Improving software quality in bioinformatics groups through temawork

This manuscript (<u>permalink</u>) was automatically generated from <u>ferenckata/SQSeminarPaper@0092807</u> on September 21, 2023.

Authors

- Katalin Ferenc

 ✓

Centre for Molecular Medicine Norway (NCMM), Nordic EMBL Partnership, University of Oslo, 0318 Oslo, Norway · Funded by Grant XXXXXXXX

- Ieva Rauluseviciute
 - **□** 0000-0001-9253-8825 · **□** ievarau

Centre for Molecular Medicine Norway (NCMM), Nordic EMBL Partnership, University of Oslo, 0318 Oslo, Norway

- Ladislav Hovan
 - © 0000-0001-8847-9295 · ♥ ladislav-hovan

Centre for Molecular Medicine Norway (NCMM), Nordic EMBL Partnership, University of Oslo, 0318 Oslo, Norway

- Vipin Kumar
 - · princeps091-binf

Centre for Molecular Medicine Norway (NCMM), Nordic EMBL Partnership, University of Oslo, 0318 Oslo, Norway

- Anthony Mathelier [™]
 - © 0000-0001-5127-5459

Centre for Molecular Medicine Norway (NCMM), Nordic EMBL Partnership, University of Oslo, 0318 Oslo, Norway; Department of Medical Genetics, Institute of Clinical Medicine, University of Oslo and Oslo University Hospital, Oslo, Norway

Abstract

Introduction

Bioinformatics and computational biology are continuously gaining importance in biological research. These fields are heavily relying on inventions of computer science and software engineering. As Goble pointed out in 2014, about 90% of researchers used or relied on results produced by scientific software [1]. This implies that incorrect software results in invalid scientific findings.

However, bioinformaticians lack formal education in computer science and related subjects [2]. Such lack of theoretical and practical foundations hinders the adoption of good practices (e.g. continuous integration, unit tests, code reviews, planning of software architecture) established in other software-heavy endeavours. Many of these practices rely on redundancy of knowledge within team members, supported by practices such as pair programming, daily stand-up meetings and code reviews.

One main feature of academic research groups is the research projects conducted by a single PhD student or post-doc. In fact, academia prides itself on enabling individual achievements, and generally measures the success of individual effort while disregarding team effort. This often results in very few people maintaining the software product, which might be used by thousands of researchers worldwide [3,4]. Lack of funding contributes to the poor maintenance status of much of scientific software [1], which has been recognized and addressed by the Chan Zuckerberg Initiative Essential Open Source Software for Science fund [5].

The concept of a team is therefore different in a science project from a software project. While group members help each other with scientific suggestions, most often there is a single responsible person for the design and implementation of the code base. As official guidelines are rarely available on coding practices, the actual craft of software engineering is treated as secondary and up to individual judgment. We can observe (can we?) that these guidelines naturally emerge in larger groups (Seurat?), this requires a form of team organization not intrinsic to academic groups. In line with this, Hannay et al. [2] found that researchers tended to rank software engineering concepts higher if they worked in a team. Russel et al. [4] showed that high-profile code bases feature larger development teams than low-profile ones, and their activities are also associated with longer commit messages as a sign of better communication and documentation of changes. We asked ourselves whether a form of team structure organized around individual software products could improve the quality of our work.

In this work we review the findings of software engineer researchers on bioinformatics practitioners, and their suggestions for improving the overall quality of scientific software in bioinformatics and related fields (citations?). Inspired by the findings, especially related to process quality, our research group incorporated weekly discussions on code quality. We noticed that good practice requires investment in time and effort that may not pay off on an individual basis. Sharing the onboarding effort of new tools, and normalizing discussions on implementation details resulted in an overall better perceived process quality and code quality of the members.

We implemented a form of code review that fits our specific needs and context. We used our experience from these sessions, with special focus on team-based software development practices ensuring good code quality during in the update of the curation analysis software of the 2024 release of the JASPAR database (citation). In this project, as well as in some of our own individual research projects, we introduced user stories when documenting the assumptions, current features and new ideas, and we relied on development tools such as Jira and git. We note that the usage of these tools is not necessarily aligned with industry practices, due to the experimental nature of scientific software. Nevertheless, as bioinformatics becomes a more and more software-heavy field, we believe a good direction is to collectively lower the barrier to adapting to new technologies.

