

STATE OF PRACTICE FOR MEDICAL IMAGING SOFTWARE

ASSESSING THE CURRENT STATE OF THE PRACTICE FOR MEDICAL IMAGING
SOFTWARE

By
AO DONG.

A Thesis
Submitted to the School of Graduate Studies
in Partial Fulfillment of the Requirements
for the degree
Master of Engineering in Computing and Software

McMaster University

© Copyright by Ao Dong, August 2021

MASTER OF ENGINEERING (2021)
(Computing and Software)

McMaster University
Hamilton, Ontario

TITLE: Assessing the Current State of the Practice for Medical Imaging Software

AUTHOR: Ao Dong

SUPERVISOR: Dr. Spencer Smith

NUMBER OF PAGES: viii, ??

Abstract

Abstract here

Acknowledgments

acknowledgements here

Contents

| | |
|---|------------|
| Abstract | iii |
| Acknowledgments | iv |
| 1 Introduction | 2 |
| 1.1 Motivation | 2 |
| 1.2 Purpose | 2 |
| 1.3 Scope | 2 |
| 2 Background | 3 |
| 2.1 Software Categories | 3 |
| 2.1.1 Open Source Software | 3 |
| 2.1.2 Freeware | 4 |
| 2.1.3 Commercial Software | 4 |
| 2.1.4 Scientific Computing Software | 4 |

| | | |
|----------|--|-----------|
| 2.2 | Software Quality Definitions | 5 |
| 2.3 | Analytic Hierarchy Process | 6 |
| 3 | Methods | 9 |
| 3.1 | Domain Selection | 10 |
| 3.2 | Software Product Selection | 11 |
| 3.2.1 | Identify Software Candidates | 11 |
| 3.2.2 | Filter the Software List | 12 |
| 3.3 | Grading Template | 14 |
| 3.4 | Empirical Measurements | 15 |
| 3.5 | Measuring Qualities | 16 |
| 3.6 | Interview Methods | 17 |
| 3.6.1 | Interviewee Selection | 17 |
| 3.6.2 | Interview Question Selection | 18 |
| 3.6.3 | Interview Methods | 19 |
| 4 | Measurement Results | 20 |
| 4.1 | Selected Software List | 21 |
| 4.2 | Installability | 21 |
| 4.3 | Correctness & Verifiability | 21 |
| 4.4 | Reliability | 21 |

| | | |
|----------|-------------------------------------|-----------|
| 4.5 | Robustness | 21 |
| 4.6 | Usability | 21 |
| 4.7 | Maintainability | 21 |
| 4.8 | Reusability | 21 |
| 4.9 | Understandability | 21 |
| 4.10 | Visibility & Transparency | 21 |
| 5 | Interviews with Developers | 22 |
| 5.1 | Summary of Answers | 22 |
| 5.2 | Discussions | 22 |
| 6 | Threat to Validity | 23 |
| 7 | Recommendations | 24 |
| 8 | Conclusions | 25 |
| | Bibliography | 26 |
| A | Full Grading Template | 30 |
| B | Summary of Measurements | 31 |
| C | Interview Answers | 32 |

Chapter 1

Introduction

introduction

scientific computing (SC), also known as scientific computation or computational science

1.1 Motivation

1.2 Purpose

1.3 Scope

Chapter 2

Background

2.1 Software Categories

In this section, we discuss three common software categories that are mentioned in Section 3.2, and also SC software.

2.1.1 Open Source Software

For Open Source software (OSS), its source code is openly accessible, and users have the right to study, change and distribute it under a license granted by the copyright holder. For many OSS projects, the development process is based on the collaboration of different contributors worldwide [4].

2.1.2 Freeware

Freeware is software that can be used free of charge. Unlike with OSS, the authors of freeware typically do not allow users to access or modify the source code of the software [13]. The term *freeware* should not be confused with *free software*, which is similar to OSS but with a few differences.

2.1.3 Commercial Software

“Commercial software is software developed by a business as part of its business” [7]. Typically speaking, the users are required to pay to access all of the features of commercial software, excluding access to the source code. However, some commercial software is also free of charge [7].

2.1.4 Scientific Computing Software

Software development in Scientific Computing (SC) depends on the knowledge of three areas - the knowledge of a specific engineering or science domain, the ability to mathematically build models and applying algorithms, and the capability to implement theoretical models and algorithms with computational tools. SC software is built with mathematical and computational tools to serve the purpose of solving scientific problems in a domain [11].

