
Abstract

Software is contributing a substantial part of new functionality and innovations in safety-critical systems. These systems put a huge demand on software reliability, because a minor error can produce failure of a complete system. The evolution of software requires continuous development and maintenance. With size and complexity of safety-critical software growing as time goes, additional challenges arise, including implicit assumptions of *technical debt*. *Technical debt* refers to the sum of compromises that are made in software development and maintenance to meet a short-term business goal. *Design debt* is an instance of technical debt. As software systems evolve, their design tends to decay over time, hence accumulating design debt. Consequently, software design becomes more difficult to maintain. Therefore, the developers need to understand the reason for design debt accumulation so they can take proactive steps that may potentially reduce the debt in the future.

The main goal of this thesis is to empirically investigate design debt in safety-critical systems. The goal is reflected in our attempt to answer the following research questions:

RQ1: How can design debt be identified?

RQ2: What kind of design debt can be found in embedded systems?

RQ3: What are the effects of design debt?

RQ4: How to pay design debt?

A case study has been conducted to answer our research questions. The case study involves an analysis of a safety-critical system developed by Autronica Fire and Security AS. The system is written in C/C++. We used object-oriented metrics to identify classes that are most likely to pose problems for the system. Quantitative data were collected and analyzed using descriptive statistics. A set of thresholds for the metrics were derived to identify classes that have higher metric values than its threshold values. In addition, automatic static analysis tools were applied in order to detect code smells.

This work contributes mainly to improvement in software metrics and software quality. The stated contributions of this work are:

- **C1:** Empirical knowledge about design debt identification in safety-critical systems by object-oriented metric analysis and code smell detection.
 - **C1.1:** A set of threshold values for the measured object-oriented metrics.
- **C2:** Empirical knowledge about the effects of having design debt in safety-critical systems.
- **C3:** Empirical knowledge about the different types of design debt in safety-critical systems.
- **C4:** Empirical knowledge about paying design debt.

Sammendrag

Programvare bidrar en vesentlig del i utvikling av ny funksjonalitet og innovasjoner i sikkerhetskritiske systemer. Slike systemer er avhengig av programvarens pålitelighet, fordi en liten feil kan medføre svikt i et system. Evolusjon av programvaresystemer krever kontinuerlig utvikling og vedlikehold. Nye utfordringer oppstår ettersom størrelsen og kompleksiteten til sikkerhetskritiske vokser. Dette inkluderer implisitte forutsetninger om teknisk gjeld. Teknisk gjeld oppstår når man i utviklingsarbeidet velger suboptimale måter å løse problemer på. Design gjeld er en instans av teknisk gjeld. Designet til et system har en tendens til å forfalle over tid ettersom et programvaresystem utvikler seg. Dette fører til en oppsamling av design gjeld. Utviklere er dermed nødt til å forstå hvorfor design gjeld akkumulerer slik at de kan ta gradvise steg for å redusere gjelden.

Hovedmålet med oppgaven er å undersøke design gjeld i sikkerhetskritiske systemer. Målet med oppgaven er reflektert i vår forsøk på å svare på følgende forskningsspørsmål:

FS1: Hvordan kan design gjeld identifiseres?

FS2: Hva slags type design gjeld finnes i sikkerhetskritiske systemer?

FS3: Hva slags virkinger har design gjeld i sikkerhetskritiske systemer?

FS4: Hvordan kan design gjeld tilbakebetales?

Et case studie har blitt gjennomført i et forsøk på å svare på våres forskningsspørsmål. Case studiet involverer en analyse av et sikkerhetskritisk system utviklet av Autronica Fire and Security AS. Systemet er utviklet i C/C++. Vi har brukt objekt-orienterte metrikker for å identifisere klasser som mest sannsynlig vil påføre problemer for systemet. Kvantitative data ble samlet ved hjelp av verktøy, og analysert ved hjelp av deskriptiv statistikk. Et sett med . I tillegg

Resultatene av arbeidet bidrar hovedsakelig til forbedring av programvaremetrikker og programvarekvalitet. Bidragene for dette arbeidet er som følger:

- **B1:** Empirisk kunnskap om identifikasjon av design gjeld i sikkerhetskritiske systemer ved analyse av objekt-orienterte metrikker og code smells.
- – ????
- Empirisk kunnskap om virkningene til design gjeld i sikkerhetskritiske systemer.
- Empirisk kunnskap om ulike typer av design gjeld i sikkerhetskritiske systemer.
- Empirisk kunnskap om tilbakebetaling av design gjeld.

