

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	4
2.1.2 Freeware	4
2.1.3 Commercial Software	5
2.1.4 Scientific Computing Software	5

2.2	Software Quality Definitions	6
2.3	Analytic Hierarchy Process	7
3	Methods	10
3.1	Domain Selection	10
3.2	Software Product Selection	11
3.2.1	Identify Software Candidates	12
3.2.2	Filter the Software List	12
3.3	Grading Software	13
3.3.1	Grading Template	13
3.3.2	Empirical Measurements	14
3.3.3	Technical Details	15
3.4	Interview Methods	16
3.4.1	Interviewee Selection	16
3.4.2	Interview Question Selection	16
3.4.3	Interview Process	17
3.5	Applying the Method to Medical Imaging (MI)	17
4	Measurement Results	20
4.1	Overall Scores	22
4.2	Installability	23

4.3	Correctness & Verifiability	25
4.4	Surface Reliability	27
4.5	Surface Robustness	28
4.6	Surface Usability	30
4.7	Maintainability	31
4.8	Reusability	33
4.9	Surface Understandability	34
4.10	Visibility/Transparency	35
5	Interviews with Developers	37
5.1	Current and Past Pain Points	38
5.1.1	Pain Points in Group 1	40
5.1.2	Pain Points in Group 2	41
5.1.3	Pain Points in Group 3	44
5.2	Documents in the Projects	45
5.3	Contribution Management and Project Management	47
6	Threat to Validity	52
7	Recommendations	53
8	Conclusions	54

Bibliography	55
A Full Grading Template	61
B Summary of Measurements	62
C Other Interview Answers	63
D Ethics Approval	81

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

This chapter introduces several different categories of software, based on which we designed the processes to select the software domain and software candidates in Chapter 3. It also covers the software quality definitions used in the grading template in Appendix A and an overview of the Analytic Hierarchy Process (AHP), a tool we used to compare and grade software products.

2.1 Software Categories

We usually target specific software categories to narrow down the scopes when selecting software domains and software packages. In this section, we discuss three common software categories that are mentioned in Section 3.2, and also SC software. Section 3.1 and Section 3.2 present more details about why we prefer some of these categories.

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 relies on the collaboration of different contributors worldwide [6]. Accessible source code usually exposes more “secrets” of a software project, such as the underlying logic of software functions, how developers achieve their works, and the flaws and potential risks in the final product. Thus, it brings much more convenience to the researchers analyzing the qualities of the project.

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 [24]. The term *freeware* should not be confused with *free software*, which is similar to OSS but with a few differences. To the end-users, the differences between freeware and OSS often do not bother them. The fact that these products are free of charge is likely to make them popular with many users. However, software developers, end-users who wish to modify the source code, and researchers looking for inner characteristics may find the inaccessible source code to be a problem.

2.1.3 Commercial Software

“Commercial software is software developed by a business as part of its business” [9]. 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 [9]. Based on our experience, most commercial software products are not OSS.

For some specific software, the backgrounds of commercial software developers often differ from the ones of non-commercial OSS. In such a case, the former is usually the product of software engineers, and the latter is likely to have developers who work in the domain and are also end-users of the products. One example is software in Scientific Computing (SC), since the developers need to utilize their domain-specific during the development process [32].

2.1.4 Scientific Computing Software

Software development in SC depends on the knowledge of three areas - the inside of a specific engineering or science domain, the ability to mathematically build models and apply algorithms, and 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 [22]. However, the majority of scientists developing their software are self-taught programmers [32], so there may be a bigger Room for

software quality improvement in SC domains.

2.2 Software Quality Definitions

The definitions of software qualities are from Smith et al. [29]. 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

To generate grading scores for a group of software packages, we use the AHP to pairwise compare them. Thomas L. Saaty developed this tool, and people widely used it to make and analyze multiple criteria decisions [31]. The AHP organizes multiple criteria factors in a hierarchical structure and pairwise compares the alternatives to calculate relative ratios [27].

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 $1/A_{ij}$ [27],

- $A_{ij} = 1$ 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 [27],

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

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

In this project, w is approximated with the approach classic *mean of normalized values* [14],

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

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 we favor one alternative over the other. We use the same method 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

We designed a general process for evaluating the state of the practice of domain-specific software, that we instantiate to SC software for specific scientific domains.

Our method involves: 1) choosing a software domain (Section 3.1); 2) collecting and filtering software packages (Section 3.2); 3) grading the selected software (Section 3.3); 4) interviewing development teams for further information (Section 3.4).

Details of how we applied the method on the MI domain are expanded upon in Section 3.5.

3.1 Domain Selection

Our methods are generic, but we have only applied them to scientific domains due to the objective of our research.

When choosing a candidate domain, we prefer one with a large number of active OSS. The reason is that we aim to finalize a list of 30 software packages [30] after the screening step. For example, we remove the ones without recent updates or specific functions. Thus, we need enough software candidates from the beginning. Besides, we prefer OSS projects because our grading method works better on them. Moreover, we prefer domains with active communities developing and using SC software, so it is easier to invite enough developers for interviews.

We prefer 30 software packages providing similar functions or falling into different sub-categories depending on our research purpose. So the domain needs to have enough candidates in one sub-category or enough sub-categories to cross-compare.

We also prefer domains in which our team has expertise. We invite domain experts to join and support our projects. They help us in many aspects, such as vetting the software list and interview questions.

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 [30].

3.2.1 Identify Software Candidates

We start with finding candidate software in publications of the domain. Then, we search various websites, such as GitHub, swMATH and the Google search results for software recommendation articles. Meanwhile, we include the suggested ones from domain experts [30].

3.2.2 Filter the Software List

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

The only “mandatory” requirement is that the software must be OSS, as defined in Section 2.1.1. Because to evaluate all software qualities, our methods need to access the source code.

The other factors to filter the list are optional, and we consider them according to the number of software candidates and the objectives of the research project.

One factor is the functions and purposes of the software. For example, we can choose a group of software with similar functions, so that the comparison is between the software; or we can cross-compare sub-categories in the domain, then we need to select candidates from each sub-category.

The empirical measurement tools listed in Section 3.3.2 only work on projects using Git, so we prefer software with Git. Some manual steps in empirical measurement depend on a few metrics of GitHub, which makes projects held on GitHub more favored [30].

Some of the OSS projects may experience a lack of recent maintenance. So we eliminate packages without recent updates, unless they are still popular and highly recommended by the domain users [30].

With domain experts in the team, we value their opinions on the filtering process. For example, if a software package is not OSS and has no updates for a long while, but the domain experts identify it as a valuable product, we still consider it to keep.

3.3 Grading Software

We grade the selected software using a template (Section 3.3.1) and a specific empirical method (Section 3.3.2). Some technical details for the measurements are in Section 3.3.3.

3.3.1 Grading Template

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

We use the first section of the template to collect general information, 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 understand the projects better and may be helpful for further analysis, but it does not directly affect the grading scores.

We designed the following nine sections in the template for the nine software qualities mentioned in Section 2.2. For each quality, we ask several questions and the typical answers are among the collection of “yes”, “no”, “n/a”, “unclear”, a number, a string, a date, a set of strings, etc. Each quality needs an overall score between 1 and 10 based on all the previous questions. For some qualities, the empirical measurements also affect the score.

All the last three sections are about the empirical measurements. We use two command-line software tools *GitStats* and *scc* to extract information from the source code repositories. For projects held on GitHub, we manually collect additional metrics, such as the stars of the GitHub repository, and the numbers of open and closed pull requests. Section 3.3.2 presents more details about the empirical measurements.

3.3.2 Empirical Measurements

We use two command-line tools for the empirical measurements. One is *GitStats* that generates statistics for git repositories and display outputs in the format of web pages [8]; the other one is Sloc Cloc and Code (as known as *scc*) [4], aiming to count the lines of code, comments, etc.

Both tools measure the number of text-based files in a git repository and lines of text in these files. Based on our experience, most text-based files in a repository contain programming source code, and developers use them to compile and build software products. A minority of these files are instructions and other documents. So we roughly regard the lines

of text in text-based files as lines of programming code. The two tools usually generate similar but not identical results. From our understanding, this minor difference is due to the different techniques to detect if a file is text-based or binary.

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

These empirical measurements help us from two aspects. Firstly, they help us with getting a project overview faster and more accurately. For example, the number of commits over the last 12 months shows how active this project has been, and the number of stars and forks may reveal its popularity. Secondly, the results may affect our decisions regarding the grading scores for some software qualities. For example, if the percentage of comment lines is low, we double-check the *understandability* of the code; if the ratio of open versus closed pull requests is high, we pay more attention to the *maintainability*.

3.3.3 Technical Details

To test the software on a “clean” system, we create a new virtual machine (VM) for each software and only install the necessary dependencies before measuring. We make all 30 VM on the same computer, one at a time, and destroy them after measuring.

We spend about two hours grading one package, unless we find technical issues and need more time to resolve them. In most of the situation, we finish all the measurements

for one software on the same day.

3.4 Interview Methods

3.4.1 Interviewee Selection

For a software list with a length of roughly 30, we aim to interview about ten development teams. Interviewing multiple individuals from each team gives us more comprehensive information, but a single engineer well-knowing the project is also sufficient.

