

LEWIS UNIVERSITY

Structural Quality & Software Evolution

A Thesis

By

Alison Major

Department of Computer and Mathematical Sciences

Submitted in partial fulfillment of the requirements

for the degree of
Master of Science in Computer Science,

Concentration in Software Engineering

February 13, 2022

The undersigned have examined the thesis entitled ‘**Structural Quality & Software Evolution**’ presented by **Alison Major**, a candidate for the degree of **Master of Science in Computer Science (Concentration in Software Engineering)** and hereby certify that it is worthy of acceptance.

TODO: Add lines for signatures here. Justify above paragraph.

Abstract

Some software engineering projects fail to evolve, which makes them obsolete. This topic is interesting and important to developers because the software that fails to evolve will fail to generate user engagement, leading to revenue loss. We review a number of projects and resources to understand the correlation of software structure quality and its impacts on a system's ability to evolve. With this understanding, we explore ways to improve the evolution of a software system through tools and suggestions.

TODO: Update the abstract with each iteration of this paper.

Acknowledgements

TODO: Update this section with my own acknowledgements.

Gratitude is a great virtue, though revenge is profitable

It's customary and good manners to say thank you however, where do you draw the line? In some of the theses that I've read, and I write this after having read thousands, literally, the following and more have been acknowledged: God, one's advisor, one's better half, parents, children, friends, classmates, lab-mates, lab technicians, lab assistants, pets, fav. Prof, neighbors, physicians, exercise trainer(s), wiki, the maintenance guy, landlord, the school hockey team, secretary, department head, driver, dentist, chauffeur, the police, fav. presidential candidate, one's chef, Led Zeppelin, the pastor, one's biggest crush, the cable man, the mani/pedi girl, hair stylist, the best/worst/fav bar tender(s), the janitor, one's obs/gyn, one's mentor, and in a more recent thesis, Michael Phelps (8 gold medals at the 2008 Olympic games in Beijing, China, way to go...)

Keep in mind that one has to use one's own words when writing an acknowledgement. Plagiarism is unauthorized.

Contents

| | |
|--|----------|
| Abstract | ii |
| Acknowledgements | iii |
| List of Tables | vi |
| List of Figures | vii |
| 1 Introduction | 1 |
| 1.1 Maintainability Index and Pylint Refactor Scores | 3 |
| 1.2 Paper Structure | 4 |
| 2 Background and Literature Review | 5 |
| 2.1 Keeping Users Engaged Long Term | 6 |
| 2.1.1 Why does software evolution matter? | 6 |
| 2.1.2 How do we ensure software evolution? | 7 |
| 2.2 The Impact of Structural Quality | 8 |
| 2.2.1 Software Maintenance | 8 |
| 2.2.2 Software Evolution | 10 |
| 2.2.3 Measuring Maintainability | 14 |

| | | |
|----------|---|-----------|
| 2.2.4 | Maintainability Scores | 15 |
| 2.2.5 | Other Maintainability Characteristics | 18 |
| 2.2.6 | Documentation and Maintainability | 20 |
| 2.3 | Related Work | 23 |
| 2.3.1 | Considering Data Sets | 23 |
| 2.3.2 | Design Patterns and Software Quality | 24 |
| 2.3.3 | Software Architecture and Maintainability | 25 |
| 3 | Methodology | 27 |
| 3.1 | Initial Repository Set | 28 |
| 3.2 | Filtered Repository Set | 29 |
| 4 | Results | 31 |
| 5 | Conclusions and Recommendations | 32 |
| | References | 40 |

List of Tables

| | | |
|-----|---|----|
| 3.1 | Criteria used to filter down the initial set of repositories. . . . | 30 |
|-----|---|----|

List of Figures

| | | |
|-----|---|----|
| 2.1 | A simplified visual of Lehman’s fifth law, “Conservation of Familiarity.” | 13 |
| 5.1 | A snapshot of the badges from <i>SymPy</i> ’s repository. | 33 |

Chapter 1

Introduction

When building software systems, we have several areas of concern: cost, delivery timeline, quality, etc. The cost and time-to-market are often the two problems given the highest priority in a project. However, engineers must consider the software quality to preserve the system's longevity. Despite its importance, the code and architecture quality can be challenging to understand and measure.

When we think about projects, we can assume that as time goes on and changes and additions occur within a system's source code, the complexity of that system will grow. However, when we manage the code structure, we can keep the complexity in check, allowing systems to evolve. Developers can maintain this structure through simple steps like having readable code and more complex considerations, like how coupled and cohesive a system is.

One way to understand the quality around a system is to discuss its

“maintainability,” the ease of receiving new features or resolving bugs. For example, developers may find that adjusting one area to add a new feature requires touching several other code areas in tightly coupled systems. Some code measuring systems provide a Maintainability Index (MI), a well-known quality measure. However, its effectiveness in quantifying software quality is debated [1].