We aim to provide guidelines on how to get started with improving the process quality of bioinformatics groups, even without members who have formal training in computer science or software engineering.

Status of scientific and bioinformatics software from the SE perspective

A landmark paper in 2007 by Diane F. Kelly [6] discussed the separation ("chasm") between software engineering and scientific-computing community. She points out the need to bridge the one-size-fits all software engineering solutions and the particularities of the scientific software development which relies heavily on domain knowledge. Without such bridge, scientific computations keep on being performed using error-prone development practices and reaching suboptimal solutions and poor software quality. Her predictions seem to hold true even after almost 20 years.

In the past two decades, software engineering researchers have been surveying bioinformatics software, and more broadly scientific software from the software engineering perspective[[2]; more citations]. There are multiple guidelines and suggestions to improve the quality of scientific code, many of which would target students of scientific disciplines [1]. Recently, an extensive literature review has been published which collects known issues and suggested solutions [7]. Yet, these guidelines seem to not have reached the majority of the bioinformatics community, which is found to be "still in its infancy compared with the majority of other scientific disciplines" [7]. Indeed, the guidelines suggested include following agile practices, the DRY (don't repeat yourself) principle, requirements gathering and unit testing, none of which concept is intuitive or well known in the bioinformatics community. Without a shift in coding culture within bioinformatics, these concepts might remain in the terrain of unknown unknowns.

Undoubtedly, scientific software development has its own challenges. However, it cannot be an excuse for skipping good practices: as Carole Globe [1] puts it "in Hannay's survey [2], only 47 percent of scientists had a good understanding of testing, and just 34 percent thought any formal training was important. This is strange because presumably they wouldn't use and trust the results of a microscope or telescope that hadn't been built by qualified engineers or tested. Yet software is the most prevalent of all the instruments used in modern science". Software engineering emerged and has been developing to address issues naturally arising from poorly planned development, such as project failures, delays, incorrect functionality or defects [8], none of which is unknown to the scientific community. In fact, the crisis of scientific software is fairly well known and suggestions are being made from both inside and outside the community [9,10].

One key challenge is the limited funding and the lack of recognition for development and maintenance of scientific software. Currently, as Alexander Szalay puts it "The funding stops when they (researchers) actually develop the software prototype" [11]. This would be a working system, if future researchers would not want to build on each other's findings, or even tools that are meant to be reused instead of reinventing the wheel or update outdated code.

Another obstacle is the non-trivial nature of testing of scientific software. In a recent review paper [12] two key aspects of scientific software testing has been highlighted: the oracle problem and the cultural differences between scientists and software engineers. Software behaviour can be tested against an expected output, but often in science we use software to find new knowledge. This results in an oracle problem, when scientists actually do not know *a priori* how the software should behave, thus straight forward verification is impossible. According to the authors, scientists also view their scientific model and the implementation as a single entity. Therefore, scientists tend to test the validity of the model but not verify the code which produces it. Uncovered faults can and do lead to

incorrect scientific insights as shown in multiple examples [here we can have fun finding such cases or cite from the paper of this paragraph].

Nevertheless, software quality standards (Figure 1) defined by the International Organization for Standardization [13] are universally applicable. Depending on the application of the scientific software, whether it is a tool or a data analysis pipeline, the authors may prioritize different quality attributes. For example, in the world of big data, performance and efficiency gain importance. Shown in a previous study reviewing mappers, individual tools have varying level of compatibility, usability, and portability [3]; quality attributes which directly impact user experience. Frameworks, such as Snakemake [14] or Nextflow [15] support usability, reliability, and maintainability. Anaconda [16] and container solutions [17,18] help achieve portability. These are also compatible with Snakemake and Nextflow, making these frameworks staple for reproducible data analysis.

Figure 1: ISO25000

Review of existing suggestions for improving software quality in biomedical sciences

We reviewed the existing literature that focuses on providing guidelines for programming practices for scientists without extensive training in computational sciences. We used several key phrases to search for papers: "guidelines for bioinformatics software", "rules for biologists learning bioinformatics". Such papers are quite abundant and have a number of things in common. For example, they can focus on specific suggestions, often referred to as rules or "tips & tricks", or they can more broadly direct towards good practices of coding, which are put together into "guidelines".