Ideally, we select projects after the grading measurements and prefer the ones with higher overall scores. However, if we do not find enough participants, we contact all teams on the list.

We try to find the contacts of the teams on the projects' websites, such as the official web pages, repositories, publications, and bio pages of the teams' institutions. Then, we send at most two emails to one contact asking for its participation before receiving any replies.

3.4.2 Interview Question Selection

We have a list of 20 questions to guide our interviews, which can be found in Section 5 and Appendix C.

Some questions are about the background of the software, the development teams, the

interviewees, and how they organize the projects. We also ask about their understandings of the users. Some questions focus on the current and past difficulties, and the solutions the team has found or will try. We also discuss the importance and current situations of documentation. A few questions are about specific software qualities, such as *maintainability*, *understandability*, *usability*, and *reproducibility*.

The interviews are semi-structured based on the question list, and we also ask follow-up questions when necessary. Based on our experience, the interviewees usually bring up some exciting ideas that we do not expect, and it is worth expanding on these topics.

3.4.3 Interview Process

Since the members of the development teams are usually around the world, we organize these interviews as virtual meetings online with Zoom. After receiving consent from the interviewees, we also record our discussions to transcribe them better.

3.5 Applying the Method to Medical Imaging (MI)

Based on the principles in Section 3.1, we selected the MI domain and the viewer sub-category. We also included MI domain experts in our team.

By using the method in Section 3.2.1, we identified 48 MI software projects as the candidates from publications [3] [5] [10], online articles related to the domain [7] [13]

[23], forum discussions related to the domain [28], etc.

Among the 48 software packages, there were several ones that we could not find in their source code, such as *MicroDicom*, *Aliza*, and *jivex*, etc. These packages are likely to be freeware defined in Section 2.1.2 and not OSS. So following guidelines in Section 3.2.2 we removed them from the list. Next, 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*, *dcm4che*, etc. We also eliminated these from the list. After that, there were 29 software products left on the list. We still preferred projects using git and GitHub and being updated recently, but did not apply another filtering 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 chose GitHub. Furthermore, 27 packages had the latest updates after the year 2018, and 23 after 2020.

Then we followed the steps in Section 3.3 to measure and grade the software. 27 out of the 29 packages are compatible with 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.

Going through the interview process in Section 3.4, we started with the teams with higher scores on our list, and eventually contacted all of them. As a result, develop-

ers/architects from 8 teams have participated in our interviews so far. Before contacting any interviewee candidate, we received ethics clearance from the McMaster University Research Ethics Board (Appendix D).

Chapter 4

Measurement Results

A summary of the measurement results can be found in Appendix B, and the detailed data is in the repository <https://github.com/smiths/AIMSS>.

All of the 29 software packages that we measured are OSS with functions of displaying medical images. The initial release dates of these packages are between the years 1998 and 2018. We considered 2 of them as “dead” projects since they had the latest updates more than 18 months earlier than the time of measuring, but both had updates after 2017.

We found out that at least 8 of the projects received funding for at least a certain period, but we were not sure about the rest of them due to limited information.

After analyzing the repositories and official websites of the projects, we found out that the number of contributors for a project varied between 1 to roughly 100. We considered anyone who made at least one accepted commit to the source code as a contributor, so

it does not mean that any development teams had long-term members as many as 100. Many of these projects received change requests and code from the community, such as pull requests and git commits on GitHub. At least 27 of the software projects had code commits from more than a single developer, and as far as we could find out, at least 13 of them had no less than 10 contributors.

According to our measurements, 27 out of the 29 packages supported multiple operating systems, and 25 of them could work on all three Windows, macOS, and Linux systems. However, there was a significant difference in the philosophy to achieve cross-platform compatibility. The majority of the 27 were desktop software products built for each operating system separately. On the other hand, 5 were web applications, which made them naturally platform-independent and compatible with almost all operating systems with compatible browsers. There are certain advantages and disadvantages to developing either a native or a browser-based application, and some of the details can be found from the interview answers in Section ??.

Most of the projects used more than one programming language, including a primary language that the developers used the most. The most popular primary language was C++, chosen by 11 projects, and another 3 teams selected C. The second most used one was JavaScript, with 6 projects using it as the primary language. 4 of the projects chose Java, and 3 selected Python.

4.1 Overall Scores

As described in Section 2.3, for our AHP measurements, there are 9 criteria which are the 9 software qualities and 29 software packages as the alternatives. We decided to make all 9 qualities equally important, so the score of each quality affects the overall scores on the same scale.

Figure 4.1 shows the overall scores of all 29 software packages in descending order. Since we produced the scores from the AHP process, the total sum of the 29 scores is precisely 1.

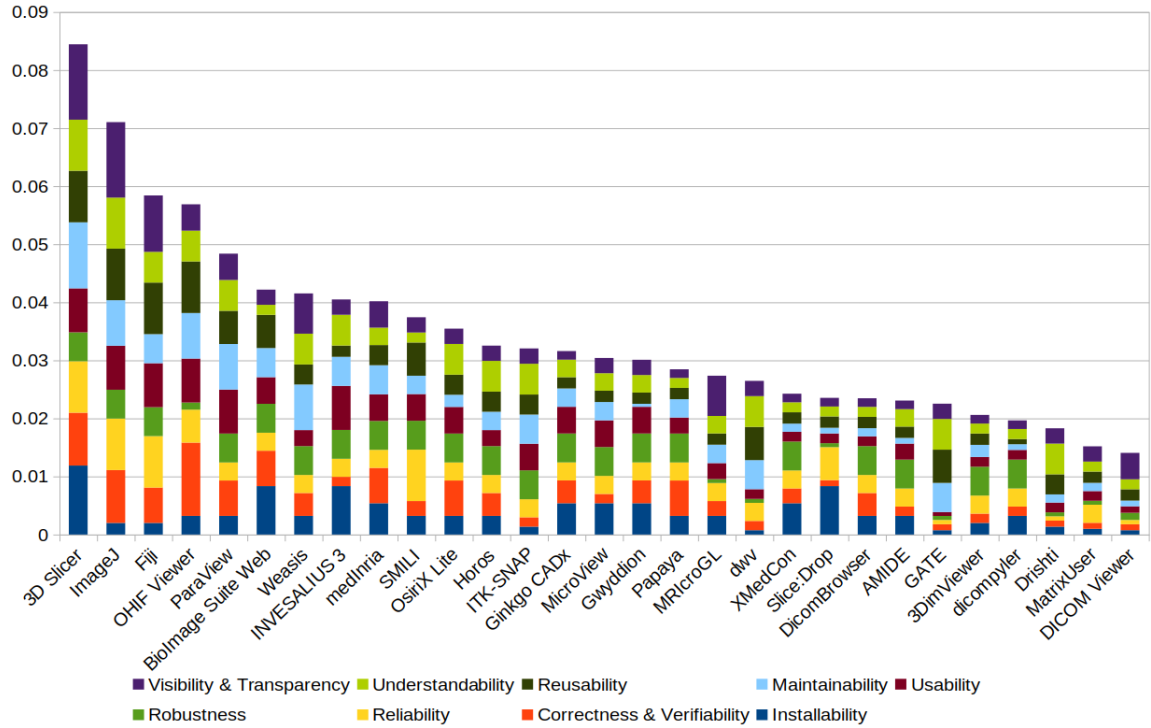


Figure 4.1: Overall AHP scores for all 9 software qualities

The top two software products had higher scores in each criterion. For example, *3D Slicer* [17] ranked at the first or second place for 8 of the 9 qualities; *ImageJ* [26] was within the top three products for 7 qualities. There were 3 software packages that we could not install or build correctly. Among them, *DICOM Viewer* [1] was the only one that did not have an online test version, so that we could not finish all measurements for it. Consequently, we might underestimate its score.

4.2 Installability

When measuring the *installability*, we checked the existence and quality of installation instructions. The user experience was also an important factor, such as the ease to follow the instructions, number of steps, automation tools, and the prerequisite steps for the installation. If any problem interrupted the process of installation or uninstallation, we gave a lower score to this quality. We also recorded the operating system for the installation test and whether we could verify the installation.

Figure 4.2 lists the scores of *installability*, with a total sum of 1.