Preface

This thesis was written as a part of my MSc degree at Norwegian University of Science and Technology (NTNU). It was written in the spring of 2016 at the Department of Computer and Information Science (IDI), NTNU, Trondheim.

This thesis is an extension of the work done in the specialization project - "Managing Technical Debt in Embedded Systems" that was carried out by the same author in the fall of 2015. The supervisor for both projects was Carl-Fredrik Sørensen.

Trondheim, June 10, 2016

CONTENTS

Abstract	i
Sammendrag	iii
Preface	v
Table of Contents	viii
List of Tables	ix
List of Figures	xi
Abbreviations	xiii
1 Introduction	1
1.1 Motivation	1
1.2 Research Context	2
1.3 Research Design and Questions	3
1.4 Contribution	3
1.5 Thesis Structure	4
2 State-of-the-Art	5
2.1 Technical Debt	5
2.1.1 Definitions of Technical Debt	6
2.1.2 Classification of Technical Debt	7
2.1.3 Causes and Effects of Technical Debt	7
2.1.4 Identification of Technical Debt	9
2.1.5 Strategies and Practices for Managing Technical Debt	13
2.2 Design Debt	13

2.3	Software Quality	14
2.4	Object-Oriented Metrics	14
2.5	Software Evolution and Maintenance	15
2.6	Software Reuse	16
2.7	Refactoring	17
2.8	Embedded Systems	17
2.8.1	Safety-Critical Systems	18
3	Research Methodology	21
3.1	Research Methods in Software Engineering	21
3.2	Choice of Research Method in this Thesis	22
3.2.1	Case Study Method	22
3.3	Case Context	23
3.4	Research Process	24
3.4.1	Determine and Define the Research Questions	24
3.4.2	Select the Cases and Determine Data Gathering and Analysis Techniques	25
3.4.3	Prepare To Collect Data	26
3.4.4	Data Collection	28
3.4.5	Evaluate and Analyze the Data	30
3.4.6	Prepare the Report	30
3.5	Summary of the Research Design	30
4	Results	31
4.1	Object-Oriented Metrics in Project "Firmus"	31
4.1.1	Object-Oriented Metrics for the Components	36
4.2	Identifying Code Smells using Automatic Static Analysis Tools	54
5	Discussion	59
5.1	Analysis of Object-Oriented Metrics by Applying Threshold Values	59
5.2	Research Evaluation	64
5.3	Threats To Validity	69
5.3.1	Internal Validity	69
5.3.2	External Validity	69
5.3.3	Conclusion Validity	69
6	Conclusion	71
6.1	Future Work	72
	References	73

LIST OF TABLES

2.1	Types of TD	8
2.2	Code Smell Taxonomy	11
2.3	Software quality attributes, criteria, and description (ISO/IEC 9126) . .	19
3.1	System Metrics for Project "Firmus"	24
3.2	Research Questions	25
4.1	OO-metrics and descriptive statistics for Project "Firmus"	32
4.2	OO-metrics and descriptive statistics for Component A	37
4.3	OO-metrics and descriptive statistics for Component B	39
4.4	OO-metrics and descriptive statistics for Component C	41
4.5	OO-metrics and descriptive statistics for Component Ex	44
4.6	OO-metrics and descriptive statistics for Component G	46
4.7	OO-metrics and descriptive statistics for Component L	48
4.8	OO-metrics and descriptive statistics for Component N	50
4.9	OO-metrics and descriptive statistics for Component P	52
4.10	OO-metrics and descriptive statistics for Component S	54
4.11	Number of Code Smells detected	55
4.12	Duplicated Code in Project "Firmus"	55
4.13	Speculative Generality in Project "Firmus"	56
4.14	Dead Code in Project "Firmus"	57
5.1	Thresholds for OO-software metrics	60
5.2	Connection between contributions, and research questions	64

LIST OF FIGURES

2.1	Fowler's TD Quadrant	6
2.2	TD Landscape	7
4.1	Frequency chart of the LCOM metric	33
4.2	Frequency chart of the CBO metric	33
4.3	Frequency chart of the DIT metric	34
4.4	Frequency chart of the NOC metric	34
4.5	Frequency chart of the RFC metric	35
4.6	Frequency chart of the WMC metric	36
4.7	Frequency distribution of OO-metrics in Component A	38
4.8	Frequency distribution of OO-metrics in Component B	40
4.9	Frequency distribution of OO-metrics in Component C	42
4.10	Frequency distribution of OO-metrics in Component Ex	45
4.11	Frequency distribution of OO-metrics in Component G	47
4.12	Frequency distribution of OO-metrics in Component L	49
4.13	Frequency distribution of OO-metrics in Component N	51
4.14	Frequency distribution of OO-metrics in Component P	53