Papers often choose a main focus or a theme and then distinguish a set of rules or guidelines. That focus can be very specialized, such as particular data analysis in a single disease [19]. Or the theme can be more general, for example setting the rules for next-generation sequencing (NGS) data analysis or outlining tips on how to start on computational analysis of the experimental data [20,21,22]. Both types of papers emphasize the need to learn how to analyse data properly and provide good suggestions to do that, based on the chosen topic. However, one needs to be aware when reading specialized focus guidelines that some of them can be much less applicable if the context of the analysis is different. Finally, the reference point matters. There is a difference, if the guidelines are read by a person with less experience in computational data analysis, where even a small set of rules can mean a steep learning curve, because some skills that are dependencies for the rules might be common knowledge for those with more experience. For example, a guideline to use version control systems, such as git, is a very important one, but for that a reader should most likely know the basics of the unix command line, which is not necessarily a rule on itself. One of the first things that should be done while suggesting your set of guidelines is clearly defining your target audience, even if they are meant to be broader suggestions.

As most guidelines tend to be written for researchers or students with minimal coding experience, suggestions often overlap. Most commonly highlighted are documentation and version control [23,24]. These are basic aspects of the software quality, which are still often overlooked in the publications making it hard to reproduce the research or use a published software [7]. Even the people with the least experience in computational analyses should always document their code by writing extensive README documentation. Documentation practice is very strong in the wet lab, where having a lab journal is unquestionable. This should not be different when working with computational tools. Version control of the code improves research reproducibility and the usage of the software. The most popular way to do that is by using git and remote repositories, for example, GitHub or Bitbucket [19,24,25]. Less commonly emphasized are testing of the developed code or

software and optimization of pipelines [24]. These are very important rules to follow to end up with a proper collection of code or a stable software [7]. However, it might be considered to be too advanced for "beginners" guidelines. Unfortunately, not being aware of such requirements can lead to a code that either takes a long time to run or is buggy, which accumulates the technical debt and, in a way, encourages the bad practices. To avoid that it is important to at least be aware of the caveats of your code and document them for fellow researchers that might use your code in the future.

As code testing and organization can be overwhelming for a less experienced computational scientist, there are multiple available tools to help organize the coding tasks and the code itself. For example, using task management through Jira or GitHub issues, which are commonly used in team projects, where multiple people work on the same code. However, these resources are not emphasized in the guideline's literature as this literature is most commonly focused on the individual persons and their personal practices. Often "other people" are only mentioned when guidelines suggest where to seek help when encountering a problem with the code. This includes consulting with colleagues, finding a mentor or participating in online communities (for example, Stack Overflow or Biostars) [24]. More recently, new tools are introduced that integrate artificial intelligence (AI) as a helper in coding. Tools such as GitHub copilot or editor extensions with access to ChatGPT allow to ask the models to write a piece of code to address a problem. To our knowledge bioinformatics literature almost never presents suggestions how to code in a team setting and utilize multiple people's expertise on software development. In contrary to software engineering-oriented literature, where there is a lot of focus on coding in a team practices [26,27]. Hagan et al. described Code Clubs - the practice in their research lab, where group members are collectively engaged in software development through code reviews and pair coding and software engineering education through workshops or seminars [28]. The authors give tips on how to organize such meetings and what should be the ground rules. Sharing your coding experience with others helps minimize the isolation, allows individuals to learn from their peers, and finally helps to write a better quality software.

Our experiences for development processes involving teams {.page_break_before}

The preceding sections mention a lot of possible approaches to improving software quality. Given the abundance of opinions on this topic, and the variety of challenges bioinformaticians face, we believe that everyone should find out what works best for them. Here, we describe the practices that we have settled on.

The software development practices that we have adopted can be broadly separated into three categories: code reviews, what we have called software quality meetings, and resource sharing.

Code reviews

Code reviews are not a new invention and many people have discussed their benefits [28,29]. Here we would just like to briefly summarize how being made to present your code and receiving feedback leads to improvement in the process of creating software.