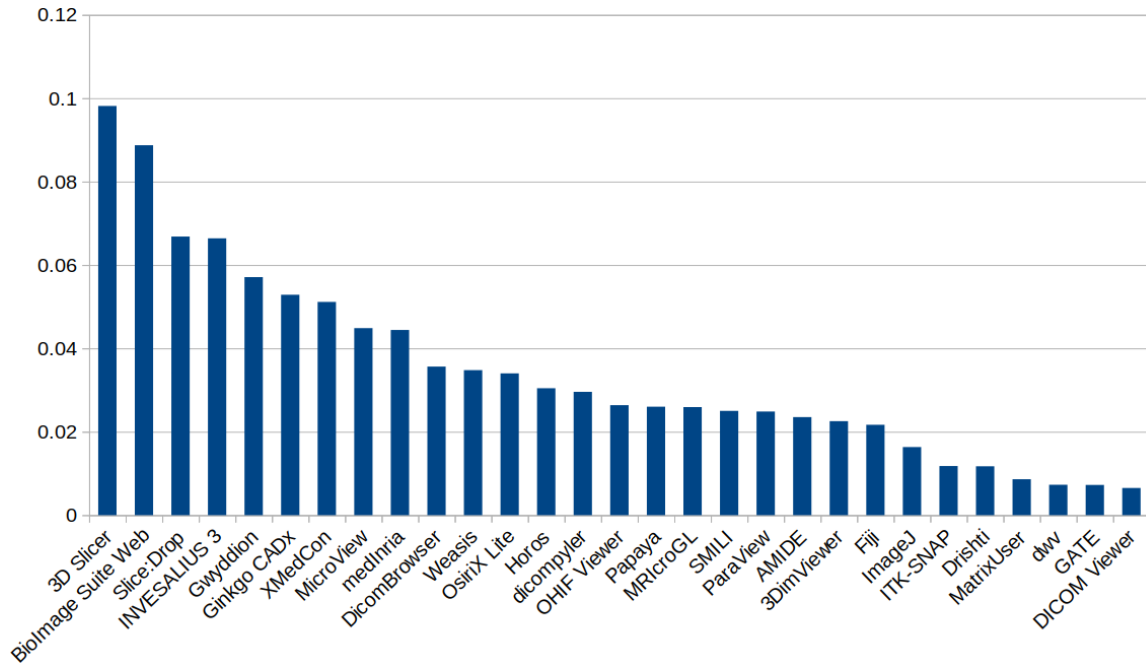


Figure 4.2: AHP installability scores

The projects with higher scores had easy-to-follow installation instructions, and the installation processes were automated, fast, and frustration-free, with all dependencies automatically added. There were also no errors during the installation and uninstallation steps.

The ones with the lowest scores often showed severe problems. We could not install some of them due to some issues that we could not solve. We spent a reasonable amount of time on these problems, then considered them major obstacles for normal users if we still did not figure out any solutions. For example, we needed to build some software from

the source code, and we failed the building for two of them following the instructions. Furthermore, *DICOM Viewer* depended on a cloud platform, and we could not successfully install the dependency. We suspect that only a part of the users faced the same problems, and given a lot of time, we might be able to find solutions. However, the difficulties greatly impacted the installation experiences, and we graded these software packages with lower scores.

Although we could not locally build two software packages, there were deployed online versions for them. With that, we finished all the measurements for them.

4.3 Correctness & Verifiability

For *correctness*, we checked the projects to identify the techniques to ensure this quality, such as literate programming, automated testing, symbolic execution, model checking, unit tests, etc. We also examined that whether the projects used continuous integration. As for *verifiability*, we went through the documents of the project to check the requirements specifications, theory manuals, and getting started tutorial. If a getting started tutorial existed and provided expected results, we followed it and checked if the outputs match.

The scores of *correctness & verifiability* are shown in Figure 4.3. Generally speaking, the packages with higher scores adopted more techniques to improve *correctness*, and had better documents for us to verify it.

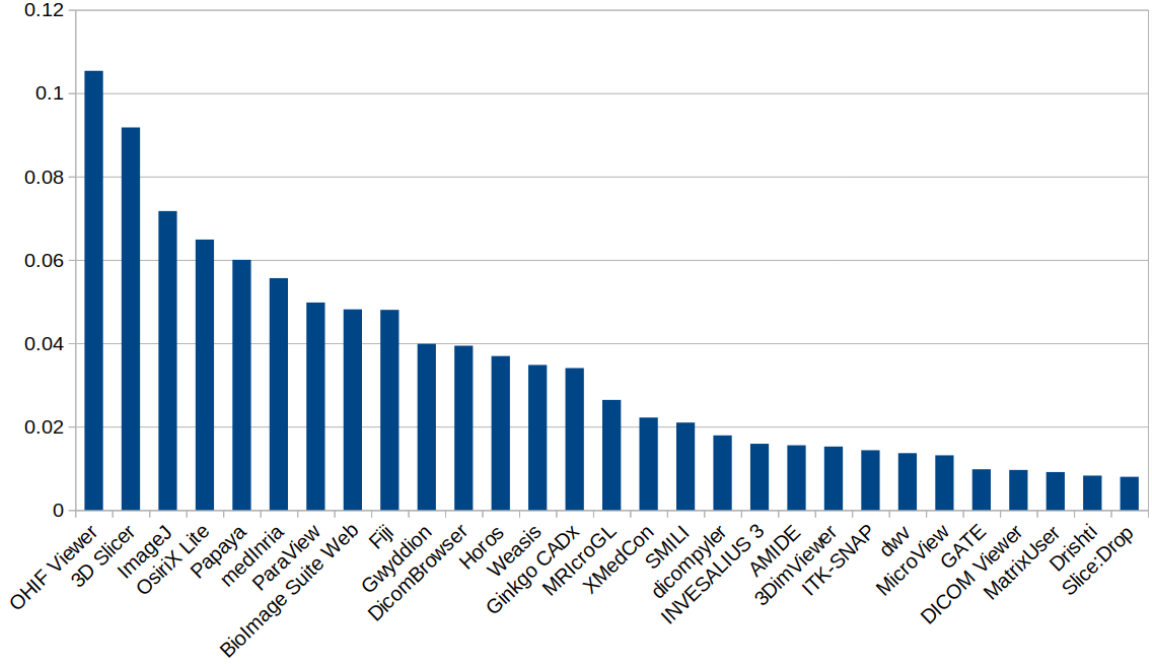


Figure 4.3: AHP correctness & verifiability scores

After examining the source code, we could not find any evidence of unit testing in more than half of the projects. Unit testing benefits most parts of the software's life cycle, such as designing, coding, debugging, and optimization [12]. It can reveal the bugs at an earlier stage of the development process, and the absence of unit testing may cause worse *correctness & verifiability*. We also could not identify the requirements specifications or theory manuals for most of the projects. It seems that even for some projects with well-organized documents, requirements specifications and theory manuals were still missing.

4.4 Surface Reliability

For *Surface Reliability*, we checked that whether the software broke during the installation and tutorial, whether there were descriptive error messages, and if we could recover the process after the errors.

As described in Section 4.2, we could not build two software packages, however, since there were online versions of them, which might be proof that successful deployments were possible. So the failure of installation did not affect their scores in *Surface Reliability*.

Figure 4.4 shows the AHP results.

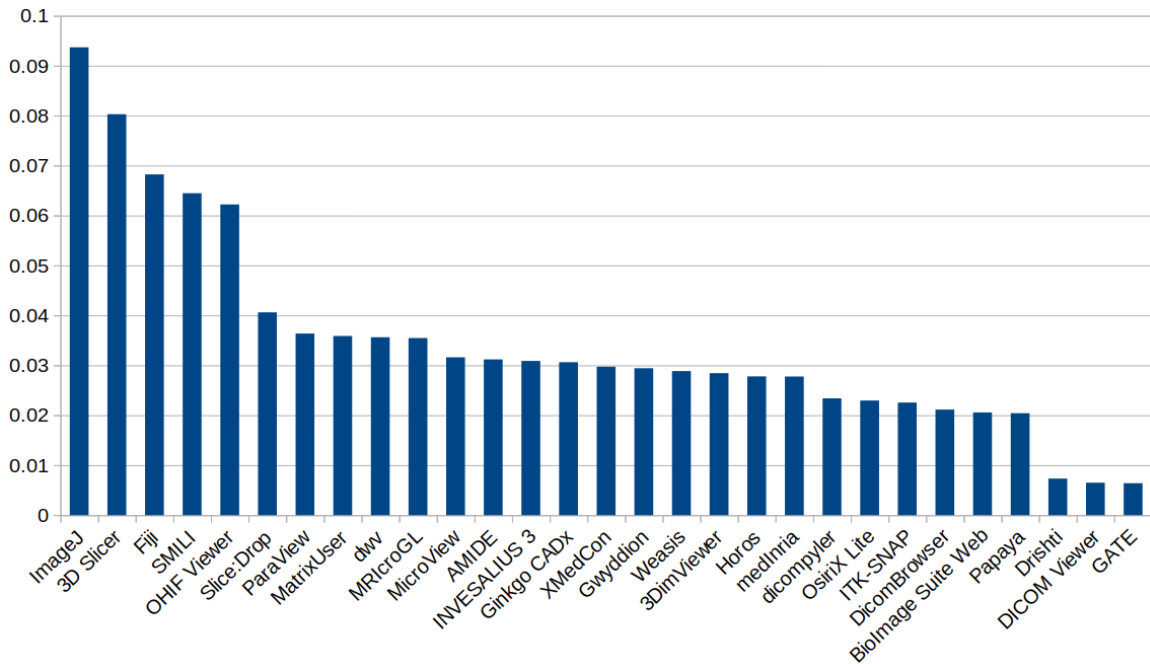


Figure 4.4: AHP surface reliability scores

For the software with the lowest scores, *Drishti* [19] and *GATE* [15] crashed during our operations, and we might underestimated the score of *DICOM Viewer* since we could not install or test it.