Abbreviations

ASA	=	Automatic Static Analysis
CBO	=	Coupling Between Objects
CC	=	Cyclomatic Complexity
DD	=	Design Debt
DIT	=	Depth in Inheritance Tree
LCOM	=	Lack of Cohesion in Methods
NIM	=	Number of Instance Methods
NIV	=	Number of Instance Variables
NOC	=	Number Of Children
OO	=	Object-Oriented
RFC	=	Response For Class
TD	=	Technical Debt
WMC	=	Weighted Methods per Class

CHAPTER 1

INTRODUCTION

This chapter provides an introduction to this masters thesis. We begin with outlining the motivation for this research. Then a brief description of the research context and the research questions is presented. Lastly, we present the thesis outline.

1.1 Motivation

Successful embedded systems continuously evolve in response to external demands for new functionality and bug fixes [1]. Evolution occurs as software changes over time. One consequence of such evolution is an increase of issues in design, development, and maintainability [2]. Software code often ends up not contributing to the mission of the original intended software architecture or design.

The main challenge with software evolution is the **Technical Debt** (TD) that is not paid by the organization during software development and maintenance. TD addresses the debt that software developers accumulate by taking shortcuts in development in order to meet the organizations business goals. For example, a deadline may lead developers to create "non-optimal" solutions to deliver on time. As TD keeps accumulating, software systems may potentially become unmanageable and eventually unusable. Even more resources have to be spent during software maintenance on paying off the interest, i.e., the cost of having the debt. According to Gartner [3], the cost of dealing with TD threatens to grow to \$1 trillion globally by 2015. That is the double of the amount of TD in 2010. Furthermore, many embedded systems are getting interconnected within existing Internet infrastructure, known as the Internet of Things. Such devices are threatened by security issues. For example, Matthew Garret got access to the electronic equipment in

every hotel room in a hotel located in London¹. These equipments were connected to a network. Additionally, two research discovered the possibility to start a Tesla Model S using a laptop².

Several studies have classified the metaphor of TD into different types of debt that are associated with the different phases of software development [4–9]. **Design Debt** (DD) is an instance of TD. DD accumulates when compromises are made in the software design. Software design plays a significant role in the development of large systems [10]. Unlike code-level debt, DD usually has more significant consequences on software evolution by making the system more complex and harder to maintain over time [11, 12]. The quality of the software is very important, and without it, there will be more system failures that may lead to accidents. This is even more crucial for the case of safety-critical software, which failure can endanger human lives.

1.2 Research Context

This master thesis builds upon our previous study "Managing TD in Embedded Systems" [13], a prestudy that was carried out in the fall of 2015. The written assignment for the specialization project had the following definition:

Managing TD in embedded systems

"This task is related to management of software in embedded systems, as well as evolution of such software over time. Embedded systems have often a long lifetime and it is thus important to find out best practices and tools for this management since it is necessary to cope with architectural and design decisions which were made perhaps decades ago, as well as clearly find out how present decisions may affect future maintenance and operation. This is called TD since all decisions will have a future cost related to them. Such decisions are often not documented, the people that made the software is not available 10-20 years after the implementation, the Internet of Things make all kind of embedded systems accessible from the Internet and thus posing security threats.

The project may take different directions based on the students interests and motivation. Industrial companies are very interested in this topic, so it is possible to study industrial systems, both past and current., make suggestions and implement them, make tools, make processes, make best practice etc."

In our previous research, we did a pre-study on the field of TD in embedded systems. We investigated the reasons for organizations to incur TD, and the different strategies for managing it. Research data was collected by conducting semi-structured interviews. After completing the study, we had a desire to look into a more narrow field of the concept

¹The Internet of dangerous, broken things: <http://www.zdnet.com/article/the-insecurity-of-the-internet-of-things/>

²Researchers Hacked a Model S, But Tesla's Already Released a Patch: <http://www.wired.com/2015/08/researchers-hacked-model-s-teslas-already/>

TD by conducting a deeper case study for our upcoming master thesis. A big interest was to study an industrial system.

The work in our master thesis has been done in collaboration with Autronica Fire and Security AS, a global provider of safety solutions, including fire safety equipment, marine safety monitoring, and surveillance equipment. Their main office is located in Trondheim, one of the largest cities in Norway. We have performed a deeper study on one of the company's fire detection systems software for approximately six weeks. The case study explored their source code in order to investigate DD.

1.3 Research Design and Questions

Being able to identify DD helps the software engineering field to get one step closer on solving the problems that are faced by the software industry today. The goal of the research is to investigate DD in safety-critical systems. DD is a problem for many software projects today. Although DD is recognized in many studies, there is lack of focus on DD in safety-critical systems.