2.2 Software Quality Definitions

The definitions of software qualities are from Smith et al. [16]. The order of the qualities follows the grading template in Appendix A.

- **Installability** The effort required for the installation, uninstallation, or reinstallation of a software or product in a specified environment.
- **Correctness & Verifiability** A program is correct if it behaves according to its stated. Verifiability is the extent to which a set of tests can be written and executed, to demonstrate that the delivered system meets the specification.
- **Reliability** The probability of failure-free operation of a computer program in a specified environment for a specified time, i.e. the average time interval between two failures also known as the mean time to failure (MTTF).
- **Robustness** Software possesses the characteristic of robustness if it behaves “reasonably” in two situations: i) when it encounters circumstances not anticipated in the requirements specification, and ii) when the assumptions in its requirements specification are violated.
- **Usability** The extent to which a product can be used by specified users to achieve specified goals with effectiveness, efficiency, and satisfaction in a specified context of use.

- **Maintainability** The effort with which a software system or component can be modified to i) correct faults; ii) improve performance or other attributes; iii) satisfy new requirements.
- **Reusability** The extent to which a software component can be used with or without adaptation in a problem solution other than the one for which it was originally developed.
- **Understandability** (To be completed)
- **Visibility & Transparency** The extent to which all of the steps of a software development process and the current status of it are conveyed clearly.

2.3 Analytic Hierarchy Process

The Analytic Hierarchy Process (AHP) was developed by Thomas L. Saaty, and it has been widely used to make and analyze multiple criteria decisions [18]. The AHP organizes multiple criteria factors in a hierarchical structure and pairwise compares the alternatives to calculate relative ratios [14].

For a project with m criteria, we can use a $m \times m$ matrix A to record the relative importance between factors. By pairwise compare criterion i and criterion j , the value of A_{ij} is decided as follows, and the value of A_{ji} is $1/A_{ij}$ [14],

- $A_{ij} = 1$ if criterion i and criterion j are equally important;
- $A_{ij} = 9$ if criterion i is extremely more important than criterion j ;
- A_{ij} equals to an integer value between 1 and 9 according the the relative importance of criterion i and criterion j .

The above process assumes that criterion i is not less important than criterion j , otherwise, we need to reverse i and j and determine A_{ji} first, then $A_{ij} = 1/A_{ji}$.

The priority vecotr w can be calculated by solving the following equation [14],

$$Aw = \lambda_{max}w, \quad (2.1)$$

where λ_{max} is the maximal eigenvalue of A .

In this project, w is approximated with the approach classic *mean of normalized values* [10],

$$w_i = \frac{1}{m} \sum_{j=1}^m \frac{A_{ij}}{\sum_{k=1}^m A_{kj}} \quad (2.2)$$

Suppose there are n alternatives, for criterion $i = 1, 2, \dots, m$, we can create an $n \times n$ matrix B_i to record the relative preferences between these choices. The way of generating B_i is similar to the one for A . However, unlike comparing the importance between criteria, we pairwise decide how much one alternative is more favored than the other. The same method is used to calculate the local priority vector for each B_i .

In this project, the 9 software qualities mentioned above are the criteria ($m = 9$), while 29 software packages ($n = 29$) are compared. The software are evaluated with the grading template in Appendix A and a subjective score is given for each quality. For a pair of qualities or software, i and j , such that i is not less significant than j , the pairwise comparison result of i versus j is converted from $\min((score_i - score_j) + 1, 9)$.

Chapter 3

Methods

[Where to include the domain experts? I plan to include them in Domain Selection and Software Product Selection, and also mention that they'll participate in interviews. —AD]

In this project, we aimed to study and research the SC software within scientific domains. However, we tried to design the whole process of this project in a more general way, with the hope that it should not be limited to only one or several specific software categories.

The way how we chose a domain is listed in Section 3.1. Within a specific domain, we used the process documented in Section 3.2 to collect and filter the software packages for our study. Then, all the software products were measured by using the grading template in Section 3.3 and the empirical measurement method in Section . Section 3.5 presents more details about how we applied the above process.

3.1 Domain Selection

Although the methods designed by us may not have such limitations, we limited our selection of a domain within scientific domains to fulfill the objective of our research.

One major factor in properly choosing a candidate domain for the study is the ease to select software packages within it. As described in Section 3.2, the software product selection process is most likely to have a screening step, and the quantity of final-decision packages may not meet our initial expectation if we do not have enough software candidates to choose from. Section 3.2 also explains why our process prioritize the OSS. Consequently, a scientific domain with a large number of active OSS is proffered. This also often indicates that the domain has an active community developing and using SC software, making it easier to conduct interviews mentioned in Section 3.6.