On the other hand, code smells are used extensively by practitioners to identify low-quality spots in the software system. These areas would need the teams’ attention and are good candidates for refactoring.

1.1 Maintainability Index and Pylint Refactor Scores

Pylint is a static analysis tool that identifies several classes of code quality concerns. Particularly relevant to our study are refactor violations, which report on various code smells. We can assume that there must be some correlation between Maintainability Index and the type and number of code smells in a software system, quantified by the Pylint refactor score.

This study explores such assumptions and systematically investigates any correlation between the Maintainability Index metric and the Pylint Refactor score. Furthermore, we perform analysis on specific refactor violations to reveal and shed light on the relative effectiveness of the different refactor violations and their relationship to Maintainability Index.

The structural quality of a software system will impact the software evolution. If the project has poor structural quality, the architecture will minimize its ability to evolve, and the software system will eventually “die-off” so to speak.

We will look at many open-source Python systems using Pylint and attempt to correlate the data from the Pylint scores to the level of ease in adding new features to the system. This will determine if a system is more maintainable with better Pylint scores.

1.2 Paper Structure

In Chapter 2, we will dig into a deeper background of the topic, exploring ideas of why software systems need to keep users engaged long term (Section 2.1). We will explore automated measurements that provide evaluation scores of software systems. By using some of these quality and maintainability scores, we can see how structure impacts evolution (Section 2.2). Additionally, we'll explain how maintainability is measured, as well as the different attributes that can factor into maintainability. We will also review related works (Section 2.3).

With more background on the problem, we can then review the methodology for our research in Chapter 3. Here we will review where we found our initial data set (Section 3.1) and what criteria we used to filter it to a manageable size for our tests (Section 3.2).

Chapter 4 will review the results of our research, using the methodology previously explained.

We will then provide final conclusions and recommendations in Chapter 5.

Chapter 2

Background and Literature Review

In this chapter we will gain a better understanding of the topic at hand, first by understanding long term user engagement. From there we will explore the impact of structural quality within a software system. Once we understand the problem thoroughly, we will determine a useful dataset that we can apply our theories to, as well as explore related works.

2.1 Keeping Users Engaged Long Term

When developing a new system or a new software idea, getting the project off the ground and in front of users is one thing. However, keeping that project alive with a thriving community of engaged users is another.

The systems we create could be customer-facing web applications, games, or internal applications used to carry out tasks. Regardless of the type of system, the product will no longer provide usefulness without evolving with the user's needs. Even in a corporate setting with internal business systems, over time, users will need change; how a system can adapt to those needs requires a level of flexibility.

“Software evolution is the continual development of software after its initial release to address changing stakeholder and/or market requirements.” [2]

2.1.1 Why does software evolution matter?

When a system cannot evolve, the impact is primarily felt by the users. However, this impact will eventually get back to those who created and continue to support the system. With users that are either unsatisfied or unable to use the system any longer, the engagement levels will drop. The decline in users will ultimately result in a loss of income, as the system can no longer deliver to the needs of its audience.

Because organizations invest large amounts of money in the software systems that they create, they depend on the software’s continued success. Software evolution will allow the system to adapt to new or changing business requirements, fix bugs and defects, and integrate with other systems that have changed and evolved that may share the same software environment.

As a system is used, inevitably, users will stumble into situations that even the best quality assurance testers will miss. When defects are found, they will require fixing.

To keep a system up-to-date, we must add new features. For example, there may be a need to improve a system’s performance or reliability, especially if the user base expands.

Security can also impact the need for a system to be maintained. New ways to infiltrate a system can be uncovered, so it is important to stay on top of newest versions of dependencies and technologies in order to avoid potential breaches of data and experience.

2.1.2 How do we ensure software evolution?

Because the maintainability of a system can ultimately influence the ability to generate revenue, we must find ways to ensure that a project will evolve. One of these ways could be to ensure that a project continues to be considered “maintainable” throughout its lifetime. This system characteristic will ensure that bugs can be fixed quickly, but new features should be easy to add as the users’ needs evolve.

2.2 The Impact of Structural Quality

2.2.1 Software Maintenance

The structural quality of a software system will impact the software evolution. If the project has poor structural quality, its ability to evolve will be minimized, and the software system will eventually “die-off” so to speak.

There is much planning involved in all software creation projects in what the product will be, will do, who it is for, etc. One of the things that should also be on the planning list is long-term maintenance and growth. That is, how do we build a thing that will be easier to add features to down the road?

Let us define maintainability in the context of software. For example, a system would be considered easy to maintain if it is easy to debug and easy to add new features. These new features are generally considered minor features, and may often be reported as bugs by users, when in reality, they are looking for functionality enhancements [3].

