

# The Changing Nature of Computational Science Software

Huy Tu, Rishabh Agrawal, Tim Menzies

Computer Science, NC State, USA

hqtu@ncsu.edu, ragrawa3@ncsu.edu, timm@ieee.org

## ABSTRACT

Recently, the use of software-based Computational Science for scientific discovery has increased. How should SE adapted for Computational Science? To answer that question, we need to understand more about the nature of Computational Science software.

A recent trend in that community is that, increasingly, those codes are being stored in public domain repositories such as Github. Hence, it is now possible to explore the nature of Computational Science projects using that public domain data. Accordingly this paper we seek quantitative evidence (from dozens of Computational Science projects housed in Github) for 13 previously published conjectures about scientific software development in the literature. In all, we explore three groups of beliefs about (1) the nature of scientific challenges; (2) the implications of limitations of computer hardware; and (3) the cultural environment of scientific software development. We find that four cannot be assessed with respect to the Github data. Of the others, only three can be endorsed.

Our conclusion will be that the nature of Computational Science software development is changing. Hence the tools we should develop to support software-based-science, needs to change as well.

## KEYWORDS

Computational Science, Software Engineering

### ACM Reference Format:

Huy Tu, Rishabh Agrawal, Tim Menzies. 2020. The Changing Nature of Computational Science Software. In *Proceedings of The 28th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE 2020)*. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

Computational Science (hereafter, CS) uses software to explore astronomy, astrophysics, chemistry, weather Modeling Assumptions, economics, genomics, molecular biology, oceanography, physics, political science, and many engineering fields Computational Science (hereafter, CS) is becoming more dependent on software. For example, in 2013 a Nobel Prize went to chemists using computer models to explore chemical reactions during photosynthesis. In the press release of the award, the Nobel Prize committee wrote:

*Today the computer is just as important a tool for chemists as the test tube.*

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).

ESEC/FSE 2020, 8 - 13 November, 2020, Sacramento, California, United States

© 2020 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/Y/MM...\$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

If we can better understand the nature of the software development practices in CS then we would be (1) better able to improve those methods; and (2) alleviate some of the effort involved in sustainability, verifiability, reproducibility, understanding, and utilization of that software. Improving (say) the maintainability of CS code is desirable since this would improve:

- The adaptation of scientific projects simulations to new and efficient hardware (multi-core and heterogeneous systems);
- The ability for larger teams to co-ordinate (through integration with interdisciplinary teams);
- How well we can model complex phenomena.

Note that we are not the first to stress the need for better SE for CS. The quality of scientific software through the “Climategate” scandal [23] uncovered a lack reproducibility of CS results. Improving the verifiability of CS code would hence increase the credibility of results from CS research.

Table 1 lists some of the prior results where empirical software engineering researchers have explored computational science (this table comes from the work of Carver, Heaton, Basili, and Johanson [3, 6, 8, 13, 16], and others). Johanson et al. [16] argues that SE practices will only be integrated into CS if the honor the 13 beliefs listed in Table 1.

Given the importance of CS and prominence of the beliefs of Table 1, we assert that it is wise to test those beliefs. A recent trend is that CS researchers store their code on opens source repositories (such as Github). This paper mines the code and comments of dozens of those repositories, to assess the current relevance of the Table 1 beliefs. For each such belief:

- We predict what would be observable (if the belief was true).
- We endorse/doubt that belief if that observation is present/absent within the Github CS repositories.

An important part of this analysis is that all our reasoning is repeatable/ refutable/ improvable. One issue with much of the analysis of Table 1 is that that analysis is often applied to just a handful of projects (or even, just one) and is often qualitative in nature. Hence, that prior work does not lend itself to being double-checked by subsequent studies that survey numerous projects. Our results, on the other hand, scales to multiple projects and can be tested by anyone with access to Github.

In summary, based on an analysis of 59 CS projects, we find that:

- Four of the beliefs cannot be assessed using Github data.
- Of the other beliefs, we find we can endorse three of them.
- But we must doubt six of them.

This is not to say that six of the beliefs of Table 1 are wrong. Rather we would say that the nature of computational science is something that is ever evolving and that old beliefs needs to be rechecked whenever new data is available. The nature of CS software development is changing. Hence we should also evolve the tools we use to support software-based-science.

Category	Characteristics	Citations	Conclusion
1. nature of scientific challenge	a) Requirements are Not Known up Front	[3, 8, 10, 30, 33]	Endorse
	b) Verification and Validation are Difficult and Strictly Scientific	[3, 7, 8, 19, 25]	Doubt
	c) Overly Formal Software Processes Restrict Research	[8, 10, 31, 33]	No-Evidence
2. limitations of computer hardware	a) Development is Driven and Limited by Hardware	[10, 36]	No Evidence
	b) Use of “Old” Programming Languages and Technologies	[3, 8, 14, 25, 26]	Doubt
	c) Intermingling of Domain Logic and Implementation Details	[36]	Endorse
	d) Conflicting Software Quality Requirements	[3, 7, 8]	No Evidence
3. limitations of cultural differences	a) Few Scientists are Trained in Software Engineering	[3, 6, 10, 28, 32]	Doubt
	b) Different Terminology	[10, 21, 36]	Endorse
	c) Scientific Software in Itself has No Value But Still It is Long-Lived	[10, 21, 32, 36]	Doubt
	d) Creating a Shared Understanding of a “Code” is Difficult	[7, 15, 28, 31]	Doubt
	e) Little Code Reuse	[3, 7, 25, 31]	No Evidence
	f) Disregard of Most Modern Software Engineering Methods	[6, 7, 12, 21, 25]	Doubt

**Table 1: Thirteen beliefs from prior studies about Computational Science. From Johanson et al. [16]. These beliefs divide into the three categories shown in the left-hand column. In the far right column, anything marked as “no evidence” refers to beliefs we could not check using our Github data.**

## 2 PRELIMINARIES

### 2.1 Why Study CS (Computational Science)?

There are good reasons for the growing dependency of science on computational methods software, not the least of which is that it is faster, cheaper, and safer to explore software models than actually explore the physical effects they represent: For example, CS software can explore thousands of hurricanes scenarios without risk to human life or property.

We assert that it is important to study CS software since that software has a widespread social impact. For example:

- Weather forecasts generated from CS software can predict the path of hurricanes. This, in turn, allows (e.g.) effected home owners to better protect themselves from damaging winds.
- CS explores the properties of new materials. Synthesizing new materials is very expensive so standard practice is to use software to determine those properties (e.g. via a finite element analysis). This, in turn, enables (e.g.) the faster transition of new materials to industry.
- Moreover, more quality software would guarantee the CS work more credible and more reproducible. Therefore, better SE improves computational science software, which would lead to better (e.g.) weather Modeling Assumptions and the faster creation of new industries based on new materials.