4.5 Surface Robustness

To test *Surface Robustness*, we checked that how the software handle unexpected/unanticipated input. Since all 29 MI software packages had functions to load image files, we prepared broken image files for the software to open. For example, a text file (.txt) with a modified extension name (.dcm) can be used as an unexpected/unanticipated input. We first tried to load a few correct input files with the software to ensure the function was working correctly, then with the unexpected/unanticipated ones. Figure 4.5 presents the scores.

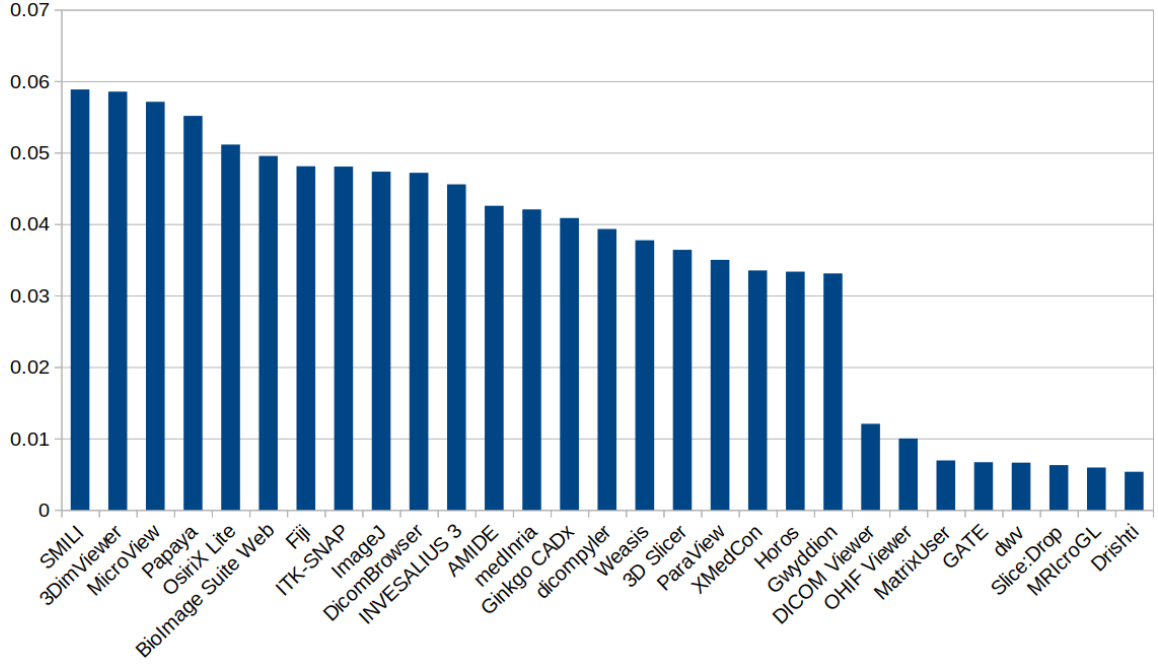


Figure 4.5: AHP surface robustness scores

The packages with higher scores elegantly handled the unexpected/unanticipated inputs, normally by showing a clear error message. Again, we could not test *DICOM Viewer*. We also might underestimate the score of *OHIF Viewer* [33] since we needed further customization to load data and the test was not complete. *MatrixUser* [20], *dww* [21], *Slice:Drop* [11], and *MRICroGL* [18] ignored the incorrect format of the input files, and displayed blank or meaningless images. *Drishti* successfully detected the unexpected/unanticipated inputs, but the software crashed as a result. For unknown reasons, *GATE* failed to load both correct and incorrect inputs.

4.6 Surface Usability

We measured *Surface Usability* by examining the documents of the projects and considered software with a getting started tutorial and a user manual easier to use. Meanwhile, we checked if users had any channels to get supports. Our impressions and user experiences when testing the software also affected the scores a lot. For example, easy-to-use graphical user interfaces gave us better experiences.

Figure 4.6 lists the AHP scores for all packages. The ones with higher scores usually provided both comprehensive document guidance and a good user experience. *INVESAL-IUS 3* [2] set an good example of detailed and clear user manual. *GATE* also provided a large number of documents, but we think that they conveyed the ideas poorly, so we had trouble understand and use them.

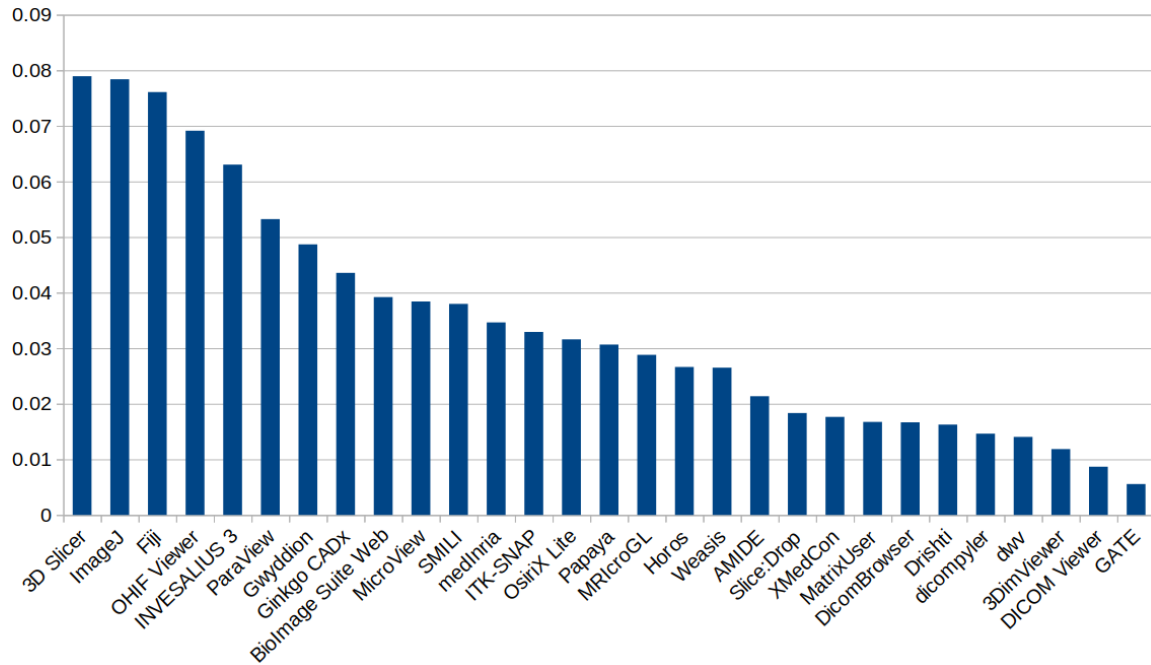


Figure 4.6: AHP surface usability scores

4.7 Maintainability

Regarding *maintainability*, we tried to search the projects' documents and identify the process of contributing and reviewing code. We believe that the artifacts of a project - including source code, documents, building scripts, etc. - can significantly affect its *maintainability*. Thus we checked each project for its artifacts, such as API documentation, bug tracker, release notes, test cases, build files, version control, etc. We also checked which tools supported issue tracking and version control and the percentages of closed issues and comment

lines in code. Figure 4.7 shows the AHP results for *maintainability*.

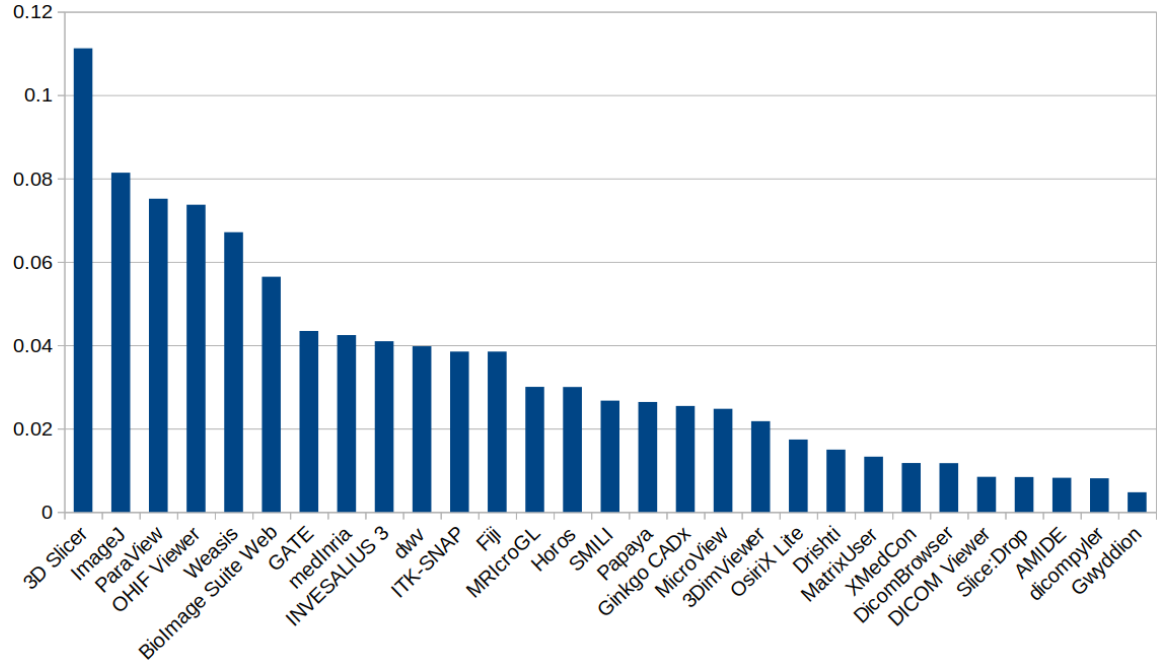


Figure 4.7: AHP maintainability scores

We marked *3D Slicer* with a much higher score than others because it did very well at closing the identified issues, and more importantly, we found it to have the most comprehensive artifacts. For example, as far as we could find out, only 5 out of the 29 projects had a project plan, 6 of them had a developer's manual, and 7 with API documentation, and *3D Slicer* included all these documents.