“Software maintenance in software engineering is the modification of a software product after delivery to correct faults, to improve performance or other attributes.” [3]

It may be easier to understand what characteristics define a system with poor maintainability. These types of systems will have poor code quality, leading to defects. For example, there could be undetected vulnerabilities or vulnerabilities that have been ignored. It may be that the system is overly

complex. In addition to the complexity, it could be hard to read due to poor naming or dead (unused) code throughout the source code.

A project is known to have good maintainability when there is an enforced set of clean and consistent standards for the code. This often involves having human-readable names for functions, methods, and variables. Any complex code is minimized, and methods are small and focus on a single thing. Parts of the system are decoupled and organized, making it easy to work on different parts with low impact on unrelated parts. For example, the code is DRY (there is limited redundancy in the code), unused code has been removed, and there is a level of documentation that supports an easy understanding of the system.

Why should we care about whether the code is maintainable? It is assumed that a large amount of the cost over the lifetime of a project is attributed to maintainability. Fred Brooks, in his book “The Mythical Man-Month” even claimed that over 90% of the costs for a typical software system come up in the maintenance phase [4]. Once the bulk of the system is off the ground and live worldwide, how well the team can improve the system with new features and fix bugs, even working on different parts in parallel, can be impacted by its maintainability. Any successful piece of software will inevitably need to be maintained.

2.2.2 Software Evolution

There is a distinction to be made between **software maintenance** and **software evolution**. We will refer to software maintenance as bug resolution and for minor functional improvements. For example, we can consider this routine maintenance when we must fix a broken route in the application or provide a subtle enhancement on the user experience. However, when we look at upgrades to the system, adaptations to the changing and growing needs of the user, or migrating the system to a new technology, we can refer to this as evolution of the software.

The evolution of software can result from new laws that have come into being. As technology itself changes, governing bodies must continually revisit data collection and information sharing policies. Changes in technology and laws may lead to adaptations in the software systems.

It is also fair to say that systems will change because we can never fully determine a user's needs at the start of a project. It would be safe to say that the user's needs will change over time themselves. This leads to a never-ending project that will always need some form of enhancement.

Meir “Manny” Lehman and László “Les” Bélády contributed to a list of laws involving software evolution known as Lehman's Laws that describe a balance between forces that drive new developments while also slowing progress. These laws apply to programs that were written to perform some real-world activity, where its behavior is linked to the environment in which it runs; additionally, this program category assumes that the program needs

to adapt to varying requirements and circumstances in that environment. Eight laws were created and are listed below. [5]

1. **Continuing Change** *(1974)*
2. **Increasing Complexity** *(1974)*
3. **Self Regulation** *(1974)*
4. **Conservation of Organisational Stability** *(1978)*
5. **Conservation of Familiarity** *(1978)*
6. **Continuing Growth** *(1991)*
7. **Declining Quality** *(1996)*
8. **Feedback System** *(1996)*

The first law, “Continuing Change,” tells us that if a system does not adapt, it will become progressively less satisfactory. The second, “Increasing Complexity,” explains that as a system evolves, unless work is done to maintain or reduce complexity, the complexity will increase. This can be due to the added volume of the code from new features or even an increasing number of developers that have edited the code. Unless this phenomenon of increased complexity is actively addressed during changes, it can impact the maintainability (and the ability of a project to continue evolving) in the future.

Lehman’s fifth law, “Conservation of Familiarity,” explains how the average incremental growth does not change over time as a system evolves. The people interacting with the system, such as the developers, business persons, or users, must still continue using and working within the system at the same “level of mastery.” If the system grows and changes excessively, the mastery will drop, slowing down the next set of changes. This could be because the source code or architecture has become more complex (impacting the developers’ ability to adapt and enhance the system) or because the user features have changed so that the system audience needs time to master the new interfaces or new tools. Because of this natural “slow-down” for excessive change, the average incremental growth will remain steady. We can see a simplified visual in “Fig. 2.1” showing that when the number of changes spikes (that is to say, when there is excessive growth in a system), it will be followed by an iteration of fewer changes, leading to a nearly consistent average of incremental growth (the thick, horizontal line) over time.

In Lehman’s sixth law, “Continuing Growth,” we see that the system user’s satisfaction will not be maintained without continually increasing the functional content. Along a similar idea, the law pertaining to “Declining Quality” states that if the operational environment for the system does not change, the system’s quality will appear to decline. Therefore, we must continue adapting for even the appearance of the maintained quality of a system.