### 2.2 Data Collection

To check our beliefs on CS projects, we proceeded as follows. Using our contacts in the CS community (from the Molecular Sciences Software Institute (MOLSSI), and the Science Gateways Community institute (SGCI)) we found 678 CS projects. Researchers warn against using all the Github data [1, 4, 18, 24] since many of these projects are simple one-person prototypes. Following their advice, we applied the sanity checks of Table 3 to select 59 projects with sufficient software development information (listed in Table 2).

Figure 4 summarizes the projects that past our sanity checks.

### 2.3 Labelling

When code is share within a software repository, an important event is the *commit comments*. These comments are all the remarks developers make to document and justify some update to the code. Code repository systems such as Github store tens of millions of such comments. These remarks are a rich source of information about a project.

In order to understand scientific development process, we manually categorised the commit comments seen within CS software. Using the power of free pizza, we assembled a team of 10 computer science graduate students That team spent 320 hours (in total) categorizing the commit comments from the projects of Table ??.

In order to allow other researchers to reproduce this work, we set the following a resource limit on this analysis. According to Tu et al. [35], two humans can manually read and categorize and cross-check 400 commit comments per day (on average). Hence, for this study, for each each project, we categorized 400 commits (selected at random) from each project.

Our population of reviewers labelled commits using the following guidelines:

- *Science enhancement*: any core science (e.g. an equation of Pascal triangle) that is being implemented or modified.
- *Engineering enhancement*: any other enhancements that related to code complexity (e.g. data structures, data types, input/output, etc)
- *Bug fixes*: Fixing software faults reported or found within the development.
- *Testing*: evaluate the functionality of a software application (e.g. scientific calculations to output/input formats).
- *Other*: not core changes, e.g. renaming or formatting changes

Each commit was labelled by two reviewers, neither of which had access to the other’s labels. Moreover, the reviewers not only looking at the commit message but also the code contribution associated with the commit (e.g. to determine if the nature of some enhancement was “scientific” or “engineering” in nature). The level of labelling disagreement was low (just 19%). When labels disagreed, the commit was given to our most experienced reviewer who made an executive decision about what was the correct label.

**Table 2: CS projects that satisfy the sanity checks of Table 3. This list has been audited by a domain expert from CS (Dr Robert Sinkovits, San Diego SuperComputer Center (<https://www.sdsc.edu/sinkovit/>)) who commented that many of these projects account for the majority of the supercomputer usage in computational science. While some of these focus on computational chemistry, they also include numerous widely-used support tools (e.g. elasticsearch) or simulation tools that are cross-disciplinary (e.g. the classical simulation tools used by molecular biologists). Also, there are also tools here used in material science (e.g. LAMMPS).**

	Language	# Developers	Duration(Years)	# Commits	# Stars	# Issues	# Releases	Analyzed in Figure 2
BLIS	C	20	4.8	1242	413	142	25	
cctools	C	43	5.5	8881	72	666	159	
keplerproject	C	22	5.25	329	26	66	18	
AMBER	C++	12	4	838	3	24	3	✓
changa	C++	19	3.5	1458	13	16	8	
cyclus	C++	20	6.5	6579	36	625	47	
dealii	C++	120	18	41514	382	1604	26	
GooFit	C++	12	5.5	1508	57	50	12	
HooMD-blue	C++	36	2.5	945	5	33	25	✓
irods	C++	36	5	6267	236	3820	34	
LAMMPS	C++	74	5.13	1581	38	29	91	✓
LIBMESH	C++	55	6	1713	24	44	59	✓
MADNESS	C++	31	4.5	5193	71	184	3	
metpy	C++	34	7.75	2199	332	481	20	
OpenMm	C++	38	7	5838	324	959	22	
OpenMX	C++	11	5	6993	26	79	62	
PCMSolver	C++	8	4	184	1	8	16	✓
PLUMED	C++	23	5.5	8075	92	282	35	
Psi4	C++	79	5.5	12178	247	504	7	
SCIRun	C++	18	6.5	8887	45	1487	79	
TRILINOS	C++	179	3	79520	310	2063	141	
ABINIT	Fortran	23	2.3	679	5	1	96	✓
OpenMolcas	Fortran	29	1	565	29	52	2	
MPQC	C++, Fortran	12	5	6362	28	44	57	
NWChem	Fortran	29	1	26013	70	33	5	
OpenMPI	C, Fortran	147	4	28680	627	1424	99	
quantum_package	Fortran	11	4.5	2721	18	92	5	
elasticsearch	Java	1103	8.5	42349	37757	16918	223	
learnsphere	Java	9	2.5	646	11	34	1	
orca	Java	11	3.5	1103	1	143	19	
trellis	Java	3	2	892	26	171	10	
Xenon	Java	11	9	231	1	37	21	✓
abaco	Python	8	3.5	1112	13	29	7	
APBS	Python	19	5	6642	67	501	8	
forcebalance	Python	10	5	1562	48	47	6	
foyer	Python	10	3.5	343	19	66	8	
hydroshare	Python	30	4	9387	63	1708	55	
Luigi	Python	35	6	3628	10348	628	37	
mast	Python	22	5.5	5050	8	471	69	
MDAnalysis	Python	76	3.5	512	23	108	46	✓
mdtraj	Python	45	6	2971	189	682	21	
openforcefield	Python	8	1	1360	37	70	5	
openmmtools	Python	10	4	1156	40	154	31	
parsl	Python	11	2	1268	63	258	15	
pymatgen	Python	107	7	14989	322	359	228	
pyscf	Python	36	4.5	4666	217	101	41	
radical-pilot	Python	17	5	56900	26	1269	137	
RMG-Py	Python	43	1	754	11	67	17	✓
signac	Python	8	2	5000	8	89	5	
signac-flow	Python	8	2	1000	5	28	3	
TauDEM	Python	11	5.5	298	102	132	10	
Use Galaxy	Python	188	3.5	36005	507	2269	51	
yank	Python	8	5	2728	41	557	34	
yt	Python	93	1.5	23923	138	1216	37	

**Table 3: Sanity checks (designed using [17])**

Check	Condition
Personal purpose (# Developers)	$\geq 7$
Collaboration (Pull requests)	$> 0$
Issues	$> 10$
Releases	$> 1$
Commits	$> 20$
Duration	$> 1$ year

**Table 4: Labels of development types, generated via manual cross-inspection on 1000 random bug-fixing commits.**

	Absolute	Percentage
Bug Fixes	113	11%
Scientific Enhancement	370	37%
Engineering Enhancement	281	28%
Testing	113	11%
Other	132	13%

## 2.4 Terminology

The sanity checks of Table 3 uses the following terms:

- **Commit:** in version control systems, a commit adds the latest changes to [part of] the source code to the repository, making these changes part of the head revision of the repository.
- **Release:** (based on Git tags) mark a specific point in the repository's history. Number of releases defines different versions published, which signifies considerable amount of changes done between each version.
- **Duration:** length of project from its inception to current date or project archive date. It signifies how long a project has been running and in active development phase.