Prior to a scheduled code review, the author is forced to write their code in a way that it will be explainable and understood by others, which is always desirable. In a large distributed project this may be trivial, but because the bioinformatic projects are often handled by a single person, it is very possible to make the code needlessly complex and obfuscated. We also observed that during data analysis parts of the code are re-run in an ad-hoc manner (e.g. by commenting out parts), making it increasingly difficult to reproduce the same analysis.

During the code review, the author has to explain some aspect of their code clearly (e.g. structure, algorithm implementation, performance related decisions), which depends on them understanding it. Trying to explain your code to someone is shown to help with understanding (rubber duck method [30]). The feedback obtained can help fix existing or potential future issues, improve the implementation, and produce cleaner, more concise code. The other participants may not be deeply familiar with the particular project, but they have their unique knowledge and point of view. We agree with the ten simple rules described by Hagan et al. [28], and note that many of those naturally emerged as a code of conduct after a few rounds of trial and error.

After the review, the received suggestions should be implemented swiftly to improve the code before advancing the project. The success of code review is highly dependent on its frequency (long time between reviews - a lot of new code, hard to cover all changes in a single session, potentially a lot of rewrite post review), and hence they should be as regular and frequent as reasonably possible.

As a positive additional outcome, we noticed an increasing understanding in each other's projects that naturally emerged through talking about the analysis code. This enabled us to give more involved comments during subsequent group meetings too. We noted however, that the focus can easily shift from the code to the biological question at hand. This we believe is more of a feature than a bug, as each code review session is led by the person bringing the code and the rest of us are there to support to the best of our abilities. Especially after the general level of coding style and quality increased to a good baseline. E.g. after about half a year, it was trivial for everyone involved that code organized into functions is preferred over spaghetti. The shared knowledge base and standards also allow us to make new group members adopt good coding practices more quickly.

Software quality meetings

Within the framework of software quality meetings, we have established larger-scale knowledge transfer between the participants. Presentations and demonstrations of new techniques and tools that are not necessarily tied to a specific project help broaden our knowledge base and awareness. In this sense, they form almost a substitute for a more formal computer science education, which most bioinformaticians lack [7]. Topics can arise from code reviews, own projects, or effectively be a reproduction of a useful talk or seminar given elsewhere.

The presenters benefit as well by having to research the topic further and present it coherently. It is not necessary to have these meetings be as regular as code reviews. The time investment is higher, given that a preparation is needed unlike just writing code as for code reviews.

During the software quality meetings, we have also explored the possibility of collaborative projects and pair programming, but have not managed to implement it successfully yet outside the scope of preparation for the JASPAR 2024 release (reference). The main reason for this is that we experimented with collaboration on a software tool not directly used by any of the members. As researchers, we could not afford to invest time in a hobby project.

The outcome of these sessions are manifold. A few examples: 1) a shared vocabulary that enables quick discussion about implementation details and code structures (e.g. design patterns, software architecture, data structures and algorithms), 2) a kind of toolkit and set of recordings we can sample from and build on in our own research projects (e.g. planning with UML diagrams, git features to ease and quicken software development), 3) awareness of previously unknown algorithms and packages, improving software performance and quality (e.g. dynamic programming, heap, Python packages such as bioframe).

Resource sharing

Resource sharing is a basic thing, but it boils down to making sure that useful online resources are brought to the attention of all participants easily.

Resource sharing could be discussed from two perspectives: external open-access resources (forums, repositories, packages and libraries) and internal (within-group resources with tools). The latter is very important as it allows for team contribution that can benefit the individual project development. A simple example of this could be a shared repository of various computational tools that were developed by members of the group. Such tools are universal enough and fit the group's research questions, so all people in the group can re-use them. In addition, each tool can be potentially developed and reviewed by multiple group members.

During software meetings, we aimed to set aside time to improve these tools from perspectives identified by the members. We observed that many of these tools do not have a clear scope and are rather a small script for a sub-task from a previous project. Based on this observation, we noted that there is a difference between a script and a standalone tool that can be inserted into various projects. The latter requires exploration of use cases related to the tool, handling of unexpected inputs, and extensive documentation, to name a few tasks. This understanding was actually quite relevant in a code review discussion when the expected usage modes of a new tool was the main focus.