With all of these characteristics surrounding the evolution of software, we

Number of Changes vs. Time

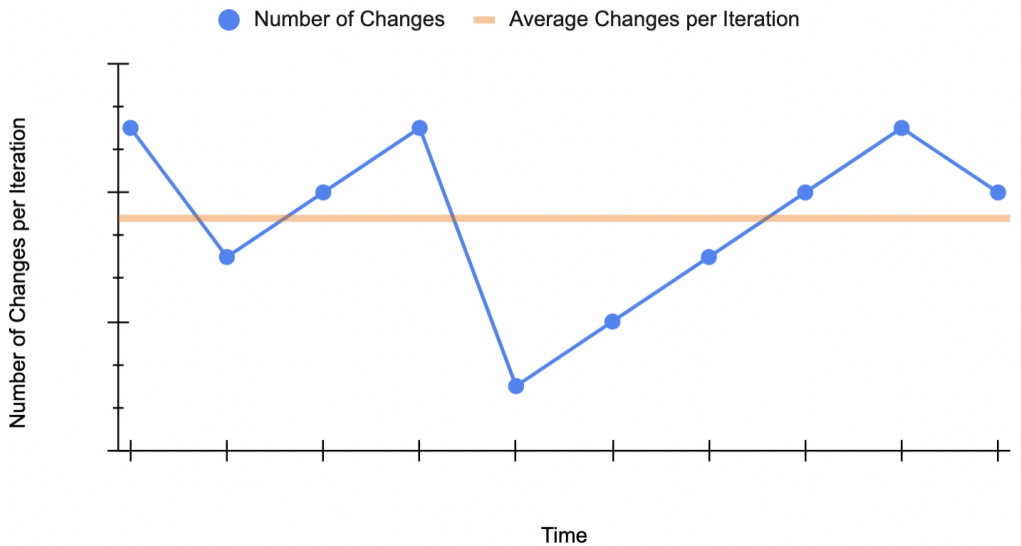


Figure 2.1: A simplified visual of Lehman’s fifth law, “Conservation of Familiarity.”

benefit from the Internet that has positively improved the experience. Two common resources currently available to developers have impacted software evolution [2]:

1. The rapid growth of the World Wide Web and Internet Resources make it easier for users and engineers to find related information.
2. Open source development where anybody could download the source codes and modify it has enabled fast and parallel evolution (through forks).

These two suggestions are very evident in modern development. For ex-

ample, a developer may regularly use resources like StackOverflow to find solutions to problems and use open-source tools that the developer and their team can contribute to or adjust to their specific needs.

2.2.3 Measuring Maintainability

Despite the nuanced differences between *maintainability* and *evolution*, the two characteristics run parallel to each other. If a system is easy to maintain, it will also be easier to evolve. If we can measure our system's maintainability, we can also determine if our system is in a good position to continue evolving to meet our future needs.

Several tools attempt to provide some value around these ideas. In this paper, we will focus on the metrics that Pylint provides, specifically looking into the Refactor score of Pylint.

We will look at many open-source Python systems using Pylint and attempt to correlate the data from the Pylint scores to the level of ease in adding new features to the system. This will determine if a system is more maintainable with better Pylint scores.

FUTURE EDITION: To do this, we will measure the locality of the changes by the number of files that are edited in a commit. We will also focus on commits that represent new features, not on commits that are bug fixes.

2.2.4 Maintainability Scores

First, let us consider our original understanding of software maintainability. While this definition focuses primarily on bug fixes and minor enhancements, maintainable projects should also have ease in their ability to evolve. Therefore, we can study the impact maintainability (structural quality) has on software evolution by reviewing the scores provided by automated code review tools.

In this study, we will be using Pylint and will be focused on the values of the Refactor score regarding a set of open-source Python systems. To understand the scores we will be working with, we must understand what Pylint itself is doing.

Through the documentation of Pylint, we can understand how to use it and the scores it will provide [6]. The Pylint score itself is calculated by the following equation [7]:

$$10.0 - ((\text{float}(5 * e + w + r + c) / s) * 10)$$

Numbers closer to 10 reflect systems that have fewer errors, fewer warnings, and have overall better structure and consistency. In the above equation, we are using the following values [8]:

- **statement** (**s**): the total number of statements analyzed
- **error** (**e**): the total number of errors, which are likely bugs in the code

- **warning** (**w**): the total number of warnings, which are python specific problems
- **refactor** (**r**): the total number of refactor warnings for bad code smells
- **convention** (**c**): the total number of convention warnings for programming standard violations

The Refactor score is of special interest to us and considers many features that are meticulously outlined on the Pylint site [9]. These types of warnings include a number of checks, such as when a boolean condition could be simplified, or a useless `return`, and so on. This score, in particular, will be part of our focus.

To calculate the Refactor score, Pylint will check the code for code smells based on the definitions for checks that have been documented. For every infraction, the score increases by one count.

“In computer programming, a **code smell** is any characteristic in the source code of a program that possibly indicates a deeper problem.” [10]

We can use these Refactor scores to help us spot architecture smells. After all, code smells can point the way to deeper problems in our system. There are fundamental design principles that have been established that we should consider when creating software; code smells alert us to areas that

have deviated from these principles. These smells are drivers for refactoring and, when addressed, can help us maintain the integrity of our architecture rather than creating a patchwork construction.