Later in this paper, we will also use these additional terms:

- **Open & Closed Issues:** Users and developers of a repository on Github use issues as a place to track ideas, enhancements, tasks, or bugs for work.
- **Tags:** These are references that point to specific points in the Git version control history. Tagging is generally used to capture a point in history that is used for a marked version release (i.e. v1.0.1).
- **Stars:** signifies how many people "liked" a project enough to create a bookmarks to follow the future progress of that project.
- **Developers:** These are the contributors to a project, who work on some code, and submit the code using commit to the codebase. The number of developers signifies the interest of developers in actively participating in the project and volume of the work.
- **Watchers:** These are GitHub users who have asked to be notified of activity in a repository, but have not become collaborators. This is a representative of people actively monitoring projects, because of possible interest or dependency.
- **Forks:** A fork is a copy of a repository. Forking a repository allows user to freely experiment with changes without affecting the original project. This number is an indicator of how many

people are interested in the repository and actively thinking of modification of the original version.

- **Heroes:** This is term used widely in the literature [1, 11, 27, 34] to denote those programmers within a project that do most of the work<sup>1</sup>. The usual definition of "heroes" are the 20% of developers who write more than X% of the code. The usual threshold for "X" is  $X \geq 80\%$  [27] but recently Majumder et al. [22] found that such heroes are so common in open source projects that we will use their threshold of  $X \geq 95\%$ .

## 2.5 "Endorse" and "Doubt", but not "Reject"

Using the above definitions, the rest of this paper makes the following conclusions about the beliefs of Table 1

In this paper, we "endorse" or "doubt" a belief (and we never say that a belief is "rejected"). This section explains why we that is so.

The reasoning of this paper makes modeling assumptions in order to bridge between the terminology of the belief and the terms in the Github data. For example, consider the belief "Verification and validation in software development for CS is difficult". None of our data is conveniently labelled with the tag "verification and validation". Instead, based on our reading of the commits, we could assign labels showing whether or not developers were fixing existing bugs or reporting the results of running tests. Hence, to explore that belief we had to make the following modeling assumptions to bridge between the terminology of the belief and the terms in the Github data.

- V&V is associated with testing and and bug fixing;
- The amount of testing and bug fixing is an indicator for the amount of V&V activity.

Formally, this means that our conclusions are based on what Schouten et al. describe as *indicators* [29] rather than direct measures. Indicator-based reasoning is often required in large scale data collection, especially when the goals being pursued are not directly measured. For example, in the software engineering literature, Mylopoulos et al. [2] relied on such indicators for their work on business intelligence. More generally, in the business management literature, there is much agreement that some strategy is required to bridge between the goals of the investigation and what is observable in the data. For example, in the original paper on balanced score cards (which, at the time of this writing, has over 9800 citations in Google Scholar [20]), Kaplan and Norton offer a four-layer "perspectives diagram" that implements the bridge from high-level business goals down to observable entities.

Due to the approximate nature of these indicators, it can be inappropriate to report results using specificity of a statistical hypothesis test. Hence, in this paper:

- We only offer conclusions when we can find large differences between different sets of observations.
- We eschew the term "reject", preferring instead to say "endorse" or "doubt".

<sup>1</sup> See [9] for a critical review of this particular term. While we agree with those criticisms, the term is highly prevalent in the literature. Hence, reluctantly, we use the term "hero".

### 3 BELIEFS WE CANNOT EXPLORE (USING GITHUB DATA)

Github stores data about code and the comments seen during code reviews and pull requests. While this is useful for assessing most of the beliefs of Table 1, it does mean that at least three of the thirteen beliefs, summarized by Johanson et al. [16], cannot be explored by this paper:

- *Overly Formal Software Processes Restrict Research:* Computational scientists perform many tasks, only one of which is developing software. For example, they must write grants, do presentations, traveling, keeping up with the fast-developing fields, etc. Hence, measuring (1) the formality of software process and (2) research efforts would be outside of the scope for Github.
- *Development is Driven and Limited by Hardware:* We found it hard to access information about hardware platforms from our Github data. Hence, we cannot reason about this belief.
- *Conflicting Software Quality Requirements:* As with issues relating to hardware, we found it hard to extract from our Github data information about competing requirements such as functional correctness vs performance or portability or maintainability. Note that performance issues conflict with portability and maintainability since these are often achieved via hardware-specific optimizations. As before, we found it hard to access Github information about these kinds of requirements decisions.

### 4 BELIEFS ABOUT THE NATURE OF THE SCIENTIFIC CHALLENGE

All characteristics of software development in computational science that are listed in this section result from the fact that scientific software is an integral part of a discovery process. When you develop software to explore previously unknown phenomena, it is hard to specify exactly up front what the software is required to do, how its output is supposed to look like, and how to proceed during its development.

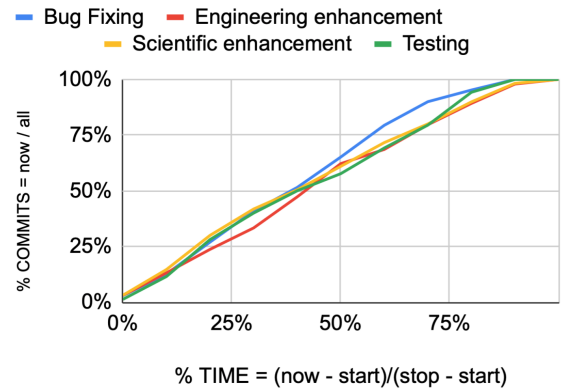
#### 4.1 Requirements

Our analysis of this first belief will conclude that CS code is built in an exploratory manner, rather than in response to some pre-defined requirements. While this first conclusion is hardly surprising, it is does offer a simple example of how this paper uses Github data to reason about CS projects.

*Belief:* Project requirements are not known up front [3, 8, 10, 30, 33].

*Rationale:* Many authors, including Carver [8] and Easterbrook [10] comment that CS code is not written in order to satisfy some pre-existing set of requirements. Rather, it is written an exploratory fashion in order to better understand some effect. This would make CS software very different to code developed using (e.g.) a waterfall model where the requirements are all known at the start of the development.