(Explanation of figure) (Figure 2)

Figure 2: wall_climbing

Figure 2: wall_climbing

References

1. Better Software, Better Research

Carole Goble

IEEE Internet Computing (2014-09) https://doi.org/vjz

DOI: 10.1109/mic.2014.88

2. How do scientists develop and use scientific software?

Jo Erskine Hannay, Carolyn MacLeod, Janice Singer, Hans Petter Langtangen, Dietmar Pfahl, Greg Wilson

2009 ICSE Workshop on Software Engineering for Computational Science and Engineering (2009-05) https://doi.org/bw966x

DOI: 10.1109/secse.2009.5069155

3. Empirical study on software and process quality in bioinformatics tools

Katalin Ferenc, Konrad Otto, Francisco Gomes de Oliveira Neto, Marcela Dávila López, Jennifer Horkoff, Alexander Schliep

Cold Spring Harbor Laboratory (2022-03-13) https://doi.org/grx4jr

DOI: 10.1101/2022.03.10.483804

4. A large-scale analysis of bioinformatics code on GitHub

Pamela H Russell, Rachel L Johnson, Shreyas Ananthan, Benjamin Harnke, Nichole E Carlson *PLOS ONE* (2018-10-31) https://doi.org/gskr8b

DOI: <u>10.1371/journal.pone.0205898</u> · PMID: <u>30379882</u> · PMCID: <u>PMC6209220</u>

5. **CZI - Essential Open Source Software for Science**

Chan Zuckerberg Initiative

https://chanzuckerberg.com/eoss/

6. A Software Chasm: Software Engineering and Scientific Computing

Diane F Kelly

IEEE Software (2007-11) https://doi.org/cbrmv5

DOI: 10.1109/ms.2007.155

7. Improving bioinformatics software quality through incorporation of software engineering practices

Adeeb Noor

PeerJ Computer Science (2022-01-05) https://doi.org/gsm3hg

DOI: 10.7717/peerj-cs.839 · PMID: 35111923 · PMCID: PMC8771759

8. **Bridging the Chasm**

Tim Storer

ACM Computing Surveys (2017-08-25) https://doi.org/gftvrn

DOI: 10.1145/3084225

9. Hunting for the best bioscience software tool? Check this database

Matthew Hutson

Nature (2023-01-12) https://doi.org/gsnnww

DOI: 10.1038/d41586-023-00053-w

10. Why science needs more research software engineers

Chris Woolston

Nature (2022-05-31) https://doi.org/gsnnwt

DOI: <u>10.1038/d41586-022-01516-2</u>

11. Ex-Google chief's venture aims to save neglected science software

David Matthews

Nature (2022-07-13) https://doi.org/gsnnwv

DOI: 10.1038/d41586-022-01901-x

12. Testing Scientific Software: A Systematic Literature Review

Upulee Kanewala, James M Bieman arXiv (2018) https://doi.org/gsrxg5
DOI: 10.48550/arxiv.1804.01954

13. **ISO 25010** https://iso25000.com/index.php/en/iso-25000-standards/iso-25010?limit=3%20

14. Sustainable data analysis with Snakemake

Felix Mölder, Kim Philipp Jablonski, Brice Letcher, Michael B Hall, Christopher H Tomkins-Tinch, Vanessa Sochat, Jan Forster, Soohyun Lee, Sven O Twardziok, Alexander Kanitz, ... Johannes Köster

F1000Research (2021-01-18) https://doi.org/gjjkwv

DOI: 10.12688/f1000research.29032.1 · PMID: 34035898 · PMCID: PMC8114187

15. Nextflow enables reproducible computational workflows

Paolo Di Tommaso, Maria Chatzou, Evan W Floden, Pablo Prieto Barja, Emilio Palumbo, Cedric Notredame

Nature Biotechnology (2017-04) https://doi.org/gfj52z

DOI: 10.1038/nbt.3820

16. Anaconda | The World's Most Popular Data Science Platform

Anaconda

https://www.anaconda.com/

17. **Home**

Docker Documentation (2023-08-22) https://docs.docker.com/

18. **Home** https://apptainer.org/

19. Guidelines for bioinformatics of single-cell sequencing data analysis in Alzheimer's disease: review, recommendation, implementation and application