Because of the relation of refactor scores to the code structure itself, we will be spending much of our focus on this particular value. The most common Refactor error returned in our data set was the `no-else-return` message. This particular message highlights when an unnecessary block of code follows an if-statement that contains a `return`. The second most common Refactor message was `too-few-public-methods`, which reminds the developer to consider whether that class is appropriate to create.

Finally, in respect to Python, it is also helpful to be familiar with PEP 8, as this is the default set of standards that Pylint uses to judge Python code [11]. This standard can be used to make code more readable and more consistent, which may contribute to the code being more maintainable. These standards cover things like indentation spacing, maximum line length, where to break lines, how to handle imports, and more. By defining a set of standards, teams can ensure they have a defined set of rules so that any contributors to the code understand the expectations (and so that automated systems like Pylint can enforce those standards to maintain readability and consistency).

2.2.5 Other Maintainability Characteristics

The authors of “Measurement and refactoring for package structure based on complex network” recently reviewed a similar idea focusing on cohesion and coupling over time for a project [12]. In a software system, we desire low coupling (allowing for changes to one area to remain independent of changes to another area) and high cohesion (indicating reduced complexity in modules, which improves maintainability). Through a few experiments on open-source software systems, the authors determined that their algorithm that calculated metrics was capable of improving package structures to have high cohesion and low coupling. Their study gives us confidence that metrics around the software’s structure can provide value in keeping systems in a maintainable state, which allows for software evolution.

Another variable that may impact the maintainability of code is readability. For example, in the article “How does code readability change during software evolution?” the authors have addressed this concern and found that most source codes were readable within the sample they reviewed. Additionally, a minority of commits changed the readability; if a file was created as less readable, it was likely that it remained that way and did not improve [13]. This variable (readability) in the maintainability of a software system can influence how easy or difficult it is to make a change. The authors also found that big commits, usually associated to adaptive changes (a form of software evolution), were the most prone to reduce code readability [13]. This assumes that smaller commits are almost always better and can lead to more readable

code.

Piantadosi et al. found that changes in readability, whether improvements or disintegrations, often occurred unintentionally [13]. By enforcing the PEP 8 standard, we know that Pylint is encouraging systems to remain readable. Therefore, projects that use some form of automated system in their pipeline benefit from keeping their project on track in this regard, limiting the effects of readability on a software’s potential for evolution.

The paper “Standardized code quality benchmarking for improving software maintainability” provides additional insights into how the code’s maintainability is impacted by the technical quality of source code [14]. Within their paper, the authors seek to show four key points: (1) how easy it is to determine where and how the change is made, (2) how easy it is to implement the change, (3) how easy it is to avoid unexpected effects, and (4) how easy it is to validate the changes. Their approach has shown that some tools and methods can be used to improve and maintain technical quality within their projects, allowing systems to continue to evolve at a reasonable pace.

2.2.6 Documentation and Maintainability

Our assumption is that the Refactor score in projects should correlate to the evolution of the system. The first pass through the data is not conclusive in this particular detail, as the projects reviewed have many other factors contributing to the evolution of the project (number of contributors, size of the code system, etc.). Our assumption is that the correlation between software quality and software evolution would indicate that the better-scoring code systems are readable in themselves. In addition, it would be helpful to understand whether there are any similarities in how a system is documented that could contribute to improved software evolution of a system.

The textbook, “Software Architecture in Practice,” chapter 18 provides some insight in documentation around architecture [15]:

“If you go to the trouble of creating a strong architecture, one that you expect to stand the test of time, then you *must* go to the trouble of describing it in enough detail, without ambiguity, and organizing it so that others can quickly find and update the needed information.”

The book describes how documentation holds the results of significant design decisions, providing valuable insights into decisions down the road. While not directly related to the Pylint Refactor score and not within the source code itself, it is still helpful to remind ourselves that documentation can also influence the ability of a software system to evolve.

Our “best scores” (regarding the current Pylint Refactor score) were found to have relatively organized and useful documentation. The code repository for *Munki* provided documentation for previous versions, lending insight into design decisions as the software evolved [16]. The repository for *Raven*, however, was a deprecated version that has since been replaced by a paid platform known as *Sentry*, but had ample documentation [17]. It is possible that the “death” of that software system was not lack of evolution, but rather a business decision. *ElastAlert* was another system with good scores and easy-to-follow documentation, though it is focused more for the use of the system rather than how to enhance the system itself [18].