If CS software was written in response to some pre-existing set of requirements, then we would expect to see bug-fixing and testing to be a predominately end-stage activity. However, as seen here, enhancements, bug-fixing and testing all occur at very similar (and near constant)



**Figure 1: BELIEF1 RESULTS: Median percent of total commits seen at 10,20,30,...100% of the time these projects were documented in Github. Data from the 59 projects of Table 2. Note that all commit types occur at a similar, and near constant, rate, across the lifetime of a project.**

*Modeling Assumptions:* Projects with pre-existing list of fixed requirements can be developed in a “waterfall” style where requirements is followed by analysis, design, code, implementation and test. The observable feature of such projects is that most of the testing and bug fixing activity occurs *after* the a code base has been enhanced with the required scientific or engineering functionality.

*Prediction:* Under these assumptions, if we trace the kinds of commits seen across the lifetime of a project, we should see:

- Early in the project’s lifetime, far more enhancements than bug-fixing and testing commits;
- Later in the project’s lifetime, far more testing and bug fixing commits than enhancements.

*Result:* As shown in Figure 1, the rate of commits of different types is nearly constant across the project lifetime.

*Conclusion:* We endorse the belief that, in CS, project requirements are usually not pre-defined at the start of a project.

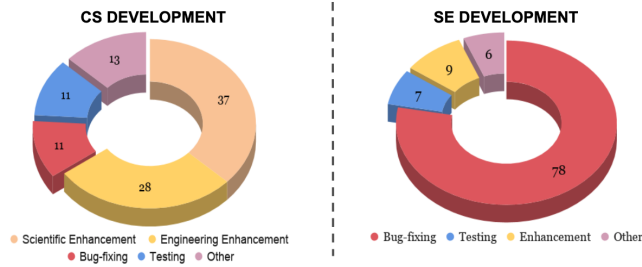
#### 4.2 Verification and validation

*Belief:* Verification and validation in software development for CS is difficult and strictly scientific [3, 7, 8, 19, 25].

*Rationale:* According to Carver et al. [8], CS verification means ensuring the mathematical model matches the real world; while CS validation means ensure the computational model matches the mathematical model. Past studies argued that verification and validation of scientific software should be difficult for several reasons:

- Lack of suitable test oracles [19],
- Complex distributed hardware environments with no comparable software [3],
- Scientists often suspect that the problems of the software is the results of their scientific theory [36],
- Lack of physical experimentation and experimental validation is impractical [8].

*Modeling Assumptions:* See above, in §2.5



**Figure 2: Distribution of development activities within CS (left) and SE (right) projects within our sample.**

**Table 5: Labels of testing type commits from the labeled Testing commits.**

	Absolute	Percentage
Science	29	45%
Engineering	11	17%
Other	25	38%

**Prediction:** Verification and validation in CS “difficult” if the observed CS effort in this area is much larger than some known baseline. As to “strictly scientific”, we should see far more “scientific enhancements” that otherwise (e.g. “bug fixes” plus “engineering enhancement”).

**Result:** It is easy to show that CS software verification and validation is heavily focused on scientific issues. Table 5 shows that “scientific testing” is the largest type of commit in our labeled Testing commits sample (at 45%). Far less effort is spend on “engineering testing” (only 17%).

As to showing the CS verification and validation is “more difficult”, the Table 4 should be compared to Figure 2, which shows the 85% percent of commits in standard SE projects associated with bug fixes and enhancements. This data comes from a recent study of the top-20 highly starred from Github that satisfy our sanity checks of Table 3<sup>2</sup>. Since 85% is much larger than 22%, we conclude that, for verification and validation, far less effort is being spent in CS projects than SE.

**Conclusion:** We endorse the belief that CS V&V is mostly concerned with scientific issues. But since CS V&V requires fewer commits that SE projects, we cannot endorse a belief that verification and validation in software development for CS is more difficult than in other disciplines.

**Discussion:** This result is somewhat strange since it runs counter to standard beliefs in the SE literature (e.g. Brookes argues that unit tests and systems tests will consume half the time of any project [5]). We conjecture that the larger V&V effort in SE is due to the nature of CS problems. CS software is more grounded in unchanging physical realities that standard SE software:

- CS software is written to correspond to physical phenomena, the nature of which may never change (e.g. the atomic weight of iron).

<sup>2</sup>TBD: need more data. 20 se projects

**Table 6: Language Usage within sample of quality CS projects.**

	Absolute Count	Percentage
Other	1	1%
Javascript	2	3%
C	4	7%
Java	5	9%
Fortran	8	14%
C++	16	27%
Python	23	39%

- On the other hand, standard SE software (e.g. the highly starred projects in Github) is written to correspond to an ever-changing ecology of platforms, tools, user expectations, and newly-arrive AI algorithms, etc etc. Hence, it is not surprising SE software requires more verification and validation effort than CS software since the problem it addresses are more dynamic.

Whatever the reason, note that this result calls for a different kind of testing device in CS. In standard SE, a “test” can be something as simple as a unit test (checking if, for example, that subtrees remain in sorted order after insertion). But in CS, “tests” need to be a higher level and refer back to some core physical properties as defined by scientific theory.

## 5 LIMITATIONS OF COMPUTER HARDWARE

In this section, we discuss characteristics of software development in computational science that are due to limitations regarding available computing resources and their efficient programming.

### 5.1 2b. Programming Languages, Technologies, and SE Methods

In this section we explore a combination of beliefs 2b and 3f. **Belief:** A widespread view is that computational scientists prefer “older”-style programming languages and technologies while disregarding most of the newer SE methods [3, 8, 14, 25, 26].

**Rationale:** CS Scientists are skeptical of modern SE methods and new technologies/languages. Given their success with older-style languages (Fortran and C). This is based on several factors:

- A decades-long commitment with these older-style languages on high-performance computing platforms [36].
- A belief that the extra features of the newer languages needlessly conflate functionality that can be more easily implemented in (e.g.) one line of “C” macros [28].
- A prejudice against the newer languages or a perception that the scientists would not find them useful [25].

#### Modeling Assumptions:

- Old languages/technologies: based on the literature review of Johanson et al. [16], we say that “C” and “Fortran” are the older, most established languages in the CS community. Everything else, we will call “newer languages”.
- Modern SE practices are associated with automating tests and deployment continuously (e.g. Travis CI)

In order to assess this combined belief (2b and 3f), the usage of language per project within the sample (Table 6) is recorded along with the usage of Travis CI.

**Prediction:** We would endorse this belief if (1) most the languages used by the CS projects are “older” style and (2) most CS projects do not adopt the usage of Travis CI.