The relevant research methods in software engineering can be survey, design and creation, case study, experimentation, action research, and ethnography [14]. In this study, a literature review and case study have been used to answer our research questions. Literature review was a part of the pre-study. It has been used to get familiar with the term DD and to define the research questions. Case studies are empirical methods used to investigate a single entity or phenomenon within a specific time space [15], which fits our desire to study an industrial system. Case studies can be both qualitative and quantitative [14, 16]. The research process we have chosen to adopt in this study follows the principles of the six steps defined by Soy [17]: *Determine and Define The Research Questions, Select the Cases and Determine Data Gathering and Analysis Techniques, Prepare to Collect Data, Data Collection, Evaluate and Analyze Data, and Prepare the Report.*

The main research questions investigated in this thesis are:

- **RQ1:** How can DD be identified?
- **RQ2:** What kind of DD can be found in embedded systems?
- **RQ3:** What are the effects of DD?
- **RQ4:** How to pay DD?

1.4 Contribution

This work contributes mainly to improvement in software metrics and software quality. The stated contributions of this work are:

- **C1:** Empirical knowledge about DD identification in safety-critical systems by object-oriented metric analysis and code smell detection.
 - **C1.1:** A set of threshold values for the measured object-oriented metrics.
- **C2:** Empirical knowledge about the effects of having DD in safety-critical systems.
- **C3:** Empirical knowledge about the different types of DD in safety-critical systems.
- **C4:** Empirical knowledge about paying DD.

1.5 Thesis Structure

The thesis is structured into several chapters with sections and subsections. The outline of the thesis is as follows:

- **Chapter 1:** Introduction contains a brief and general introduction to the study and the motivation behind it.
- **Chapter 2:** State-of-the-Art looks at important aspects of the research question.
- **Chapter 3:** Research Method describes how the literature review was carried out throughout the research, as well as a description of the case study to be performed.
- **Chapter 4:** Results presents the results from the case study, and takes a closer look at the findings from the case study.
- **Chapter 5:** Discussion contains a summarized look at the findings from the case study, and connects it with the literature review and to the research questions. An evaluation of the research is also given in this chapter.
- **Chapter 6:** Conclusion concludes the research by providing a summary of the most important points of the results and discussion chapter. Additionally, it outlines possible routes to take in the research field.

CHAPTER 2

STATE-OF-THE-ART

This chapter presents the relevant state-of-the-art for this thesis. Section 2.1 presents the metaphor of TD.

2.1 Technical Debt

The metaphor of TD was first introduced by Ward Cunningham in 1992 to communicate technical problems with non-technical stakeholders [18]. To deliver business functionality as quick as possible, *'quick and dirty'* decisions are often made. Such decisions may have short-term value, but they could affect future development and maintenance activities negatively. Cunningham was the first one who drew the comparison between technical complexity and financial debt in a 1992 experience report [18]:

“Shipping first time code is like going into debt. A little debt speeds up the development as long as it is paid back promptly with a rewrite... The danger occurs when the debt is not repaid. Every minute spent on not-quite-right code counts as interest on that debt. Entire engineering organizations can be brought to a stand-still under the debt load of an unconsolidated implementation, object-oriented or otherwise.” - Ward Cunningham, 1992.

The concept refers to the financial world where going into debt means repaying the loan with interest [19]. Like financial debt, TD accrues interest over time. Interest is defined as the extra effort that has to be dedicated in the future development in order to modify the part of the software that contains TD [4, 20, 21]. Unmanaged TD can cause projects to face significant technical and financial problems, which ultimately leads to increased maintenance and evolution costs [22].

2.1.1 Definitions of Technical Debt

Several researchers have attempted to give us a clear picture of what TD is [10, 23, 24]. Fowler [23] presented a TD quadrant which consists of two dimensions: *reckless/prudent* and *deliberate/inadvertent* [23]. The TD quadrant in Figure 2.1 shows four types of TD: *Reckless/Deliberate*, *Reckless/Inadvertent*, *Prudent/Deliberate*, and *Prudent/Inadvertent*. *Reckless/Deliberate* debt is usually incurred when technical decisions are taken intentionally without any plans on how to address the problem in the future. A team might know about good design practices, but they still choose to implement ‘quick and dirty’ solutions because they think they cannot afford the time required to write clean code. The second type is *Reckless/Inadvertent*. It is incurred when best practices for code and design are being ignored, ultimately leading to a big mess of “spaghetti code”. *Prudent/Deliberate* debt occurs when the value of implementing a ‘quick and dirty’ solution is worth the cost of incurring the debt to meet a short-term goal. The team are fully aware of the consequences, and have a plan on how to address the problem in the future. Lastly, we have *Prudent/Inadvertent debt*. This type of debt accumulates when a team realizes that the design of a valuable software could have been better after delivery. A software development process is as much as learning as it is coding.