Minghui Wang, Won-min Song, Chen Ming, Qian Wang, Xianxiao Zhou, Peng Xu, Azra Krek, Yonejung Yoon, Lap Ho, Miranda E Orr, ... Bin Zhang

Molecular Neurodegeneration (2022-03-02) https://doi.org/gptggt

DOI: <u>10.1186/s13024-022-00517-z</u> · PMID: <u>35236372</u> · PMCID: <u>PMC8889402</u>

20. A Clinician's Guide to Bioinformatics for Next-Generation Sequencing

Nicholas Bradley Larson, Ann L Oberg, Alex A Adjei, Liguo Wang *Journal of Thoracic Oncology* (2023-02) https://doi.org/gsm3mb

DOI: <u>10.1016/j.jtho.2022.11.006</u> · PMID: <u>36379355</u> · PMCID: <u>PMC9870988</u>

21. Practice guidelines for development and validation of software, with particular focus on bioinformatics pipelines for processing NGS data in clinical diagnostic laboratories

Nicola Whiffin, Kim Brugger, Joo Wook Ahn

PeerJ (2017-05-29) https://doi.org/gsm3mc

DOI: 10.7287/peerj.preprints.2996

22. Standards and Guidelines for Validating Next-Generation Sequencing Bioinformatics Pipelines

Somak Roy, Christopher Coldren, Arivarasan Karunamurthy, Nefize S Kip, Eric W Klee, Stephen E Lincoln, Annette Leon, Mrudula Pullambhatla, Robyn L Temple-Smolkin, Karl V Voelkerding, ... Alexis B Carter

The Journal of Molecular Diagnostics (2018-01) https://doi.org/gcsstd

DOI: 10.1016/j.jmoldx.2017.11.003

23. Top considerations for creating bioinformatics software documentation

Mehran Karimzadeh, Michael M Hoffman

Briefings in Bioinformatics (2017-01-14) https://doi.org/bzmp

DOI: 10.1093/bib/bbw134 · PMID: 28088754 · PMCID: PMC6054259

24. Ten simple rules for getting started with command-line bioinformatics

Parice A Brandies, Carolyn J Hogg

PLOS Computational Biology (2021-02-18) https://doi.org/gh32h2

DOI: 10.1371/journal.pcbi.1008645 · PMID: 33600404 · PMCID: PMC7891784

25. Practice guidelines for development and validation of software, with particular focus on bioinformatics pipelines for processing NGS data in clinical diagnostic laboratories

Nicola Whiffin, Kim Brugger, Joo Wook Ahn

PeerJ (2017-05-29) https://doi.org/gsm3md

DOI: 10.7287/peerj.preprints.2996v1

26. https://faculty.washington.edu/ajko/books/cooperative-software-development/

27. Studying the impact of social interactions on software quality

Nicolas Bettenburg, Ahmed E Hassan

Empirical Software Engineering (2012-04-28) https://doi.org/f4mhdp

DOI: 10.1007/s10664-012-9205-0

28. Ten simple rules to increase computational skills among biologists with Code Clubs

Ada K Hagan, Nicholas A Lesniak, Marcy J Balunas, Lucas Bishop, William L Close, Matthew D Doherty, Amanda G Elmore, Kaitlin J Flynn, Geoffrey D Hannigan, Charlie C Koumpouras, ... Patrick D Schloss

PLOS Computational Biology (2020-08-27) https://doi.org/gg92xw

DOI: 10.1371/journal.pcbi.1008119 · PMID: 32853198 · PMCID: PMC7451508

29. Walking the Talk: Adopting and Adapting Sustainable Scientific Software Development processes in a Small Biology Lab

Michael R Crusoe, CTitus Brown

Journal of Open Research Software (2016-11-29) https://doi.org/gsmb78

DOI: 10.5334/jors.35 · PMID: 27942385 · PMCID: PMC5142744

30. The pragmatic programmer: from journeyman to master

Andrew Hunt, David Thomas

Addison-Wesley (2000) ISBN: 9780201616224