more benefit to the scientific community than a “closed-source” code
29 that fosters a community of 100+ active developers and thousands of users? What would be the impact
30 on computational chemistry of destroying other teamware projects such as MOLPRO,^[6] TURBOMOLE,^[7]
31 JAGUAR,^[8] MOLCAS,^[9] PQS,^[10] or ONETEP,^[11] in the interest of satisfying some “open-source” mandate?

32 ***Practical consequences of an open-source mandate.*** One of the pillars of science funding in the
33 U.S. and elsewhere is a merit-based funding model that distributes resources based on intellectual merit,
34 productivity, and impact. In the long run, more-competitive ideas are selected over less-competitive ones,
35 and investigators are rewarded for a track record of productivity. Research that is judged to have higher
36 impact is ranked as more meritorious and more deserving of support. This model has proven successful in
37 fostering innovation and discovery, so it is worth considering what consequences a blanket requirement that
38 software be free and open-source might engender.

39 Such a requirement would, in our view, detract from the merit-based review process. When evaluating
40 grant proposals that involve software development, the questions to be asked should be:

- 41 1. What will be the quality of the software in terms of the new science that it enables, either on the
42 applications side or on the development side?
- 43 2. How will the software foster productivity? For example, how computationally efficient is it for a given
44 task? How usable will the software be, and how quickly will other scientists be able to learn to use it
45 for their own research?
46
47

48 A rigid, mindless focus on an open-source mantra is a distraction from these more important criteria. It
49 can even be an excuse to ignore them, and creates an uneven playing field in which developers who prefer
50 to work with a commercial platform are put at a disadvantage and potentially forced to adopt less efficient
51 practices.

52 Whereas Gezelter^[2] suggests that online code repositories would help young researchers to showcase
53 their work, an open-source *requirement* actually puts young researchers at a particular disadvantage. The
54 demands of tenure and promotion place special burdens to bring ideas before the community rapidly, yet also
55 to secure funding quickly. Open-source requirements potentially force a scientist to choose between pursuing
56
57

1
2
3 a funding opportunity versus implementing an idea in the quickest, most efficient, and highest-impact way.
4 A strictly open-source environment may furthermore disincentivize young researchers to make new code
5 available right away, lest their ability to publish papers be short-circuited by a more senior researcher with
6 an army of postdocs poised to take advantage of any new code. This would contribute directly to the scenario
7 that Gezelter wishes to avoid, namely, one where students leave behind “orphaned” code that will never be
8 incorporated into mainstream, production-level software. Viewed in these terms, an open-source mandate
9 degrades, rather than enhances, cyberinfrastructure.

10
11 How should the impact of software be measured? Scientific publications are a more sound metric than
12 either the price of a product or whether its source code is available in the public domain. Software is
13 meant to serve scientific research, in the same way that any other scientific instrument is intended. As such,
14 the question should not be whether software is free or open source, but rather, what new science can be
15 accomplished with it? Let us not allow political rhetoric to dictate how we are to do science. Let different
16 ideas and different models (including open source!) compete freely and flourish, and let the community focus
17 instead on the most important metric of all: what is good for scientific discovery.

18 19 Acknowledgments

20
21 The authors thank Axel Becke, Todd Martínez, and John Stanton for useful comments and supportive
22 feedback. A.I.K., J.M.H., and M.H.-G. are part owners of Q-Chem, Inc. P.P. is a part owner of PQS, Inc.
23 F.F. has an equity interest in Turbomole, GmbH, and serves as a Scientific Coordinator. The terms of this
24 arrangement have been reviewed and approved by the University of California, Irvine, in accordance with
25 its conflict of interest policies. H.J.W., P.J.K., and F.R.M. are authors of MOLPRO. R.L. is a lead developer
26 of MOLCAS. C.-K.S. is a founding author of ONETEP.