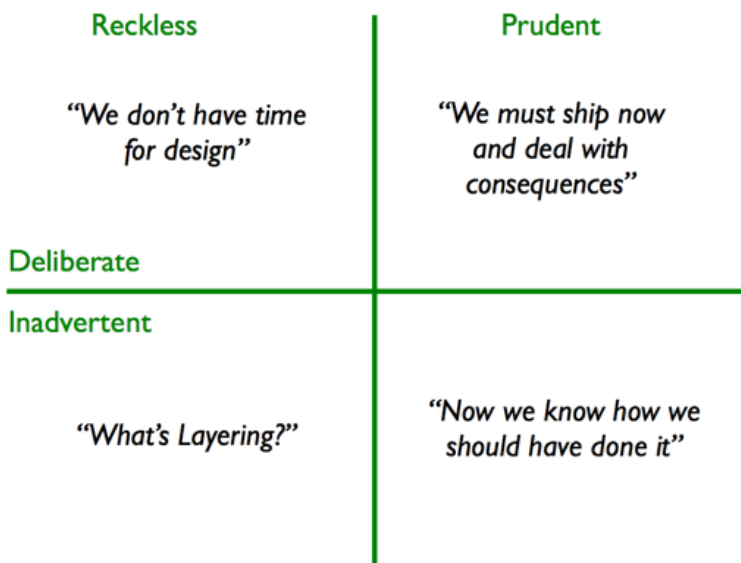


Figure 2.1: Fowler's TD Quadrant

McConnell [24] classified TD as *intentional* and *unintentional debt*. *Intentional debt* is described as debt that is incurred deliberately. For example, an organization makes a strategic decision that aims to reach a certain objective by taking a shortcut they are fully aware of. *Intentional debt* can further be viewed as “short-term” and “long-term debt” [5, 25]. Short-term debt is usually incurred reactively, for tactical reasons. Long-term debt is usually incurred pro-actively, for strategic reasons. *Unintentional debt* is described as

debt that is incurred inadvertently due to lack of knowledge or experience. For example, a junior software developer may write low quality code that does not conform with standard coding standard due to low experience.

Krutchén et al. [10] presented a TD landscape for organizing TD. They distinguished visible elements such as new functionality to add or defects to fix, and the invisible elements that are only visible to software developers. On the left side of Figure 2.2, TD affects evolvability of the system, while on the right side, TD mainly affects software maintainability.

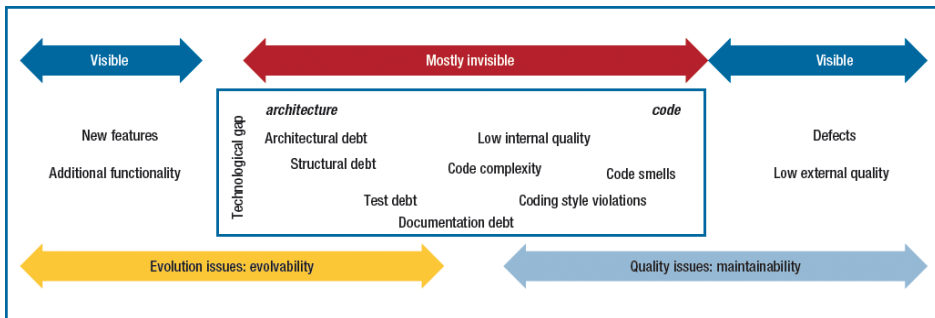


Figure 2.2: TD Landscape

2.1.2 Classification of Technical Debt

TD can accumulate in many different ways, and therefore it is important to distinguish the various types of TD. Several studies [4–9] have described several subcategories of TD based on its association with traditional software life-cycle phases; *architectural debt*, *code debt*, *defect debt*, **design debt** (DD), *documentation debt*, *infrastructure debt*, *requirements debt*, and *test debt*. Table 2.1 lists the different subcategories of TD.

2.1.3 Causes and Effects of Technical Debt

Several researchers have investigated the reasons for projects to incur TD. Klinger et al. [21] conducted an industrial case study at IBM where four technical architects with different backgrounds were interviewed. The goal of their study was to examine how decisions to incur debt were taken, and the extent to which the debt provided leverage [21]. The study revealed that the company failed to assess the impact of intentionally incurring debt on projects. Decisions regarding TD were rarely quantified. The study also revealed big organizational gaps among the business, operational, and technical stakeholders. When the project team felt pressure from the different stakeholders, TD decisions were made without quantifications of possible impacts.