**Result:** “C” and “Fortran” are just 29% of our sample while most of our projects use “newer” languages (where “new” is defined by Johanson et al. [16]).

We have other evidence that CS developers might be more open to newer technologies than suggested in the current literature. In looking over our 59 projects, we observe that 43 of them (73%) have active Travis CI connections<sup>3</sup>. Travis CI is a continuous deployment tool that run tests as a side-effect of any commit being written to Github. It is a standard tool within the continuous deployment community.

**Conclusion:** We do not support a belief that CS Scientists are skeptical of modern SE methods and new technologies/languages.

**Discussion:** While this result is at odds with numerous papers [3, 8, 14, 25, 26]. We explain our novel findings as follows. Most of the papers that endorse this view come from before the recent Silicon Valley boom. In our discussions with postdocs working on CS projects, we noted that they were very aware of the salaries they might earn if, after completing their postdoc, they moved on to software companies. They seemed very well aware that a condition of that career move would be a deep understanding of the kinds of tools used by contemporary agile software companies. Hence, it is perhaps not so surprising that we report here a widespread use of modern software techniques in CS.

## 5.2 2c. Domain Logic and Implementation Details

**Belief:** Intermingling culture of domain logic and implementation details within scientific software development [36].

**Rationale:** Scientific software development is different from traditional software development due to the inseparable relationship of the usage of older programming languages and software with the focus of scientific models performance. The developers of such software should be, but difficult to be, proficient in both aspects. This leads to researchers having difficulty in evolve one aspect independently.

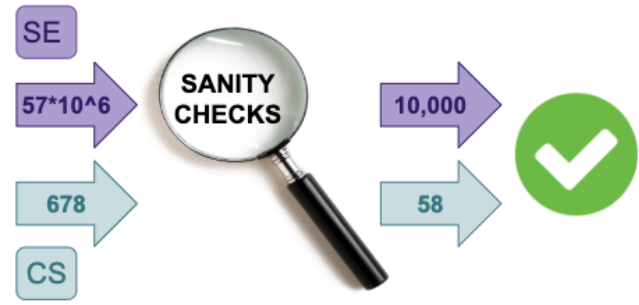
**Modeling Assumptions:** During scientific software development,

- domain logic addresses computational understanding, i.e. scientific enhancement.
- implementation details address coding/building the tool up to solve scientific problem, i.e. engineering enhancement.

**Prediction:** If domain logic and implementation details are intermingled/inseparable during the development, then both scientific and engineering enhancement contribution distribution should be the same or have a very small difference from each other.

**Result:** Across the enhancement type commits from the sample (from Table 4), 370 or 57% enhancement commits focusing on the core science while the rest of 281 or 43% enhancement commits focusing on the quality of the code. The absolute difference between two type of enhancement activities is small (14%)

<sup>3</sup>check that their travis CI accounts are active



**Figure 3: The SE (cyan) and Computational Science (magenta) projects count number of before and after sanity checks.**

**Conclusion:** the small difference indicated our endorsement for this characteristic that developers of scientific software spent their time and efforts to enhance both the core science and the code evenly in parallel.

## 6 CULTURAL ENVIRONMENT OF SCIENTIFIC SOFTWARE DEVELOPMENT

The characteristics that are listed in this section result from the cultural environment in which scientific software development takes place. This environment is shaped, for example, by the training of computational scientists and the funding schemes of scientific research projects.

### 6.1 SE Training

**Belief:** Few computational scientists are trained in SE [3, 6, 10, 28, 32].

**Rationale:** Not many of scientists have SE knowledge or skill due to (1) SE are only considered as “techniques”, a means to an end during the development of scientific software [28] and (2) learning SE is perceived as an excessive demand [21].

**Modeling Assumptions:** training in SE is indicated by

- the general quality of the software (e.g. the amount of projects that pass the sanity checks)
- the adoption of SE practices (e.g. the

**Result:** In order to assess this belief, we first checked the number of projects between SE and Computational Science community after the sanity checks (from Table 3 and mentioned in §2.1) to make sure that they are more than just “hobby” projects as Figure 3 showed. We can see that after the sanity checks, the proportion of quality projects that are in CS community is 486 times higher than the quality projects that are in SE community (the calculation justification is offered right below). In another word, CS developers are more serious about their software development.

$$\frac{CS\_post\_pre\_sanity\_rate}{SE\_post\_pre\_sanity\_rate} = \frac{\frac{CS\_post\_sanity}{CS\_pre\_sanity}}{\frac{SE\_post\_sanity}{SE\_pre\_sanity}} = \frac{\frac{58}{678}}{\frac{10,000}{57*10^6}} = 486$$

Furthermore, we summarized the Github statistics of 1,000+ projects within the “Github showcase project” as SE projects and 59 quality CS projects in Figure 4. Beside the *Developer* statistics,

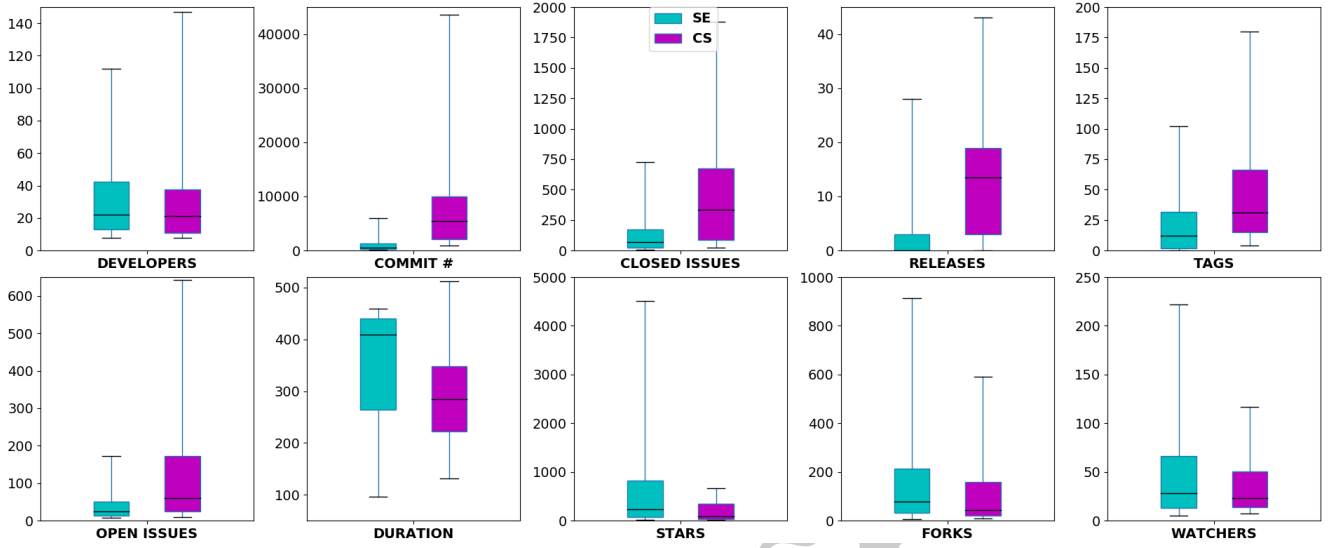


Figure 4: Github statistics comparison of 1,000+ SE projects (Cyan) & 59 CS projects (Magenta).