Even if we can find an adequate number of software packages in a domain, we may still find the software products are developed to solve various problems within the domain and should be categorized into different sub-groups. So one question needs to be asked - do we prefer a group of software all providing similar functions and features, or do we aim to cross-compare several sub-sections within the same domain? With that answered, it should be easier to determine a favored domain.

Another aspect to consider is the team carrying out the research, more specifically, the domain experts - if there is any - in the team. In our team, the projects are often led by researchers in the software engineering field and supported by experts working in other

scientific domains. Having domain experts in the team provides significant benefits in selecting software packages and designing interview questions.

In this project, Medical Imaging (MI) domain was selected. Numerous software products can be found in this domain and a great number of professionals working in this domain use and/or develop such software. We decided to focus on MI software with the viewing function, and more details about this filtering are in Section 3.2. Being able to include MI domain experts in our team also made this domain more preferred.

3.2 Software Product Selection

The process of selecting software packages contains two steps: i) identify software candidates in the chosen domain, ii) filter the list according to needs [17].

3.2.1 Identify Software Candidates

The candidate software can be found in publications in the domain. Another source is to search various websites, such as GitHub, swMATH and the Google search results for software recommendation articles. Meanwhile, we should also include the suggested ones from domain experts [17].

As for this project, 48 MI software projects were identified as the candidates, and they were found from publications [1] [3] [8], online articles related to the domain [5] [9] [12],

forum discussions related to the domain [15], etc.

3.2.2 Filter the Software List

The goal is to build a software list with a length of about 30 [17].

The only “mandatory” requirement is that the software must be OSS, as defined in Section 2.1.1. This is due to the grading process defined in Section 3.3. To evaluate all aspects of some software qualities, the source code should be accessible.

The other factors to filter the list can be optional and should be considered according to the number of software candidates and the objectives of a research project.

One of the factors is the functions and purposes of the software. For example, we can choose a group of software with similar functions, so that the cross-comparison is between each individual of them. On the other hand, if our objective is comparing sub-categories in the domain, we should select from candidates in each of the categories.

The empirical measurement tools listed in Section 3.4 are limited to projects using Git as the version control tool, so software with Git is preferred. Some manual steps in empirical measurement depend on a few metrics of GitHub, which makes projects held on GitHub more favored [17].

Some of the OSS projects may experience a lack of recent maintenance. So packages that have not been updated for a long time can be eliminated, unless they are still popular and highly recommended by the users in the domain [17].

No matter what standards are set to filter the packages, with domain experts in the team, their opinion should be valued. For example, if a software package is not OSS and has not been updated for a long while, the domain experts may still identify it as a popular and widely used product. In this case, perhaps it is valuable to keep this software on the list.

In this project, among the 48 software packages identified in Section 3.2.1, there were several ones that we could not find their source code, such as MicroDicom, Aliza, and jivex, etc. So these packages are likely to be freeware defined in Section 2.1.2 and were removed from the list since they might not be OSS.

We focused on the MI software that could work as a MI viewer. Some of the software on the list were tool kits or libraries for other software to use as dependencies but not for end-users to view medical images, such as VTK, ITK, and dcm4che, etc. These were also eliminated from the list.

After the above two steps, there were 29 software products left on the list. We still preferred projects using git and GitHub and being updated recently, but no other filtering was needed since the number of packages was already below 30. However, 27 out of the 29 software packages on the final list used git, and 24 of which were held on GitHub. Furthermore, 27 packages had the latest updates after the year 2018, and 23 after 2020.

3.3 Grading Template

The full grading template can be found in Appendix A. The template contains 101 questions that we used for grading software products.

In the first section of the template, some general information about the software is collected, such as the name, purpose, platform, programming language, publications about the software, the first release and the most recent change date, website, and source code repository of the product, etc. Information in this section helps us to understand the projects better and may be useful for further analysis, but it does not directly affect the grading scores for the packages.

The next 9 sections in the template are designed for the 9 software qualities mentioned in Section 2.2. For each quality, several questions are asked and the typical choices of answers are “yes”, “no”, “n/a”, “unclear”, a number, a string, a date, a set of strings, etc. The last question of each quality is asking for an overall score between 1 and 10, which is based on all the previous questions of this section. This score is also the grading score for this quality. For some qualities, this grading score is also affected by the empirical measurement sections on the template.