Lim et al. [26] pointed out that TD is not always the result of poor developer disciplines, or sloppy programming. It can also include intentional decisions to trade off competing

Table 2.1: Types of TD

Subcategory	Definition
Architectural debt [4–6]	Architectural decisions that make compromises in some of the quality attributes, such as modifiability.
Code debt [4, 6, 7]	Poorly written code that violates best coding practices and guidelines, such as code duplication.
Defect debt [4, 7]	Defect, failures, or bugs in the software.
Design debt (DD) [4, 6, 8]	Technical shortcuts that are taken in design.
Documentation debt [4, 6, 9]	Refers to insufficient, incomplete, or outdated documentation in any aspect of software development.
Infrastructure debt [4, 5, 7]	Refers to sub-optimal configuration of development-related processes, technologies, and supporting tools. An example is lack of continuous integration.
Requirements debt [4, 9]	Refers to the requirements that are not fully implemented, or the distance between actual requirements and implemented requirements.
Test debt [4, 6, 9]	Refers to shortcuts taken in testing. An example is lack of unit tests, and integration tests.

concerns during business pressure. They explained that TD can be used in short-term to capture market share and to collect customers feedback early. In the long-term, TD tended to be negative. The trade-offs included increased complexity, reduced performance, low maintainability, and fragile code. This led to bad customer satisfaction and extra working hours. In many cases, the short-term benefits of TD outweighed the future costs.

Guo et al. [27] studied the effects of TD by tracking a single delayed maintenance task in a real software project throughout its life-cycle, and simulated how managing TD can impact the project result. Their results show that a delay in the maintenance task would have tripled the cost if it had been done later.

Siebra et al. [28] carried out an industrial case study where they analyzed documents, emails, and code files. Additionally, they interviewed multiple developers and project managers. Their case study revealed that TD were mainly taken by strategic decisions. Furthermore, they pointed out that the use of a unique specialist could lead the development team to solutions that the specialist wanted and believe were correct, leading the team to incur debt. Their study also revealed that TD can both increase and decrease the amount of working hours.

Zazworka et al. [29] studied the effects of *God Class* code smell and TD on software quality. *God Class* are examples on bad coding, and therefore includes a possibility for refactoring [8]. The results indicated that *God Class* code smell requires more maintenance effort including bug fixing and changes to software that are considered as a cost to software project. In other words, if developers desire higher software quality, then TD needs to be addressed closely in the development process.

Buschmann [30] explained three different stories of TD effects. In the first case, TD accumulated in a software platform had grew to a point where development, testing,

and maintenance costs started to increase dramatically. Additionally, the software components were hardly usable. In the second case, developers started to use shortcuts to increase the development speed. This resulted in significant performance issues because an improper software modularization reflected organizational structures instead of the system domains. It ended up turning into economic consequences. In the last case, an existing software product experienced an increase in maintenance cost due to an architectural erosion. However, management analyzed that re-engineering the whole software would cost more than doing nothing. They decided not to do address the debt because it was cheaper from a business point-of-view.

Codabux et al. [5] carried out an industrial case study where the topic was agile development focusing on TD. They observed and interviewed developers to understand how TD is characterized, addressed and prioritized, and how decisions led to TD. Two subcategories of TD were commonly described in this case study; infrastructure and automation debt.

These studies indicate that causes and effects of TD are not always caused by technical reasons. TD can be the result of intentional decisions made by the different stakeholders. Incurring TD may have short-term positive effects such as time-to-market benefits. However, by not paying down TD can result economic consequences, or quality issues in the long-run. The allowance of TD can facilitate product development for a period, but decreases the product maintainability in the long-term. However, there are some times where short-term benefits overweight long-term costs [27].

2.1.4 Identification of Technical Debt

TD accumulation may cause an increase in maintenance and evolution costs. At worst, it may even cancel out projects. The first step towards managing TD is to properly identify and visualize TD items.

According to Zazworka et al. [31], there are four main techniques for identifying TD in source code: *modularity violations*, *design patterns and grime buildup*, *code smells*, and *automatic static analysis issues* (ASA).