Table 7: Win percentages of G-score (left) and  $P_{opt20}$  (right) from [35]. Gray cells highlight the labeling method that were top-ranked most in that project by the statistical tests (in  $P(W/N)$  format). Treatments: SE method (SE) and tuned SE method for CS (CS)

Dataset	% G-score Wins		Dataset	% $P_{opt20}$ Wins	
	SE	CS		SE	CS
PCMSOLVER	100 (1/1)	0 (0/1)	PCMSOLVER	100 (1/1)	0 (0/1)
AMBER	67 (2/3)	33 (1/3)	XENON	50 (3/6)	50 (3/6)
HOOMD	40 (2/5)	60 (3/5)	MDANALYSIS	43 (3/7)	57 (4/7)
RMG-PY	40 (2/5)	60 (3/5)	LIBMESH	14 (1/7)	57 (4/7)
ABINIT	25 (2/8)	63 (5/8)	HOOMD	40 (2/5)	60 (3/5)
LIBMESH	28 (2/7)	72 (5/7)	LAMMPS	25 (2/8)	63 (5/8)
MDANALYSIS	28 (2/7)	72 (5/7)	ABINIT	25 (2/8)	63 (5/8)
LAMMPS	25 (2/8)	75 (6/8)	AMBER	33 (1/3)	67 (2/3)
XENON	17 (1/6)	83 (5/6)	RMG-PY	0 (0/5)	80 (4/5)

it is observed that they committed, closed issues, deployed/released, and tagged more while the projects are shorter in duration. Scientific developers develop with SE process philosophies: build software faster and in a more granular fashion than SE projects.

**Conclusion:** With both of these evidences, we doubt this belief that few computational scientists are trained in SE as scientific developers (1) are more serious in developing their software and (2) built their software faster in a more granular fashion.

## 6.2 Terminology

**Belief:** Both SE and CS fields used different terminologies (or even languages) to describe the things they do, sometimes, the same activity [10, 21, 36].

**Rationale:** SE and CS fields have developed in isolation with established distinct languages. As a result, the scientific programmers

might have not realized they have re-invented or used existing SE techniques.

**Modeling Assumptions:**

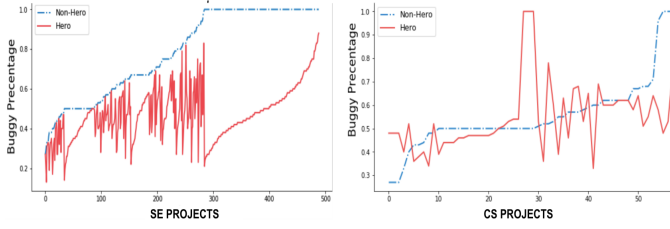
- Terminologies are regarded here as how scientists describe their work in the commits documentation.
- The difference between using off-the-shelf SE method versus tuned method that learn the CS language to label CS commits is the indicator.

**Prediction:** If the belief is hold, then the off-the-shelf SE method will perform badly when labeling scientific development commits in comparison with with tuned method that learn the CS language.

**Result:** Tu et al. [35] demonstrated that SE and CS used different languages in development documentation as different backgrounds come different terminologies and the language usage. They noted that the difference is so large that SE method for identifying bug-fixing commits cannot be adapted without prior modified to CS community. Comparing to human, after tuning the SE method to learn the CS language, they were able to reproduce better ground truths than SE method. That better ground truths generated more quality prediction for defect predictors. Specifically, Table 7 compares predictive performance using CS method (tuned SE method to learn the CS language) and standard SE method. The treatments with more superior performance are denoted by gray color. Note that, in the majority case 7 out of 9 projects for both G-score and  $P_{opt20}$ .

**Conclusion:** With this empirical study's straightforward result, we can confidently endorse that the computational science community utilize a different language when describing and documenting their work.

**Discussion:** Further, this result also raises attention for the current bad research fashion of reusing and applying established work without critiques. As seen, off-the-shell SE standard methods did not apply well for CS development data. The methods need to tuned to specialize how scientific researchers develop their software.



**Figure 5: Defect Introduction rate by hero and non-hero developers within SE projects (left) and CS projects (right).**

**Table 8: The table summarizes of Figure 5 and stratifies the data according to 25th, 50th, and 75th percentiles of buggy percentages introduced through code interaction.**

Percentile	Category					
	CS Projects			SE Projects		
	Hero	Non-Hero	Ratio	Hero	Non-Hero	Ratio
25th	52	67	1.3	46	50	1.09
50th	58	75	1.3	39	43	1.1
75th	53	100	1.9	52	60	1.15

### 6.3 Value

**Belief:** Scientific software in itself has no value but still it is long-lived [10, 21, 32, 36].

**Rationale:** It is due to the belief that software is a representation of the underlying scientific theory with no novelty value of for new scientific discoveries [28].

**Modeling Assumptions:**

- CS software’s value is associated with the popularity of the project on Github (open issues, stars, watchers, tags, and forks)
- Live of scientific software is associated with the Duration of the software on Github.
- In both case, the difference of the metric for SE versus CS is the indicator.

**Prediction:**

**Result:** According to the Figure 4, the number of *Open & Closed Issues* are higher in Computational Science projects than in SE projects. The similar distributions of *Stars, Forks*, and *Watchers* between CS and SE projects indicated similar values are placed by the audiences for both types of software. Meanwhile, the distribution of *Duration* demonstrated that even when the scientists believed their software is long-lived, the actual Duration of their software development is observed, on average, to be significantly shorter, almost 75%, of the *Duration* of the SE software. Similar to 4.3.1, from our sample, CS projects’ lives are shorter with a greater commit density than SE projects. It indicates that CS people actually understand their code and actually offer more relevant contributions to the project’s development.

**Conclusion:** From those statistics, they are some clear signals that we doubted this belief that CS softwares have no value but long-lived.