All the last 3 sections on the template are related to the empirical measurements. Two command-line software tools `git_stats` and `scc` are used to extract information about the source code from the project repositories. For projects held on GitHub, additional metrics are collected manually, such as the stars of the GitHub repository, and the numbers of open

and closed pull requests. Details of how these empirical measurements affect software quality grading scores are presented in Section 3.4.

3.4 Empirical Measurements

Two command-line tools were used for the empirical measurements in this project. One is GitStats that generates statistics for git repositories and display outputs in the format of web pages [6]; the other one is Sloc Cloc and Code (abbreviated as scc by the author) [2], aiming to count the lines of code, comments, etc.

Both tools can measure the number of text-based files in a git repository, as well as lines of text in these files. Based on our experience, most text-based files in a software project repository contain programming source code and are used to compile and build software products, and a minority of these files are instructions and other documents. So the lines of text in text-based files are roughly regarded as lines of programming code. The two tools usually generate similar but not identical results as for the above measurements. From our understanding, this minor difference is because of the different techniques they use to detect if a file is a text-based or binary file.

Additionally, we also manually collected some information for projects held on GitHub, such as the numbers of stars, forks, people watching this repository, open pull request, and closed pull request.

These empirical measurements helped us from two aspects. Firstly, with more statistical details, we could get an overview of a project faster and more accurately. For example, the number of commits over the last 12 months shows how active this project has been during this period, and the number of stars and forks may be related to the popularity of it. Secondly, the results were factors to consider when we determined the grading scores for some software qualities. For example, if the percentage of comment lines is low, we might want to double-check the understandability of the code, and if the ratio of open versus closed pull requests is high, perhaps we should pay more attention to the maintainability of the project.

3.5 Measuring Qualities

27 out of the 29 software packages can be installed on two or three different operating systems such as Windows, macOS, and Linux, and 5 of them are browser-based, making them platform-independent. However, in the interest of time, we only performed the measurements for each project by installing it on one of the platforms, most likely Windows.

In order to test the software on a “clean” system, we created a new virtual machine (VM) for each software and only installed the necessary dependencies before measuring. All the 29 VM were created on the same computer. We only started the measuring of the next software after finishing a current one, and after grading each software, the VM was

destroyed.

Generally speaking, about two hours were spent on grading one software, unless there were technical issues and more time was need to resolve them. For most of the situation, all the measurements for one software were finished on the same day.

The results were filled into the grading template alongside measuring a software in a VM. The sequence of measuring also followed the grading template.

3.6 Interview Methods

3.6.1 Interviewee Selection

For a software list with a length of roughly 30, we aim to interview about 10 development teams of these projects. Interviewing multiple individuals from each team should give us more comprehensive information about a project, but due to the difficulties of finding willing participants, a single engineer well knowing the project can also be sufficient.

Ideally, we select projects after the grading measurements are done and prefer the projects with higher overall scores. However, if not enough participants are found, we should also contact all teams on the list.

For MI domain, we started with the teams with higher score on our list, and eventually contacted all of them. There are developers/architects from 8 teams having participated our interviews so far.

The contacts of the the teams were found on the websites related to the software. For example, the official web pages, repository websites, publication websites, and bio pages of the teams' institutions. For each candidate, we sent at most two emails asking for their support and participating before receiving any replies.

3.6.2 Interview Question Selection

We have a list of 20 questions to guide our interviews with the development teams, which can be found in Appendix C.

Some of the questions are about the background of the software, the development teams, the interviewees, and the way they organize the projects. We also ask about their understandings to the users. Another part of the interview questions focuses on the major difficulties the team experience currently or in the past, and the solutions that they have found or will try in the future. We also discuss with interviewees the importance of documents to their projects and the current situations of these documents. One more proportion of the questions are about several specific software qualities, such as maintainability, understandability, usability, and reproducibility.

The interviews are supposed to be semi-structured based on the question list and follow-up questions can be asked. Based on our experience, the interviewees usually bring up some interesting topics not expected by us and we found it valuable to continue on these topics and ask for a few more details.

3.6.3 Interview Methods

Before contacting any interviewee candidate, the study had been reviewed by the McMaster University Research Ethics Board and received ethics clearance.

The members of the development teams are based around the world, so we organized these interviews as virtual meetings online. Zoom was used for the meetings. After receiving consents from the interviewees, we also recorded our discussions to better transcribe them.