4.8 Reusability

To measure *reusability*, we counted the total number of code files for each project. We considered the projects with more code files to be more reusable. As mentioned in Section 3.3.2, the tools *GitStats* and *scc* count this number with minor differences, and we used the result from *GitStats* for all projects. We also decided that the projects with API documentation can deliver better *Reusability*. The AHP results are in Figure 4.8.

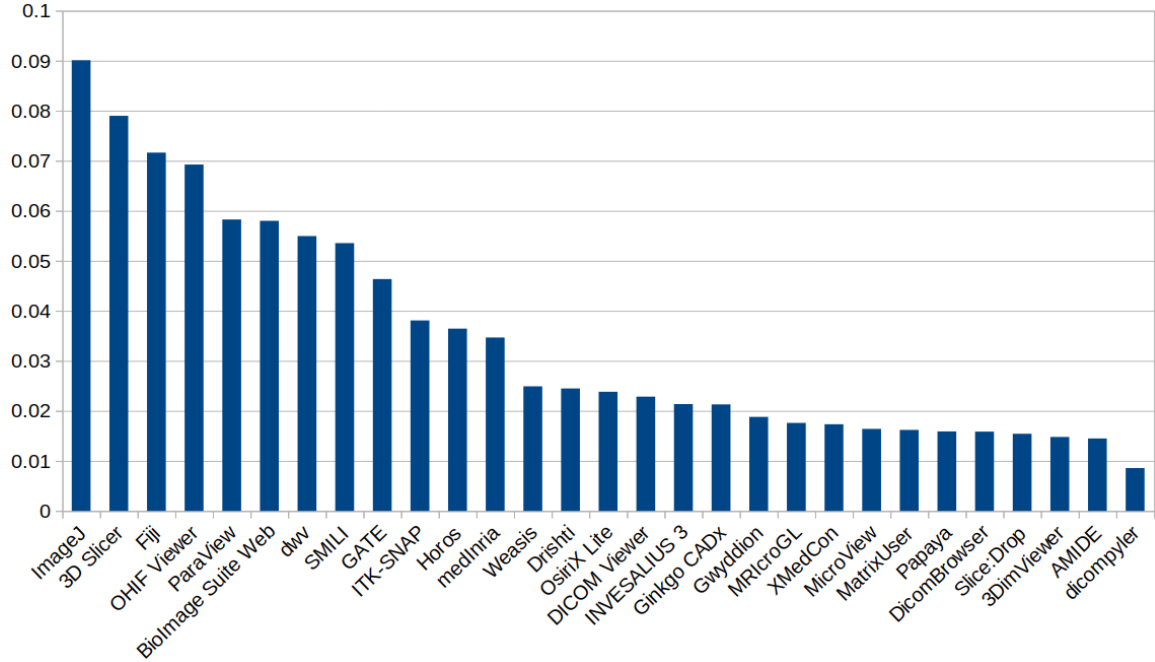


Figure 4.8: AHP reusability scores

4.9 Surface Understandability

To determine the *surface understandability* for each project, we randomly examined 10 code files. We checked the code's style within each code file, such as whether the identifiers, parameters, indentation, and formatting were consistent, whether the constants (other than 0 and 1) were hardcoded into the code if the developers modularized the code. We also checked the descriptive information for the code, such as documents mentioning the coding standard, the comments in the code, and the descriptions or links for algorithms used in the code. Figure 4.9 shows the scores.

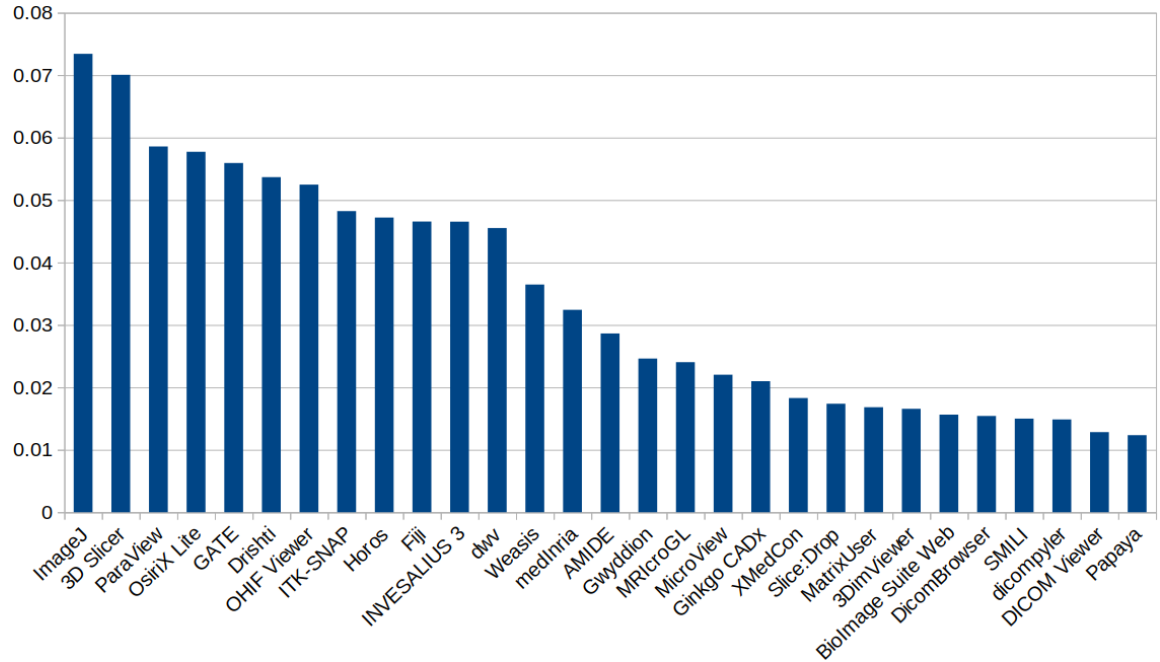


Figure 4.9: AHP surface understandability scores

Most projects had a consistent style for the code, but we only found explicit identification of a coding standard for only 3 out of the 29, which are *3D Slicer*, *Weasis* [25], and *ImageJ*.

4.10 Visibility/Transparency

To understand the *visibility/transparency*, such as all of the steps of a software development process and the current status described in Section 2.2, we tried to access this information from the existing documents. We tried to examine the development process, current status, development environment, and release notes for each project. If any information was missing or poorly conveyed, the *visibility/transparency* was not ideal.

Figure 4.10 shows the AHP scores. Generally speaking, the teams that actively documented their development process and plans scored higher because they delivered better communication to people outside the team.

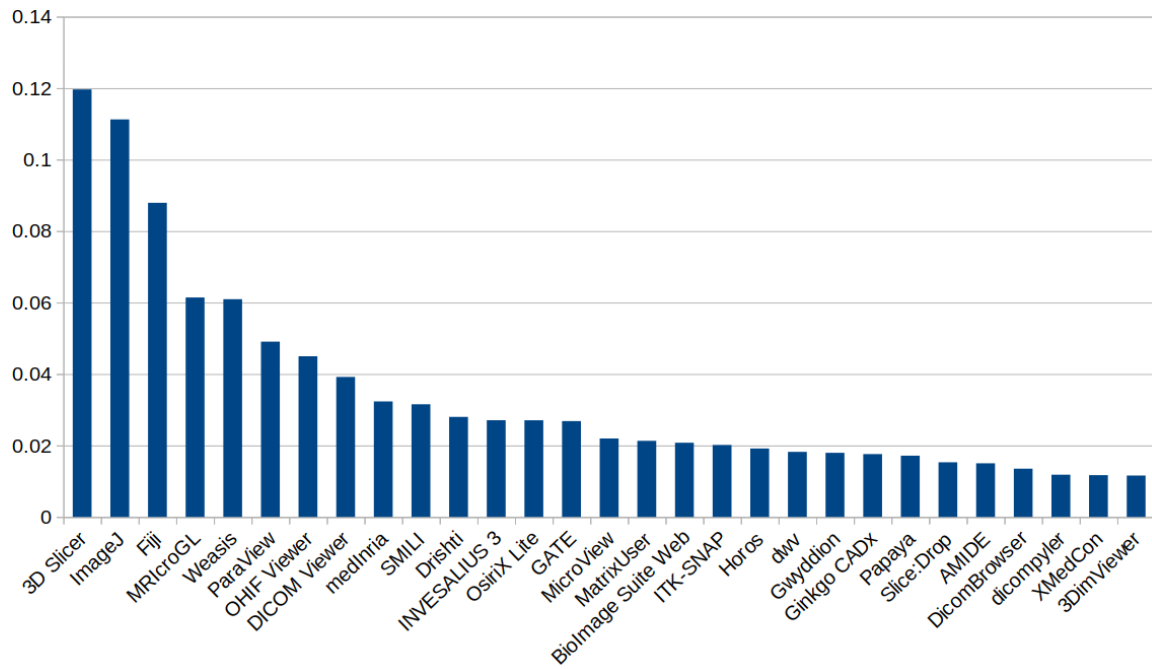


Figure 4.10: AHP visibility/transparency scores

Chapter 5

Interviews with Developers

This section summarizes some answers from the interviews with developers. We interviewed nine developers from eight of the 29 MI software projects. The eight projects are *3D Slicer*, *INVESALIUS 3*, *dvv*, *BioImage Suite Web*, *ITK-SNAP*, *MRICroGL*, *Weasis*, and *OHIF*. We spent about 90 minutes for each interview and asked 20 prepared questions. We also asked following-up questions when we thought it was worth diving deeper. One participant was too busy to have an interview, so he wrote down the answers and sent them to us.