Modularity Violation

Software modularity determines software quality in terms of evolveability, changeability, and maintainability [32], and the essence is to allow modules to evolve independently. However, in reality, two software components may change together though belonging to distinct modules, due to unwanted side effects caused by '*quick and dirty*' solutions [31, 33]. This causes a violation in the software designed modular structure, which is called a modularity violation. Wong et al. [33] identified 231 modularity violations from 490 modification requests in their experiment using Hadoop. 152 of the 490 identified violations were confirmed by the fact that they were either addressed in later versions of Hadoop, or recognized as problems by the developers. In addition, they identified 399

modularity violation from 3458 modification request of Eclipse JDT [33]. Among these violations, 161 were confirmed. Zazworka et al. [31] revealed that the average number of modularity violations per class in release 0.2.0 to release 0.14.0 of Hadoop ranged from 0.04 to 0.11. They identified 8 modularity violations in the first release of Hadoop and 37 in the last one. In addition, they revealed that modularity violations are strongly related to classes with high defect- and change-proneness.

Design Pattern and Grime Buildup

Patterns are known to be general solutions to recurrent design problems. They are commonly used to improve maintainability, flexibility, and architecture design of software systems by reducing the number of defects and faults. Twenty-three design patterns are widely used in software development. They are classified into three types of patterns [34]: *creational*, *behavioral*, and *structural*. *Creational design patterns* are all about class instantiation. This pattern can be divided into class-creation patterns and object-creational patterns. *Creational patterns* include patterns such as Abstract Factory, Builder, Factory Method, Object Pool, Prototype, and Singleton [34]. *Structural design patterns* are all about Class and Object composition. They include patterns such as Adapter, Bridge, Composite, Decorator, Facade, Flyweight, Private Class Data, and Proxy [34]. Lastly, *Behavioral design patterns* are all about class's objects communication, and are concerned with communication between objects. Behavioral patterns include patterns such as Chain of Responsibility, Command, Interpreter, Iterator, Mediator, Memento, Null Object, Observer, State, Strategy, Template Method, and Visitor [34].

Software continuously evolve in response to external demands for new functionality. One consequence of such evolution is software design decay. Izurieta et al. [12] define decay as the deterioration of the internal structure of system designs. Design pattern decay is defined as deterioration of the structural integrity of a design pattern realization. As a pattern realization evolves, its structure and behavior tend to deviate from its original intent. Changes in the code base could potentially lead to code ending up outside the pattern. This is known as design grime, non-pattern related code [12]. Moreover, design patterns can also rot, when changes break the structural and functional integrity of a pattern [12]. Izurieta et al. [35] examined the extent to which software designs actually decay, rot, and accumulate grime in three open-source systems. They found no evidence of design patterns rot, but they found evidence of pattern decay due to grime. Izurieta et al. [36] also studied the consequences of grime buildup on testability of general purpose design patterns. Their result show that as systems age, the growth of grime and the appearance of anti-patterns increase testing requirements.

Code Smell

Some forms of TD accumulate over time in the form of source code [31]. Fowler et al. [37] describe the concept of code smells as choices in **Object-Oriented** (OO) systems that do not comply with the principles of good OO-design and programming practices.

Table 2.2: Code Smell Taxonomy

Code Smell	Group
Long Method	Bloaters
Large Class	Bloaters
Primitive Obsession	Bloaters
Long Parameter List	Bloaters
Data Clumps	Bloaters
Switch Statements	OO-Abusers
Temporary Field	OO-Abusers
Refused Bequest	OO-Abusers
Alternative Classes with Different Interfaces	OO-Abusers
Parallel Inheritance Hierarchies	OO-Abusers
Divergent Change	Change Preventers
Shotgun Surgery	Change Preventers
Lazy Class	Dispensables
Data Class	Dispensables
Duplicated Code	Dispensables
Speculative Generality	Dispensables
Message Chains	Encapsulators
Middle Man	Encapsulators
Feature Envy	Couplers
Inappropriate Intimacy	Couplers
Comments	Other
Incomplete Library Class	Other

Code smell indicates that parts of the design that is inappropriate and that it can be improved. Code smells are usually removed by performing one or more refactoring [37]. For instance, one such smell is *Long Method*, a method with too many lines of code. This type of code smell can be refactored by *Extract Method*, by reducing the length of the method body [37].

Mäntylä et al. [38] proposed a taxonomy based on the criteria on code smells defined by Fowler et al. [37]. The taxonomy categories code smells into seven groups of problems: *bloaters*, *OO-abusers*, *change preventers*, *dispensables*, *encapsulators*, *couplers*, and *others*. The first class, *Bloaters*, represent large pieces of code that cannot be effectively handled. *OO-abusers* are related to cases where a solution does not exploit the the possibilities of OO-design. *Change preventers* refer to code structure that considerably hinder the modification of software. *Dispensables* represent code structure with no value. *Encapsulators* deal with data communication mechanism or encapsulation. *Couplers* refer to classes with high coupling. The last group of code smell problem is *Other*, which refer to code smells that do not fit into any of the other categories. This includes "Incomplete Library Class" and "Comments". Table 2.2 lists all the code smells that are presented by Fowler et al. [37].