Chapter 4

Measurement Results

For some qualities, it might be a good idea to cross-compare with the empirical scores.

4.1 Selected Software List

4.2 Installability

4.3 Correctness & Verifiability

4.4 Reliability

4.5 Robustness

4.6 Usability

4.7 Maintainability

4.8 Reusability

4.9 Understandability

4.10 Visibility & Transparency

Chapter 5

Interviews with Developers

5.1 Summary of Answers

- Start with one by one, with commonalities and interesting special cases.
- Shorten and summarize later.

5.2 Discussions

Any conclusions?

Chapter 6

Threat to Validity

Chapter 7

Recommendations

I think the recommendations can originate from both parts - measurements and interviews.

Chapter 8

Conclusions

No clues yet. Should be started at a later stage.

Bibliography

- [1] Kari Björn. Evaluation of open source medical imaging software: A case study on health technology student learning experience. *Procedia Computer Science*, 121:724–731, 01 2017.
- [2] Ben Boyter. Sloc cloc and code. <https://github.com/boyter/scc>, 2021. [Online; accessed 27-May-2021].
- [3] Andreas Brühshwein, Julius Klever, Anne-Sophie Hoffmann, Denise Huber, Elisabeth Kaufmann, Sven Reese, and Andrea Meyer-Lindenberg. Free dicom-viewers for veterinary medicine: Survey and comparison of functionality and user-friendliness of medical imaging pacs-dicom-viewer freeware for specific use in veterinary medicine practices. *Journal of Digital Imaging*, 03 2019.
- [4] James Edward Corbly. The free software alternative: Freeware, open source software, and libraries. *Information Technology and Libraries*, 33(3):65–75, Sep. 2014.

- [5] Steve Emms. 16 best free linux medical imaging software. <https://www.linuxlinks.com/medicalimaging/>, 2019. [Online; accessed 02-February-2020].
- [6] Tomasz Gieniusz. Gitstats. https://github.com/tomgi/git_stats, 2019. [Online; accessed 27-May-2021].
- [7] GNU. Categories of free and nonfree software. <https://www.gnu.org/philosophy/categories.html>, 2019. [Online; accessed 20-May-2021].
- [8] Daniel Haak, Charles-E Page, and Thomas Deserno. A survey of dicom viewer software to integrate clinical research and medical imaging. *Journal of digital imaging*, 29, 10 2015.
- [9] Mehedi Hasan. Top 25 best free medical imaging software for linux system. <https://www.ubuntupit.com/top-25-best-free-medical-imaging-software-for-linux-system/>, 2020. [Online; accessed 30-January-2020].
- [10] Alessio Ishizaka and Markus Lusti. How to derive priorities in ahp: A comparative study. *Central European Journal of Operations Research*, 14:387–400, 12 2006.
- [11] Hemant Kumar Mehta. *Mastering Python scientific computing: a complete guide for Python programmers to master scientific computing using Python APIs and tools*. Packt Publishing, 2015.

- [12] Hamza Mu. 20 free & open source dicom viewers for windows. <https://medevel.com/free-dicom-viewers-for-windows/>, 2019. [Online; accessed 31-January-2020].
- [13] The Linux Information Project. Freeware definition. <http://www.linfo.org/freeware.html>, 2006. [Online; accessed 20-May-2021].
- [14] Thomas L. Saaty. How to make a decision: The analytic hierarchy process. *European Journal of Operational Research*, 48(1):9–26, 1990. Decision making by the analytic hierarchy process: Theory and applications.
- [15] Ravi Samala. Can anyone suggest free software for medical images segmentation and volume? https://www.researchgate.net/post/Can_anyone_suggest_free_software_for_medical_images_segmentation_and_volume, 03 2014. [Online; accessed 31-January-2020].
- [16] Spencer Smith, Jacques Carette, Olu Owajaiye, Peter Michalski, and Ao Dong. Quality definitions of qualities. Manuscript in preparation, 2020.
- [17] Spencer Smith, Jacques Carette, Olu Owajaiye, Peter Michalski, and Ao Dong. Methodology for assessing the state of the practice for domain x. Manuscript in preparation, 2021.

- [18] Omkarprasad S. Vaidya and Sushil Kumar. Analytic hierarchy process: An overview of applications. *European Journal of Operational Research*, 169(1):1–29, 2006.

Appendix A

Full Grading Template

appendix here

Appendix B

Summary of Measurements

appendix here

Appendix C

Interview Answers

appendix here

Appendix D

Ethics Approval

appendix here