### 6.4 Code Understanding

**Belief:** Creating a shared understanding of “code” is difficult [7, 15, 28, 31].

**Rationale:** All scientists typically (1) do not produce documentation for the software they implement [28, 32] and (2) high personnel turnover rates in scientific software development [7, 31]. As a result, it renders such a knowledge and skill transfer problem.

**Modeling Assumptions:**

- Code understanding here is associated with how experts (heroes) and novices/outside (non-heroes) interacting through the code.
- The defects introduction rate gap between two types of developers and the difference between the defects introduction rate gap in CS versus SE community are the indicator for code understanding.

**Prediction:** If belief is hold, the defects introduction rate of hero developers would be a lot lower than the defect introduction rate of non-hero developers and the defects introduction rate gap of CS community would be the same or wider than the gap of SE community.

**Result:** Menzies et al. [22] checked the heroes projects for both heroes and non-heroes contribution of defects within software development. They found that non-heroes introduced 1.3-1.9 times more bugs (25th-75th percentiles) in those projects. However, when the study is replicated for CS projects, Figure 5 and Table 8 indicated that 25th-75th percentile, non-hero developers introduced slightly higher than the amount of bugs introduced by hero developers (1.09-1.15 times). While fewer people (heroes) made most of the changes within the software development, non-experts and novices are able to contribute to the software development without causing a lot of bugs (in respect to developer dynamics in SE projects development).

**Conclusion:** At least, in the aspect of defect, we doubted that shared understanding of “code” is difficult within the CS community.

**Discussion:** The even distribution of similar defects indicates that CS developers understand the project and make more relevant contributions to the software development. It is a call for the CS community to organize and encourage more novices and outsiders to contribute to the scientific software development.

## 7 THREADS OF VALIDITY

### 7.1 Construct Validity

Repeat above discussion on indicators.

### 7.2 Sampling Bias

Like any data mining paper, our work is threatened by sampling bias; i.e. what holds for the data we studied here may not hold for other kinds of data. Within the space of one paper, it is hard to avoid sampling bias. However, what researchers can do is make all their scripts and data available such that other researchers can test their conclusions whenever new data becomes available. To that end, we have made all our scripts and data available at [github.com/blinded-for-reviews/](https://github.com/blinded-for-reviews/).

### 7.3 External Validity

Methodology to understand the difference between software development domains in SE. Not our final conclusions or definitions, all are

threatened by external validities but our conclusions are reproducible and our analysis can be repeated when new data arrives.

## 8 CONCLUSION

In this paper we seek quantitative evidence (from dozens of CS projects housed in Github) for 13 previously published conjectures about scientific software development in the literature [16]. In all, we explore three groups of beliefs about (1) the nature of scientific challenges; (2) the implications of limitations of computer hardware; and (3) the cultural environment of scientific software development. Of those, we find that four cannot be assessed with respect to the Github data. Of the remaining, we can only find support for three of the nine beliefs.

We found that scientific developers are more capable of SE practices and knowledges than they believed. Because of the understanding/miscommunication gap, they have operated with SE philosophies unknownly. As the nature of the CS software development is changing, scientific software developers may not need to incorporate SE knowledge and practices but they need tuned and specialized SE methods for CS community. Software engineers can therefore help scientific software developers to tailor the existing SE practices and knowledge to better fit the needs of the scientific software developers.

It leads to a new characterization of the nature of CS software development which, in turn, suggests a new prioritization for tool development in this area. Therefore, more research on this topic is needed, especially to empirically evaluate the actual gains in productivity and quality achieved for scientific software by such SE approaches.

## 9 ACKNOWLEDGMENTS

This work was partially funded by an NSF CISE Grant #1826574 and #1931425.

## REFERENCES

- [1] A. Agrawal, A. Rahman, R. Krishna, A. Sobran, and T. Menzies. 2018. We don't need another hero?: the impact of heroes on software development. In *ICSE*.
- [2] Daniele Barone, Lei Jiang, Daniel Amyot, and John Mylopoulos. 2011. Composite indicators for business intelligence. In *International Conference on Conceptual Modeling*. Springer, 448–458.
- [3] V. R. Basili, J. C. Carver, D. Cruzes, L. M. Hochstein, J. K. Hollingsworth, F. Shull, and M. V. Zelkowitz. 2008. Understanding the High-Performance-Computing Community: A Software Engineer's Perspective. *IEEE Software* 25, 4 (July 2008), 29–36. <https://doi.org/10.1109/MS.2008.103>
- [4] C. Bird, P. C. Rigby, E. T. Barr, D. J. Hamilton, D. M. German, and P. Devanbu. 2009. The promises and perils of mining git. In *Mining Software Repositories*.
- [5] Frederick P Brooks Jr. 1995. The mythical man-month (anniversary ed.). (1995).
- [6] J. Carver, D. Heaton, L. Hochstein, and R. Bartlett. 2013. Self-Perceptions about Software Engineering: A Survey of Scientists and Engineers. *Computing in Science Engineering* 15, 1 (Jan 2013), 7–11. <https://doi.org/10.1109/MCSE.2013.12>
- [7] Jeff Carver, Lorin Hochstein, Richard Kendall, Taiga Nakamura, Marvin Zelkowitz, Victor R Basili, and Douglass Post. 2006. Observations about software development for high end computing. *CT Watch Quarterly* 2 (01 2006).
- [8] J. C. Carver, R. P. Kendall, S. E. Squires, and D. E. Post. 2007. Software Development Environments for Scientific and Engineering Software: A Series of Case Studies. In *29th International Conference on Software Engineering (ICSE'07)*. 550–559. <https://doi.org/10.1109/ICSE.2007.77>
- [9] Thomas Claburn. 2020. Rockstar dev debate reopens: Hero programmers do exist, do all the work, do chat a lot – and do need love and attention from project leaders. [https://www.theregister.co.uk/2020/01/24/developer\\_heroes\\_exist/](https://www.theregister.co.uk/2020/01/24/developer_heroes_exist/)
- [10] S. M. Easterbrook and T. C. Johns. 2009. Engineering the Software for Understanding Climate Change. *Computing in Science Engineering* 11, 6 (Nov 2009), 65–74. <https://doi.org/10.1109/MCSE.2009.193>