27 28 29 References

- 30
31 [1] National Science Foundation, “A Vision and Strategy for Software for Science, En-
32 gineering, and Education: Cyberinfrastructure Framework for the 21st Century”
33 (http://www.nsf.gov/publications/pub_summ.jsp?ods_key=nsf12113) and “Software Infrastruc-
34 ture for Sustained Innovation” (http://www.nsf.gov/funding/pgm_summ.jsp?pims_id=503489&org=NSF);
35 Department of Energy, “Scientific Discovery through Advanced Computing”
36 (<http://science.energy.gov/ascr/research/scidac>).
- 37
38 [2] Gezelter, J. D. Open Source and Open Data Should Be Standard Practices. *J. Phys. Chem. Lett.* **2015**,
39 *6*, 1168–1169.
- 40
41 [3] <http://www.informationisbeautiful.net/visualizations/million-lines-of-code>.
- 42
43 [4] http://www.nobelprize.org/nobel_prizes/chemistry/laureates/1998;
44 http://www.nobelprize.org/nobel_prizes/chemistry/laureates/2013.
- 45
46 [5] Krylov, A. I.; Gill, P. M. W. Q-Chem: An Engine for Innovation. *WIREs Comput. Mol. Sci.* **2013**, *3*,
47 317–326.
- 48
49 [6] Werner, H.-J.; Knowles, P. J.; Knizia, G.; Manby, F. R.; Schütz, M. Molpro: A General-Purpose
50 Quantum Chemistry Program Package. *WIREs Comput. Mol. Sci.* **2012**, *2*, 242–253.
- 51
52 [7] Furche, F.; Ahlrichs, R.; Hättig, C.; Klopper, W.; Sierka, M.; Weigend, F. Turbomole. *WIREs Comput.*
53 *Mol. Sci.* **2014**, *4*, 91–100.
- 54
55 [8] Bochevarov, A. D.; Harder, E.; Hughes, T. F.; Greenwood, J. R.; Braden, D. A.; Philipp, D. M.;
56 Rinaldo, D.; Halls, M. D.; Zhang, J.; Friesner, R. A. Jaguar: A High-Performance Quantum Chemistry
57 Software Program with Strengths in Life and Materials Sciences. *Int. J. Quantum Chem.* **2013**, *113*,
58 2110–2142.

- 1
2
3 [9] Aquilante, F.; De Vico, L.; Ferré, N.; Ghigo, G.; Malmqvist, P.-Å.; Neogrády, P.; Pedersen, T. B.;
4 Pitonak, M.; Reiher, M.; Roos, B. O.; Serrano-Andrés, L.; Urban, M. Veryazov, V.; Lindh, R. MOLCAS
5 7: The Next Generation. *J. Comput. Chem.* **2010**, *31*, 224–247.
6
7 [10] Baker, J.; Wolinski, K.; Malagoli, M.; Kinghorn, D.; Wolinski, P.; Magyarfalvi G.; Saebo, S.;
8 Janowski, T.; Pulay P. Quantum Chemistry in Parallel with PQS. *J. Comput. Chem.* **2009**, *30*, 317–335.
9
10 [11] Skylaris, C.-K.; Haynes, P. D.; Mostofi, A. A.; Payne, M. C. Introducing ONETEP: Linear-Scaling
11 Density Functional Simulations on Parallel Computers. *J. Chem. Phys.* **2005**, *122*, 084119:1–10.
12
13 [12] Boyd, D. B. Quantum Chemistry Program Exchange, Facilitator of Theoretical and Computational
14 Chemistry in Pre-Internet History. Chapter 8 of *Pioneers of Quantum Chemistry*, ACS Symposium
15 Series vol. 1122, pp. 221–273 (2013).
16
17 [13] See, *e.g.*, the ORCA license: <https://orcaforum.ccc.mpg.de/license.html>.
18
19 [14] Schmidt, M. W.; Baldrige, K. K.; Boatz, J. A.; Elbert, S. T.; Gordon, M. S.; Jensen, J. H.; Koseki, S.;
20 Matsunaga, N.; Nguyen, K. A.; Su, S.; Windus, T. L.; Dupuis, M.; Montgomery Jr., J. A. The General
21 Atomic and Molecular Electronic Structure System. *J. Comput. Chem.* **1993**, *14*, 1347–1363.
22
23 [15] Valiev, M. Bylaska, E. J.; Govind, N.; Kowalski, K.; Staatsma, T. P.; van Dam, H. J. J.; Wang, D.;
24 Nieplocha, J.; Apra, E.; Windus, T. L., de Jong, W. A. NWChem: A Comprehensive and Scalable
25 Open-Source Solution for Large Scale Molecular Simulations. *Comput. Phys. Commun.* **2010**, *181*,
26 1477–1489.
27
28 [16] Hutter, J.; Iannuzzi, M.; Schiffmann, F.; VandeVondele, J. CP2K: Atomistic Simulations of Con-
29 densed Matter Systems. *WIREs Comput. Mol. Sci.* **2014**, *4*, 15–25. For the CP2K license, see:
30 <http://www.cp2k.org>.
31
32 [17] See, *e.g.*, the GAMESS license: http://www.msg.ameslab.gov/gamess/License_Agreement.html.
33
34 [18] <http://www.gnu.org/licenses/gpl-3.0.en.html>
35
36 [19] Shao, Y.; Gan, Z.; Epifanovsky, E.; Gilbert, A. T. B.; Wormit, M.; Kussmann, J.; Lange, A. W.;
37 Behn, A.; Deng, J.; Feng, X.; et al. Advances in Molecular Quantum Chemistry Contained in the
38 Q-Chem 4 Program Package. *Mol. Phys.* **2015**, *113*, 184–215.
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60