When reviewing our “worst offenders” in current Refactor scores, it was noted that even with poor scores, these repositories were able to continue to see engagement from developers. While further inspection will be needed to understand whether the code itself is evolving or just has engagement from a maintenance level, it is interesting to note that there is decent documentation provided. *SymPy* goes as far as documenting the architecture for the software as well as design decisions, enabling developers to better understand the structure as they make contributions [19].

“Our study has shown that the primary studies provide empirical evidence on the positive effect of documentation of designs pattern instances on programme comprehension, and therefore, maintainability.”

“...developers should pay more effort to add such documentation, even if in the form of simple comments in the source code.”

In research done by Wedyan and Abufakher (quoted above), it was found that documenting design patterns was useful in enhancing code understanding [20]. In turn, the comprehensibility impacts the maintainability of the code in a positive way, which continues to reinforce the impact that documentation can have and how it ties well into considerations for software structure.

2.3 Related Work

In our research, we are running with the assumption that Maintainability Index (MI) is our primary indicator and we will look for the correlation between the MI and other Pylint scores. We hope to find which correlations align and which are the most important.

2.3.1 Considering Data Sets

When exploring the correlations of maintainability and refactoring, there are many sources available for research. Some researchers have looked at proprietary systems as they evolve over time, while others have chosen open-source code available to the general public.

A study conducted by Baishakhi Ray, Daryl Posnett, Premkumar Devanbu, and Vladimir Filkov begins by programmatically collecting a sample set of projects in GitHub that vary in languages. Then the group of projects is appropriately culled, resulting in a final set used for the review. The results are then studied for the impact different programming languages may have on the code quality [21]. Through their research, they were able to determine which languages were more prone to defects, and that individual languages are more related to individual bugs rather than bugs overall.

The authors of “Predicting Maintainability with Object-Oriented Metrics - An Empirical Comparison” performed a very similar study to what we are doing here. Their study focuses on object-oriented software (specifi-

cally C/C++ and Java) and a correlation analysis between object-oriented metrics and software maintainability, looking for the best metrics to predict maintainability [22]. This particular study focuses on a few hand-picked software systems with an analysis of the change logs. Our study, however, will be of a larger scale (about 50 software systems) and focused solely on Python-heavy projects.

2.3.2 Design Patterns and Software Quality

In the paper “Impact of design patterns on software quality: a systematic literature review” the authors compared the use of design patterns to software evolution and maintainability. They found that design patterns provided clear flexibility when they reviewed changes that extended (evolved) software [20].

“Changes performed in a class can be corrective, adaptive, perfective, or preventive. These changes can occur due to new requirements, debugging, changes that propagate from changes in other classes and refactoring.”

Wedyan and Abufakher found that there were two reasons that a class had more frequent changes [20]:

1. The class was easy to extend.

2. The class correlated to other classes (raising alarms about class modularity).

With these findings in mind, we intentionally aim to focus our research on changes for system extensions and adaptations rather than bug fixes that appeared to be larger change due to high coupling. Within this paper, we were able to focus on Refactor scores (code smells) rather than Error scores (bugs) within the system.

2.3.3 Software Architecture and Maintainability

The research done in “Software Architecture Metrics: A Literature Review”, the authors discuss how early detection of issues within the software’s architecture is key to mitigating the risk of poor performance and can lower the cost of repairing faults [23]. While most developers have had access to these types of metrics for several decades, the industry and open-source community have not really latched onto their use for keeping code in easy-to-work-with condition.

The review done by Coulin et al. called out five important qualities of software architecture [23]:

1. Maintainability
2. Extensibility
3. Simplicity, Understandability

4. Re-usability

5. Performance

Focusing on these qualities can narrow down the choice between different design options to end up with the most ideal solution. Keeping these five qualities in top-of-mind for new (and changed) code allows for easier future development and evolution of the software system.

Chapter 3

Methodology

TODO: Chapter 3 - how I did it

3.1 Initial Repository Set

We have established that we have a problem with projects that fail to evolve, resulting in a loss of revenue. We also understand that evolving software is essential in order to keep users engaged; without it, there is an appearance in the decline of quality and the program becomes less satisfactory to the user, as well as potential for competitors to outpace us with features available. We must now understand how we can ensure that our systems evolve. For this, we will look to understand how the system’s structural quality impacts software evolution.

The work done by Dr. Omari and Dr. Martinez involves collecting a subset of Python projects that we can use for further research. The bulk of the effort they have provided is determining which classifiers to use to pare down the public set of Python systems into a good collection for further analysis [24]. The work that they have provided was used to select appropriate Python systems for review by collecting meta-data on these code systems.

From their subset of repositories, we were then able to collect current Pylint scores from each of our 129 systems. This gives us a sampling of data that we can now dig deeper into, comparing similar systems (similar size, similar number of contributors, etc.) and their evolution process by reviewing past commits rather than merely the current state of the system, as we have done here.