In this section, we only discuss the answers which we think are the most exciting and essential. We summarize the rest part of the interviews in Appendix C. The interviewees may have provided multiple answers to each question. When counting the number of interviewees for each response, we count every interviewee with this answer.

5.1 Current and Past Pain Points

By asking questions 9, 10, and 12, we tried to identify the pain points during the development process in the eight projects. The pain points include current and past obstacles. We also asked the interviewees how they would solve the problems.

Questions 9, 10, and 12 are as follows,

Q9. Currently, what are the most significant obstacles in your development process?

Q10. How might you change your development process to remove or reduce these obstacles?

Q12. In the past, is there any major obstacle to your development process that has been solved? How did you solve it?

Table 5.1 shows the number of times the interviewees mentioned the current and past obstacles in their projects.

Group	Obstacle	Number of interviewees with the answer	
		as a current obstacle	as a past obstacle
1	Lack of fundings	3	
	Lack of time to devote to the project	2	1
2	Hard to keep up with changes in operating systems and libraries	1	
	Hard to support multiple operating systems	2	
	Hard to support lower-end computers	1	2
3	Lack of access to real-world datasets for testing	3	2
Others	Hard to have a high level roadmap from the start	1	
	Not enough participants for usability tests	1	
	Codebase is so big and complicated and only a few people fully understand it	1	
	Hard to transfer to new technologies		2
	Hard to understand users' needs		1
	Hard to maintain good documentations		1

Table 5.1: Current and past obstacles by the numbers of interviewees with the answers

The interviewees provided some potential and proven solutions for the problems in

Table 5.1. We group these pain points into three major groups, and put the less mentioned ones into the group *Others*. Sections 5.1.2, 5.1.2, and 5.1.3 include further discussion about the three major groups of pain points and their solutions.

5.1.1 Pain Points in Group 1

We summarize the pain points in Group 1 as **the lack of fundings and time**.

We also summarize the potential and proven solutions as follows.

Potential solutions from interviewees:

- shifting from development mode toward maintenance mode when developers are not enough;
- licensing the software to commercial companies that integrate it into their products;
- better documentation to save time for answering users' and developers' questions;
- supporting third-party plugins and extensions.

Proven solutions from interviewees:

- GitHub Actions, which is a good continuous integration and continuous delivery (CI/CD) tool to save time.

Many interviewees thought lack of fundings and lack of time were the most significant obstacles. The interviewees from *3D Slicer* team and *OHIF* team pointed out that it was

more challenging to get fundings for software maintenance as opposed to researches. The interviewee from the *ITK-SNAP* team thought enough fundings was a way to solve to problem of lacking time, because they could hire more dedicated developers. On the other hand, the interviewee from the *ITK-SNAP* team did not think that fundings could solve the same problem, since he still would need a lot of time to supervise the project.

No interviewee suggested any solution to solve the funding problem fundamentally. However, they provided ideas to save time, such as better documentation, third-party plugins, and good CI/CD tools.

5.1.2 Pain Points in Group 2

We summarize the pain points in Group 2 as **the difficulty to balance between four factors: cross-platform compatibility, convenience to development & maintenance, performance, and security**. They are also related to the choice between native application and web application.

We also summarize the potential and proven solutions as follows.

Potential solutions from interviewees:

- web applications that use computing power from computers GPU;
- to better support lower-end computers, adopting a web-based approach with backend servers;

- to better support lower-end computers, using memory-mapped files to consume less computer memory;
- more fundings;
- maintaining better documentations to ease the development & maintenance processes;

Proven solutions from interviewees:

- one interviewee saw the performance problem disappeared over the years when computers became more and more powerful.

Table 5.2 shows the teams' choices between native application and web application. In all the 29 teams on our list, most of them chose to develop native applications. For the eight teams we interviewed, three of them were building web applications, and the *MRICroGL* team was considering web-based solutions. So we had a good chance to discuss the differences between the two choices with the interviewees.

Software team	Native application	Web application
3D Slicer	X	
INVESALIUS 3	X	
dvw		X
BioImage Suite Web		X
ITK-SNAP	X	
MRICroGL	X	
Weasis	X	
OHIF		X

Total number among the eight teams	5	3
Total number among the 29 teams	24	5

Table 5.2: Teams' choices between native application and web application

The interviewees talked about the advantages and disadvantages of the two choices. We summarize the opinions from the interviewees as follows.

Native application:

- Advantages
 - higher performance
- Disadvantages

- hard to achieve cross-platform compatibility
- more complicated build process

Web application:

- Advantages
 - easy to achieve cross-platform compatibility
 - simpler build process
- Disadvantages
 - lower performance (without a backend)
 - may not meet the requirements of protecting patients' privacy (with a backend)
 - extra cost for backend servers (with a backend)

The interviewees did not conclude any perfect solutions.

5.1.3 Pain Points in Group 3

The pain point in Group 3 is **the lack of access to real-world datasets for testing.**

We also summarize the potential and proven solutions as follows.

Potential solutions from interviewees:

- using open datasets

Proven solutions from interviewees:

- asking the users to provide deidentified copies of medical images if they have problems loading the images;
- sending the beta versions of software to medical workers who can access the data and complete the tests;
- if (part of) the team belongs to a medical school or a hospital, using the datasets they can access to;
- if the team has access to MRI scanners, self-building MI datasets;
- if the team has connections with MI equipment manufacturers, asking for their help on data format problems.

No interviewee provided a perfect way to solve this problem. However, connections between the development team and medical professionals/institutions could ease the pain.

5.2 Documents in the Projects

By asking questions 11 and 19, we tried to understand the interviewees' opinions on documentation. We also aimed to find out the quality of documentations in their projects.

Questions 11 and 19 are as follows,

Q11. How does documentation fit into your development process? Would improved documentation help with the obstacles you typically face?

Q19. Do you think the current documentation can clearly convey all necessary knowledge to the users? If yes, how did you successfully achieve it? If no, what improvements are needed?

Table 5.3 summarizes interviewees' opinions on documentation. Interviewees from each of the eight projects thought that documentation was important to their projects, and most of them said that it could save their time to answer questions from users and developers. However, most of them saw the need to improve their documentation, and only three of them thought that their documentations conveyed information clearly enough.

Opinion on documentation	Number of interviewees with the answer
Documentation is vital to the project	8
Documentation of the project needs improvements	7
Referring to documentation saves time to answer questions	6
Lack of time to maintain good documentation	4
Documentation of the project conveys information clearly	3
Coding is more preferable than documentation	2
Users help each other by referring to documentation	1

Table 5.3: Opinions on documentation by the numbers of interviewees with the answers

Table 5.4 lists some of the tools and methods mentioned by the interviewees, which they used for documentation.

Tool or method for documentation	Number of interviewees with the answer
Forum discussions	3
Videos	3
GitHub	2
Mediawiki / wiki pages	2
Workshops	2
Social media	2
Writing books	1
Google Form	1
State management	1

Table 5.4: Tools and methods for documentation by the numbers of interviewees with the answers

5.3 Contribution Management and Project Management

By asking questions 5, 13, and 14, we tried to understand how the teams managed the contributions and their projects.

Questions 5, 13, and 14 are as follows,

Q5. Do you have a defined process for accepting new contributions into your team?

Q13. What is your software development model? For example, waterfall, agile, etc.

Q14. What is your project management process? Do you think improving this process can tackle the current problem? Were any project management tools used?

Maybe some team had a documented process for accepting new contributions, but none talked about it during the interview. However, most of them mentioned using GitHub and pull requests to manage contributions from the community. The interviewees generally gave very positive feedback on using GitHub. Some also said they had handled the project repository with some other tools, and eventually transferred to git and GitHub. Table 5.5 shows the number of times the interviewees mentioned the methods of receiving contributions.

Method of receiving contributions	Number of interviewees with the answer	
	as a current method	as a past method
GitHub with pull requests	8	
Code contributions from emails		3
Code contributions from forums		1
Sharing the git repository by email		1

Table 5.5: Methods of receiving contributions by the numbers of interviewees with the answers

Additionally, the *3D Slicer* team encouraged users to develop their extensions for specific use cases, and the *OHIF* team was trying to enable the use of plug-ins; the interviewee from the *ITK-SNAP* team said one way of accepting new team members was through funded academic projects.