- [11] Mathieu Goeminne and Tom Mens. 2011. Evidence for the pareto principle in open source software activity. In *the Joint Proceedings of the 1st International workshop on Model Driven Software Maintenance and 5th International Workshop on Software Quality and Maintainability*. 74–82.
- [12] J. E. Hannay, C. MacLeod, J. Singer, H. P. Langtangen, D. Pfahl, and G. Wilson. 2009. How do scientists develop and use scientific software?. In *2009 ICSE Workshop on Software Engineering for Computational Science and Engineering*. 1–8. <https://doi.org/10.1109/SECSE.2009.5069155>
- [13] Dustin Heaton and Jeffrey C. Carver. 2015. Claims about the use of software engineering practices in science: A systematic literature review. *Information and Software Technology* 67 (2015), 207 – 219. <https://doi.org/10.1016/j.infsof.2015.07.011>
- [14] Michael Heroux, Roscoe Bartlett, Victoria Howle, Robert Hoekstra, Jonathan Hu, Tamara Kolda, Richard Lehoucq, Katharine Long, Roger Pawlowski, Eric Phipps, Andrew Salinger, Heidi Thornquist, R. Tuminaro, James Willenbring, Alan Williams, and Kendall Stanley. 2005. An overview of the Trilinos Project. *ACM Trans. Math. Softw.* 31 (09 2005), 397–423. <https://doi.org/10.1145/1089014.1089021>
- [15] L. Hochstein, J. Carver, F. Shull, S. Asgari, V. Basili, J. K. Hollingsworth, and M. V. Zelkowitz. 2005. Parallel Programmer Productivity: A Case Study of Novice Parallel Programmers. In *SC '05: Proceedings of the 2005 ACM/IEEE Conference on Supercomputing*. 35–35. <https://doi.org/10.1109/SC.2005.53>
- [16] A. Johanson and W. Hasselbring. 2018. Software Engineering for Computational Science: Past, Present, Future. *Computing in Science Engineering* 20, 2 (Mar 2018), 90–109. <https://doi.org/10.1109/MCSE.2018.021651343>
- [17] E. Kalliamvakou, G. Gousios, K. Blincoe, L. Singer, D. German, and D. Damian. 2014. The Promises and Perils of Mining GitHub. In *MSR*.
- [18] E. Kalliamvakou, G. Gousios, K. Blincoe, L. Singer, D. M German, and D. Damian. 2015. The Promises and Perils of Mining GitHub (Extended Version). *EMSE* (2015).
- [19] U. Kanewala and J. M. Bieman. 2013. Using machine learning techniques to detect metamorphic relations for programs without test oracles. In *2013 IEEE 24th International Symposium on Software Reliability Engineering (ISSRE)*. 1–10. <https://doi.org/10.1109/ISSRE.2013.6698899>
- [20] Robert S Kaplan and David P Norton. [n. d.]. Using the Balanced Scorecard as a Strategic Management System. *Harvard Business Review* (January-February [n. d.]).
- [21] S. Killcoyne and J. Boyle. 2009. Managing Chaos: Lessons Learned Developing Software in the Life Sciences. *Computing in Science Engineering* 11, 6 (Nov 2009), 20–29. <https://doi.org/10.1109/MCSE.2009.198>
- [22] Suvodeep Majumder, Joymallya Chakraborty, Amritanshu Agrawal, and Tim Menzies. 2019. Why Software Projects need Heroes (Lessons Learned from 1100+ Projects). *CoRR* abs/1904.09954 (2019). arXiv:1904.09954 <http://arxiv.org/abs/1904.09954>
- [23] Zeeya Merali. 2010. Computational science: Error, why scientific programming does not compute. *Nature* 467, 7317 (2010). <https://doi.org/10.1038/467775a>
- [24] N. Munaiah, S. Kroh, C. Cabrey, and M. Nagappan. 2017. Curating GitHub for Engineered Software Projects. *EMSE* (2017).
- [25] Prakash Prabhu, Thomas B. Jablin, Arun Raman, Yun Zhang, Jialu Huang, Hanjun Kim, Nick P. Johnson, Feng Liu, Soumyadeep Ghosh, Stephen Beard, Taewook Oh, Matthew Zoufaly, David Walker, and David I. August. 2011. A Survey of the Practice of Computational Science. In *State of the Practice Reports (SC '11)*. ACM, New York, NY, USA, Article 19, 12 pages. <https://doi.org/10.1145/2063348.2063374>
- [26] M. Ragan-Kelley, F. Perez, B. Granger, T. Kluyver, P. Ivanov, J. Frederic, and M. Bussonnier. 2014. The Jupyter/IPython architecture: a unified view of computational research, from interactive exploration to communication and publication.. In *AGU Fall Meeting Abstracts*, Vol. 2014. Article H44D-07, H44D-07 pages.
- [27] Gregorio Robles, Jesus M Gonzalez-Barahona, and Israel Herraiz. 2009. Evolution of the core team of developers in libre software projects. In *Mining Software Repositories, 2009. MSR'09. 6th IEEE International Working Conference on*. IEEE, 167–170.
- [28] Rebecca Sanders and Diane Kelly. 2008. Dealing with Risk in Scientific Software Development. *Software, IEEE* 25 (08 2008), 21 – 28. <https://doi.org/10.1109/MS.2008.84>
- [29] Barry Schouten, Natalie Shlomo, and Chris Skinner. 2010. Indicators for monitoring and improving representativeness of response. (2010).
- [30] Judith Segal. 2005. When Software Engineers Met Research Scientists: A Case Study. *Empirical Software Engineering* 10, 4 (01 Oct 2005), 517–536. <https://doi.org/10.1007/s10664-005-3865-y>
- [31] Judith Segal. 2007. Some Problems of Professional End User Developers. In *Proceedings of the IEEE Symposium on Visual Languages and Human-Centric Computing (VLHCC '07)*. IEEE Computer Society, Washington, DC, USA, 111–118. <https://doi.org/10.1109/VLHCC.2007.50>
- [32] J. Segal. 2007. Some Problems of Professional End User Developers. In *IEEE Symposium on Visual Languages and Human-Centric Computing (VLHCC 2007)*. 111–118. <https://doi.org/10.1109/VLHCC.2007.17>
- [33] J. Segal and C. Morris. 2008. Developing Scientific Software. *IEEE Software* 25, 04 (jul 2008), 18–20. <https://doi.org/10.1109/MS.2008.85>

- [34] MR Martinez Torres, SL Toral, M Perales, and F Barrero. 2011. Analysis of the core team role in open source communities. In *Complex, Intelligent and Software Intensive Systems (CISIS), 2011 International Conference on*. IEEE, 109–114.
- [35] Huy Tu, Zhe Yu, and Tim Menzies. 2019. Better Data Labelling with EMBLEM (and how that Impacts Defect Prediction). arXiv:cs.SE/1905.01719

- [36] M. L. Vanter, S. Faulk, S. Squires, E. Loh, and L. G. Votta. 2009. Scientific Computing's Productivity Gridlock: How Software Engineering Can Help. *Computing in Science Engineering* 11 (11 2009), 30–39. <https://doi.org/10.1109/MCSE.2009.205>