3.2 Filtered Repository Set

Once we had a narrowed set of projects from GitHub that were primarily written in Python, we culled the set more using several criteria:

1. Projects that are at least 80% Python
2. Projects with a long history of commits
3. Projects with large development teams (contributors)
4. Projects with many releases
5. Projects of a substantial age

Armed with this list, we were able to use the metadata from GitHub for each of our repositories already collected and determine a cross-section of these criteria that would result in about 50 repositories for further study.

Beginning with the languages field from GitHub, we could easily narrow down projects that had at least 80% of the code in Python. In our set, 103 repositories contained 80% or more Python code.

With this narrowed set, we then looked to see at which percentile all the remaining criteria would yield the desired number of repositories. We determined that using the value at the 20th percentile in each of the above categories would yield the size set we'd need.

| Criteria | 20th Percentile Value |
|------------------------|-----------------------|
| Number of Commits | 2,968 |
| Number of Contributors | 90 |
| Number of Releases | 44 |
| Age (in months) | 66.4 |

Table 3.1: Criteria used to filter down the initial set of repositories.

Table 3.1 shows the values found for each of our criteria. Using these values as our minimum requirements, we can narrow our repository set to 46 repositories.

Chapter 4

Results

TODO: Chapter 4 - what I found

Chapter 5

Conclusions and Recommendations

TODO: Chapter 5 - what it all means, putting the pieces together (what is my contribution to the research field)

By collecting data and drawing our conclusions from it, with help from the insights from the studies done before ours, we may better understand metrics that can be useful regarding maintainability. Good projects will inevitably continue to grow and evolve. Understanding methods to keep code refactor on a certain level makes code easy to change. We may also find that projects with worsening scores slow down with updates and have reduced engagement.

When reviewing our surface-level data with current project Refactor scores, our three worst offenders were *Ansible* [25], *SymPy* [26], and *Salt* [27]. All three projects are still quite active with development despite their poor cur-

rent scores. The projects have high download rates, which may be the reason for continued development despite potential difficulty in maintenance.

Projects that may be open source or have many contributors are especially vulnerable to maintainability degrading over the evolution of a project. Having a reliable metric can be very useful in programmatically avoiding code smells and keeping code in a state that is easy to manage through simple metric checks in deployment pipelines.

We can see an example of this in reviewing some current symptoms that *SymPy* is experiencing, with only 72% code coverage and a failing build (see “Fig. 5.1”). Despite the engagement and continued development, we suspect that real adaptations and evolution of the software may be difficult with this code.



Figure 5.1: A snapshot of the badges from *SymPy*’s repository.

We have further work to do in this study to gain better understanding. With a set of several “best” and “worst” Python software systems, we will look into the history of the projects’ commits. It would be useful to see how the Refactor scores have changed over time, and if the rate at which changes were pushed correlated to the increase or decrease in that Refactor score.

Additional data can be gathered from this set that may provide more insights than this first brush of the data provides us. Understanding the impact of structural quality on the evolution of a project can provide compelling

perspectives.

References

- [1] A. V. Deursen, “Think twice before using the “maintainability index”, ” 2014, <https://avandeursen.com/2014/08/29/think-twice-before-using-the-maintainability-index/> [Online; accessed 23-Jan-2022]. [Online]. Available: <https://avandeursen.com/2014/08/29/think-twice-before-using-the-maintainability-index/>
- [2] W. contributors, “Software evolution — Wikipedia, the free encyclopedia,” 2021, https://en.wikipedia.org/wiki/Software_evolution [Online; accessed 12-December-2021]. [Online]. Available: https://en.wikipedia.org/wiki/Software_evolution
- [3] —, “Software maintenance — Wikipedia, the free encyclopedia,” 2021, https://en.wikipedia.org/wiki/Software_maintenance [Online; accessed 17-December-2021]. [Online]. Available: https://en.wikipedia.org/wiki/Software_maintenance
- [4] F. Brooks, *The Mythical Man-Month*. Addison-Wesley, 1975.

- [5] W. contributors, “Lehman’s laws of software evolution,” 2021, https://en.wikipedia.org/wiki/Lehman%27s_laws_of_software_evolution [Online; accessed 12-December-2021]. [Online]. Available: https://en.wikipedia.org/wiki/Lehman%27s_laws_of_software_evolution
- [6] Logilab and contributors, “Pylint,” Logilab, 2020, <https://pylint.org/> [Online; accessed 14-December-2021]. [Online]. Available: <https://pylint.org/>
- [7] P. Logilab and contributors, “Pylint features,” Logilab and PyCQA, 2021, https://pylint.pycqa.org/en/latest/technical_reference/features.html#reports-options [Online; accessed 14-December-2021]. [Online]. Available: https://pylint.pycqa.org/en/latest/technical_reference/features.html#reports-options
- [8] R. Kirkpatrick, “A beginner’s guide to code standards in python - pylint tutorial,” 2016, <https://docs.pylint.org/en/1.6.0/tutorial.html> [Online; accessed 18-December-2021]. [Online]. Available: <https://docs.pylint.org/en/1.6.0/tutorial.html>
- [9] P. Logilab and contributors, “Pylint features,” Logilab and PyCQA, 2021, https://pylint.pycqa.org/en/latest/technical_reference/features.html#refactoring-checker [Online; accessed 14-December-2021]. [Online]. Available: https://pylint.pycqa.org/en/latest/technical_reference/features.html#refactoring-checker