Table 5.6 shows the software development models by the numbers of interviewees with the answers. Only two interviewees confirmed their development models. The others did not think they used a specific model, but three of them suggested that their processes were similar to Waterfall or Agile.

Software development model	Number of interviewees with the answer
Undefined/self-directed	3
Similar to Agile	2
Similar to Waterfall	1
Agile	1
Waterfall	1

Table 5.6: Software development models by the numbers of interviewees with the answers

Some interviewees mentioned the project management tools they used, which are in Table 5.7. Generally speaking, they talked about two types:

- Trackers, including GitHub, issue trackers, bug trackers and Jira;

- Documents, including GitHub, Wiki page, Google Doc, and Confluence.

Project management tools	Number of interviewees with the answer
GitHub	3
Issue trackers	1
Bug trackers	1
Jira	1
Wiki page	1
Google Doc	1
Confluence	1

Table 5.7: Project management tools by the numbers of interviewees with the answers

No interviewee introduced any strictly defined project management process. The most common way was following the issues, such as bugs and feature requests. Additionally, the *3D Slicer* team had weekly meetings to discuss the goals for the project; the *INVESALIUS* 3 team relied on the GitHub process for their project management; the *ITK-SNAP* team had a fixed six-month release pace; only the interviewee from the *OHIF* team mentioned that the team has a project manager; the *3D Slicer* team and *BioImage Suite Web* team were doing nightly builds and tests.

Most interviewees skipped the question “Do you think improving this process can

tackle the current problem?”, but the interviewee from the *OHIF* team gave a positive answer.

Chapter 6

Threat to Validity

- limited time to each software, surface quality
- impressions, subjective
- searched all information by ourselves, so maybe missed some info
- mainly measured by one person, so maybe biased

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] Aysel Afsar. Dicom viewer. <https://github.com/ayselafsar/dicomviewer>, 2021. [Online; accessed 27-May-2021].
- [2] Paulo Amorim, Thiago Franco de Moraes, Helio Pedrini, and Jorge Silva. Invesalius: An interactive rendering framework for health care support. page 10, 12 2015.
- [3] 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.
- [4] Ben Boyter. Sloc cloc and code. <https://github.com/boyter/scc>, 2021. [Online; accessed 27-May-2021].
- [5] 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.

- [6] James Edward Corbly. The free software alternative: Freeware, open source software, and libraries. *Information Technology and Libraries*, 33(3):65–75, Sep. 2014.
- [7] Steve Emms. 16 best free linux medical imaging software. <https://www.linuxlinks.com/medicalimaging/>, 2019. [Online; accessed 02-February-2020].
- [8] Tomasz Gieniusz. Gitstats. https://github.com/tomgi/git_stats, 2019. [Online; accessed 27-May-2021].
- [9] GNU. Categories of free and nonfree software. <https://www.gnu.org/philosophy/categories.html>, 2019. [Online; accessed 20-May-2021].
- [10] 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.
- [11] Daniel Haehn. Slice:drop: collaborative medical imaging in the browser. pages 1–1, 07 2013.
- [12] Paul Hamill. *Unit test frameworks: Tools for high-quality software development*. O'Reilly Media, 2004.

- [13] 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].
- [14] 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.
- [15] Sama Jan, Giovanni Santin, Daniel Strul, S Staelens, K Assié, Damien Autret, Stéphane Avner, Remi Barbier, Manuel Bardiès, Peter Bloomfield, David Brasse, Vincent Breton, Peter Bruyndonckx, Irene Buvat, AF Chatziioannou, Yunsung Choi, YH Chung, Claude Comtat, Denise Donnarieix, and Christian Morel. Gate: a simulation toolkit for pet and spect. *Physics in medicine and biology*, 49:4543–61, 11 2004.
- [16] Alark Joshi, Dustin Scheinost, Hirohito Okuda, Dominique Belhachemi, Isabella Murphy, Lawrence Staib, and Xenophon Papademetris. Unified framework for development, deployment and robust testing of neuroimaging algorithms. *Neuroinformatics*, 9:69–84, 03 2011.
- [17] Ron Kikinis, Steve Pieper, and Kirby Vosburgh. *3D Slicer: A Platform for Subject-Specific Image Analysis, Visualization, and Clinical Support*, volume 3, pages 277–289. 01 2014.

- [18] Chris Rorden's Lab. Mricrogl. <https://github.com/rordenlab/MRIcroGL>, 2021. [Online; accessed 27-May-2021].
- [19] Ajay Limaye. Drishti, a volume exploration and presentation tool. volume 8506, page 85060X, 10 2012.
- [20] Fang Liu, Julia Velikina, Walter Block, Richard Kijowski, and Alexey Samsonov. Fast realistic mri simulations based on generalized multi-pool exchange tissue model. *IEEE Transactions on Medical Imaging*, PP:1–1, 10 2016.
- [21] Yves Martelli. dwv. <https://github.com/ivmartel/dwv>, 2021. [Online; accessed 27-May-2021].
- [22] 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.
- [23] Hamza Mu. 20 free & open source dicom viewers for windows. <https://medevel.com/free-dicom-viewers-for-windows/>, 2019. [Online; accessed 31-January-2020].
- [24] The Linux Information Project. Freeware definition. <http://www.linfo.org/freeware.html>, 2006. [Online; accessed 20-May-2021].

- [25] Nicolas Roduit. Weasis. <https://github.com/nroduit/nroduit.github.io>, 2021. [Online; accessed 27-May-2021].
- [26] Curtis Rueden, Johannes Schindelin, Mark Hiner, Barry DeZonia, Alison Walter, and Kevin Eliceiri. ImageJ2: ImageJ for the next generation of scientific image data. *BMC Bioinformatics*, 18, 11 2017.
- [27] 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.
- [28] 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].
- [29] Spencer Smith, Jacques Carette, Olu Owolaiye, Peter Michalski, and Ao Dong. Quality definitions of qualities. Manuscript in preparation, 2020.
- [30] Spencer Smith, Jacques Carette, Olu Owolaiye, Peter Michalski, and Ao Dong. Methodology for assessing the state of the practice for domain x. Manuscript in preparation, 2021.

- [31] Omkarprasad S. Vaidya and Sushil Kumar. Analytic hierarchy process: An overview of applications. *European Journal of Operational Research*, 169(1):1–29, 2006.
- [32] Greg Wilson, Dhavide Aruliah, C. Titus Brown, Neil Chue Hong, Matt Davis, Richard Guy, Steven Haddock, Kathryn Huff, Ian Mitchell, Mark Plumbley, Ben Waugh, Ethan White, and Paul Wilson. Best practices for scientific computing. *PLoS Biology*, 12:e1001745, 01 2014.
- [33] Erik Ziegler, Trinity Urban, Danny Brown, James Petts, Steve D. Pieper, Rob Lewis, Chris Hafey, and Gordon J. Harris. Open health imaging foundation viewer: An extensible open-source framework for building web-based imaging applications to support cancer research. *JCO Clinical Cancer Informatics*, (4):336–345, 2020. PMID: 32324447.

Appendix A

Full Grading Template

appendix here

Appendix B

Summary of Measurements

appendix here

Appendix C

Other Interview Answers

We asked 20 interview questions to the nine interviewees from eight software projects. We discuss the answers to interview questions 5, 9, 10, 11, 12, 13, 14, and 19 in Section 5, and summarize the answers to the other questions in this section.

Q1. Interviewees' current position/title? degrees?

Six of the nine interviewees revealed their position/title, such as CEO of a company, endowed chair and professor in universities, software engineers in a commercial company and a hospital.

Most of them answered their backgrounds and degrees. Table C.1 shows the highest academic degrees the participants have, and Table C.2 shows what majors they studied. Many of the interviewees studied in multiple majors.

Highest degree	Number of interviewees
PHD	4
Master	3
Bachelor	0
Unspecified academic degree	2

Table C.1: Interviewees' highest academic degrees

Major	Number of interviewees
Computer Science	4
Physics	2
Biomedical Engineering	1
Neuroimaging	1
Geology (image analysis)	1
Media Arts and Sciences	1
Mechanical Engineering	1
Materials Engineering	1
Psychology	1

Table C.2: Interviewees' majors at university

Q2. Interviewees' contribution to/relationship with the software?

Table C.3 shows the interviewees' roles and responsibilities in the projects. One of the participants did not explicitly mention his role, but implicitly revealed that he was a primary contributor to the project.

Role in the projects	Number of interviewees
Chief Architect	2
Lead Developer	1
Core Developer	5
Unspecified	1

Table C.3: Interviewees' roles in the projects

Q3. Length of time the interviewee has been involved with this software?

Table C.4 shows the distribution of the lengths of time the interviewees had worked on the projects.

Length of time in the projects	Number of interviewees
0-1	1
2-5	0
6-10	2
11-15	3
16-20	2
21-25	1

Table C.4: Lengths of time that the interviewees worked in the projects

Q4. How large is the development group?

The size of each group grows and shrinks over the years. Most teams mentioned that the team members join and leave. Some teams said that when there was sufficient funding, they could afford more developers.

Table C.5 shows the numbers of active members at the time of interviews. The members include people working on development and project management.

Number of current members	Number of projects
1-3	5
4-6	3

Table C.5: Numbers of current members in the projects

As shown in the table, no team had a vast number of members. Some projects had more developers, such as *3D Slicer*; on the other hand, some teams such as *dwy* had only one primary developer, plus a maximum of two or three developers occasionally.

3D Slicer is a special case, because it supports third-party extensions. So there have been community members developing and maintaining these extensions. Table C.5 does not include these members.

Q6. What is the typical background of a developer?

Not all interviewees could clearly answer this question. Many of them talked about the backgrounds of members with who they were familiar. Table C.6 shows the number of times all interviewees mentioned a background.

Background of a developer	Number of interviewees with the answer
Computer Science, Information Technology, and Software Development	6
Imaging	2
Medical Imaging	1
Mathematics	1
Biomedical Engineering	1
Computer Aided Medical Procedures	1
Physician	1

Table C.6: Backgrounds of developers by the numbers of interviewees with the answers

Q7. What is your estimated number of users? How did you come up with that estimate?

None of the interviewees knew the exact number of users. Some of them provided estimations based on different facts. However, we do not think these numbers are comparable to each other.

Software	Rough estimation	Considered facts
3D Slicer	100,000	Search results on Google Scholar; number of new posts per year on slicer.org; number of downloads.
INVESALIUS 3	75,000	Number of random IDs created by new installation.
dwv	No estimation	About 20 companies integrated <i>dwv</i> in their products.
BioImage Suite Web	100 active users	The interviewee only counted the users from several Universities who were active users.
ITK-SNAP	10,000 plus	Number of downloads.
MRICroGL	No estimation	It is the top 1 downloaded software on this NITRC list https://www.nitrc.org/top/toplist.php?type=downloads
Weasis	10,000 user used it as least once	Number of profiles.
OHIF	About 5000	Some platforms integrated <i>OHIF</i> , and it was hard to know the number of end users.

Table C.7: Rough Estimations for the Number of Users

Table C.7 shows the estimations and how the interviewees made them. It is clear that some estimated only the active users, and some counted users who had used only once. So

we do not compare these numbers.

Q8. What is the typical background of a user?

All interviewees provided several different user backgrounds, and all of them mentioned medical researchers or medical professionals. Table C.8 shows the number of times all interviewees mentioned a background.

Background of a user	Number of interviewees with the answer
Medical Researchers	6
Doctors/Health care professionals/Surgeons	5
Student Researchers	4
Patients	3
Paleontologist	1
Biomechanical Engineers	1
Imaging Researchers	1
Mechanical Engineers	1

Table C.8: Backgrounds of users by the numbers of interviewees with the answers

Q15. Was it hard to ensure the correctness of the software? If there were any obstacles, what methods have been considered or practiced to improve the situation? If practiced, did it work?

Table C.9 shows the threats to *correctness* by the numbers of interviewees with the

answers.

Threat to correctness	Number of interviewees with the answer
Clinical systems produce data in various formats (e.g. DICOM), which have complicated standards.	2
The software needs to handle the complexity.	
Different types of MI machines can create data in slightly different ways, adding more complexity for the software to handle.	2
Besides viewing, the software has additional functions, which lead to extra complexity.	1
There is lack of real word image data for testing.	1
The team cannot use private data for debugging, even when the data cause problems.	1
The software uses huge datasets for testing, so the tests are expensive and time-consuming.	1
It is hard to well manage releases.	1
The project has no unit tests.	1
The project has no dedicated quality assurance team.	1

Table C.9: Threats to correctness by the numbers of interviewees with the answers

Table C.10 shows the strategies to ensure *correctness* by the numbers of interviewees with the answers. The interviewees from the *3D Slicer* and *ITK-SNAP* teams thought that the self-tests and automated tests were beneficial and could significantly save time. The interviewee from the *Weasis* team kept collecting medical images for more than ten years. These images have caused problems with the software. So he had samples to test specific problems.

Strategy to ensure correctness	Number of interviewees with the answer
Test-driven development / component tests / integration tests / smoke tests / regression tests.	4
Self tests / automated tests.	3
Two stage development process / stable release & nightly builds.	3
CI/CD.	1
Using deidentified copies of medical images for debugging.	1
Sending the beta versions of software to medical workers who can access the data and do the tests.	1
Collecting and maintaining a dataset of problematic images.	1

Table C.10: Strategies to ensure correctness by the numbers of interviewees with the answers

Q16. When designing the software, did you consider the ease of future changes? For example, will it be hard to change the system's structure, modules, or code blocks? What measures have been taken to ensure the ease of future changes and maintains?

Table C.11 shows the strategies to ensure *maintainability* by the numbers of interviewees with the answers. The modular approach is the most talked-about solution to improve *maintainability*. The *3D Slicer* team used a well-defined structure for the software, which they named as “event-driven MVC pattern”. Moreover, *3D Slicer* discovers and loads necessary modules at runtime, according to the configuration and installed extensions. The *BioImage Suite Web* team had designed and re-designed their software multiple times in the last 10+ years. They found that their modular approach effectively supported the maintainability [16].

Strategy to ensure maintainability	Number of interviewees with the answer
Modular approach / thinking the software as reusable Lego bricks / maintain repetitive functions as libraries.	5
Supporting third-party extensions.	1
Easy-to-understand architecture.	1
Dedicated architect.	1
Starting from simple solutions.	1
Documentation.	1

Table C.11: Strategies to ensure maintainability by the numbers of interviewees with the answers

Q17. Provide instances where users have misunderstood the software. What, if any, actions were taken to address understandability issues?

Table C.12 shows the threats to *understandability* by the numbers of interviewees with the answers. It separates *understandability* issues to users and developers by the horizontal dash line.

Threat to understandability	Number of interviewees with the answer
Not all users understand how to use some features of the software.	2
The team has no dedicated user experience (UX) designer.	1
The software does not make some important indicators noticeable (e.g. a progress bar).	1
Not all users understand the purpose of the software.	1
Not all users know if the software includes certain features.	1
Not all users understand how to use the command line tool.	1
Not all users understand that the software is a web application.	1
<hr/>	
Not all developers understand how to deploy the software.	1
The architecture is difficult for new developers to understand.	1

Table C.12: Threats to understandability by the numbers of interviewees with the answers

Table C.13 shows the strategies to ensure *understandability* by the numbers of interviewees with the answers.

Strategy to ensure understandability	Number of interviewees with the answer
Documentation / user manual / user mailing list / forum.	4
Graphical user interface.	2
Testing every release with active users.	1
Making simple things simple and complicated things possible.	1
Icons with more clear visual expressions.	1
Designing the software to be intuitive.	1
Having a UX designer with the right experience.	1
Dialog windows for important notifications.	1
Providing an example if the users need to build the software by themselves.	1

Table C.13: Strategies to ensure understandability by the numbers of interviewees with the answers

Q18. What, if any, actions were taken to address usability issues?

Table C.14 shows the strategies to ensure *usability* by the numbers of interviewees with the answers.

Strategy to ensure usability	Number of interviewees with the answer
Usability tests and interviews with end users.	3
Adjusting according to users' feedbacks.	3
Straightforward and intuitively designed interface / professional UX designer.	2
Providing step-by-step processes, and showing the step numbers.	1
Making the basic functions easy to use without reading the documentation.	1
Focusing on limited number of functions.	1
Making the software more streamlined.	1
Downsampling images to consume less memory.	1
An option to load only part of the data to boost performance.	1

Table C.14: Strategies to ensure usability by the numbers of interviewees with the answers

Q20. Do you have any concern that your computational results won't be reproducible

in the future? Have you taken any steps to ensure reproducibility?

Table C.15 shows the threats to *reproducibility* by the numbers of interviewees with the answers.

Threat to reproducibility	Number of interviewees with the answer
If the software is closed-source, the reproducibility is hard to achieve.	1
The project has no user interaction tests.	1
The project has no unit tests.	1
Using different versions of some common libraries may cause problems.	1
CPU variability can leads to non-reproducibility.	1
When reverse-engineering how manufacturers create medical images, the team may misinterpret it.	1

Table C.15: Threats to reproducibility by the numbers of interviewees with the answers

Table C.16 shows the strategies to ensure *reproducibility* by the numbers of interviewees with the answers. The interviewee from the *3D Slicer* team provided various suggestions. One interviewee from another team suggested that they used *3D Slicer* as the benchmark to test their *reproducibility*.

Strategy to ensure reproducibility	Number of interviewees with the answer
Regression tests / unit tests / having good tests.	6
Making code, data, and documentation available / making the software open-source / using open-source libraries.	5
Running same tests on all platforms.	1
A dockerized version of the software, insulating it from the operating system environment.	1
Using standard libraries.	1
Monitoring the upgrades of the libraries.	1
Clearly documenting the versions.	1
Bringing along the exact versions of all the dependencies with the software.	1
Providing checksums of the data.	1
Benchmark the software against other software with similar purposes.	1

Table C.16: Strategies to ensure reproducibility by the numbers of interviewees with the answers

Appendix D

Ethics Approval

appendix here