- [10] W. contributors, “Code smell — Wikipedia, the free encyclopedia,” 2021, https://en.wikipedia.org/wiki/Code_smell [Online; accessed 18-December-2021]. [Online]. Available: https://en.wikipedia.org/wiki/Code_smell
- [11] P. S. Foundation and contributors, “Pep 8 – style guide for python code,” Heroku Application, 2021, <https://www.python.org/dev/peps/pep-0008/> [Online; accessed 14-December-2021]. [Online]. Available: <https://www.python.org/dev/peps/pep-0008/>
- [12] Y. Zhou, Y. Mi, Y. Zhu, and L. Chen, “Measurement and refactoring for package structure based on complex network,” *Applied Network Science*, vol. 5, no. 50, 2020.
- [13] V. Piantadosi, F. Fierro, S. Scalabrino, A. Serebrenik, and R. Oliveto, “How does code readability change during software evolution?” *Software Qual J*, vol. 25, pp. 5374–5412, 2020.
- [14] R. Baggen, José, P. Correia, K. Schill, and J. Visser, “Standardized code quality benchmarking for improving software maintainability,” *Software Qual J*, vol. 20, pp. 287–307, 2012.
- [15] L. Bass, P. Clements, and R. Kazman, *Software Architecture in Practice: Third Edition*. Addison-Wesley Professional, 2012.

- [16] M. contributors, “Munki,” 2021, <https://github.com/munki/munki> [Online; accessed 17-December-2021]. [Online]. Available: <https://github.com/munki/munki>
- [17] S. contributors, “Sentry,” 2021, <https://github.com/getsentry/raven-python> [Online; accessed 18-December-2021]. [Online]. Available: <https://github.com/getsentry/raven-python>
- [18] E. contributors, “Elastalert,” 2021, <https://github.com/Yelp/elastalert> [Online; accessed 18-December-2021]. [Online]. Available: <https://github.com/Yelp/elastalert>
- [19] S. D. Team, “SymPy user’s guide,” 2019, <https://web.archive.org/web/20200403130424/https://docs.sympy.org/latest/guide.html> [Online; accessed 18-December-2021]. [Online]. Available: <https://web.archive.org/web/20200403130424/https://docs.sympy.org/latest/guide.html>
- [20] F. Wedyan and S. Abufakher, “Impact of design patterns on software quality: a systematic literature review,” *IET Software*, vol. 14, no. 1, 2020.
- [21] B. Ray, D. Posnett, P. Devanbu, and V. Filkov, “A large-scale study of programming languages and code quality in github,” *COMMUNICATIONS OF THE ACM*, vol. 60, no. 10, pp. 91–100, 2017.

- [22] M. Dagpinar and J. H. Janke, “Predicting maintainability with object-oriented metrics - an empirical comparison,” *Reverse Engineering - Working Conference Proceedings*, pp. 155–164, 12 2003.
- [23] T. Coulin, M. Detante, W. Mouchère, and F. Petrillo, “Software architecture metrics: A literature review,” 2019.
- [24] S. Omari and G. Martinez, “Enabling empirical research: A corpus of large-scale python systems,” 2018, [Provided by Dr. Omari].
- [25] A. contributors, “Ansible,” 2021, <https://github.com/ansible/ansible> [Online; accessed 17-December-2021]. [Online]. Available: <https://github.com/ansible/ansible>
- [26] S. contributors, “SymPy,” 2021, <https://github.com/sympy/sympy> [Online; accessed 18-December-2021]. [Online]. Available: <https://github.com/sympy/sympy>
- [27] —, “Salt,” 2021, <https://github.com/saltstack/salt> [Online; accessed 18-December-2021]. [Online]. Available: <https://github.com/saltstack/salt>

Appendix A

Type or paste your appendices here. Appendices are a place to organize and include all of the “extra” material that is important to your research work but that is too detailed for the main text. Examples can include: specific analytical methods, computer code, spreadsheets of data, details of statistical analyses, etc. But, these materials do not speak for themselves. There should be a reference to these materials from the main chapters (complete details included in Appendix A) and there should be some text at the beginning of each appendix to briefly explain what the information is and means that is included in that appendix.

TODO: Consider a list of the pylint messages for refactor scores, etc.