Several studies have been conducted to investigate the relationship between code smell and change-proneness of classes in OO source code. A study by Olbrich et al. [39] showed that different phases during evolution of code smells could be identified, and classes infected with code smells have a higher change frequency; such classes seem to need more maintenance than non-infected classes. Khomh et al. [40] investigated if classes with code smells are more change-prone than classes without smells. After studying 9 releases of Azureus and 13 releases of Eclipse, their findings show that classes with code smells are more change-prone than others.

Multiple approaches have been proposed for identifying code smells, ranging from manual approaches to automatic. Manual detection of code smells can be done by code inspections [41]. Travassos et al. [41] present a set of reading techniques that gives specific and practical guidance for identifying defects in OO-design. However, Marinescu [42] argue that manual code inspection can be time expensive, unrepeatable, and non-scalable. In addition, it is often unclear what exactly to search for when inspecting code [43]. Mäntylä [38] revealed more issues regarding manual inspection of code. He states that manual code inspection is hard to conduct due to conflicting perceptions of code smells among the developers, which could cause a lack of uniformity in the smell evaluation.

Automatic approaches for identifying code smells reduce the effort of browsing through large amounts of code during code inspection processes. Ciupke [43] propose an approach for detecting code smells in OO-systems. In this approach, code smells to be identified are specified as a query. The result of a queries is a piece of design specifying the location of the code smell in the source code. This approach was applied to several case studies, both in academical and industrial context. Their findings revealed that code smell detection can be automated to a large degree, and that the technique can be effectively applied to real-world code.

Another method for automatic detection of code smells is done by using metrics. Marinescu [42] proposed a general metric-based approach to identify code smells. Instead of a purely manual approach, the use code metrics were proposed for detecting design flaws in OO-systems. This approach were later refined, with the introduction of detection strategies [44]. Based on their case study, the precision of automatic detection of code smells is reported to be 70%. Furthermore, a study by Schumacher et al. [45] investigated how human elicitation of TD by detecting god class code smells compares to automatic approaches by using a detection strategy for god classes. Their findings show that humans are able to detect code smells in an effective way if provided with a suitable process. Moreover, the findings revealed that the automatic approach yield high recall and precision in this context.

Automatic Static Analysis Issues

The identification of software design and code issues can be done with ASA code tools. ASA tools look for violations of recommend programming practices on source code line level that might cause faults or degrade some parts of software quality [31]. Tools

are able to alert software developers of potential problems in the source code, and they may suggest refactoring solutions to avoid future problems. Several tools for detecting code and design issues have been proposed in the literature [46]. These tools have been applied by several researchers [31, 46–48], and the overall finding is that a small set of automatic static analysis issues are related to defects in the software. However, the set of issues depends on the context and type of software analyzed. Despite the fact that a tool indicates a potential problem in the source code, it takes human judgment to determine if something could be problematic down the road.

2.1.5 Strategies and Practices for Managing Technical Debt

Managing TD comprises the actions of identifying the debt and making decisions about which debt that should be repaid [8, 10, 37]. There are three high-level steps that are required to manage TD [49]: *Increasing awareness of TD*, *detecting and repaying TD*, and *prevent accumulation of TD*. *Increasing awareness of TD* is the first step towards TD management. This includes awareness of the concept of TD, its different forms, the impact of TD, and the factors that contribute to TD. Awareness will help an organization taking the right decisions to achieve their goals. The second step is *detecting and repaying TD*. This step is focused on determining the extent of TD in the software. Identifying specific instances of debt and their impact can help us to prepare a plan to recover from the debt. The final step, *prevent accumulation of TD*, ensures that TD does not increase and remains manageable in the future. Best practices such as refactoring, re-engineering, and testing is necessary to manage the debt [5]. This step also involves collaboration between all stakeholders to collectively track and monitor the debt. For example, using a backlog during release and iteration planning to list debt-related tasks [10]. Guo et al. [20] suggest the use of a portfolio for TD management. This approach collects TD items to a “*TD List*”, which could be used to pay TD based on its cost and value. Tools may also be used for TD management. *SonarQube* is an open source application for quality management [50]. It can assess the debt in a software project by performing automatic static code analysis. Each automatic static analysis issue is assigned a score based on how much work it requires to fix that error. The analysis gives the total sum of TD for the entire product.

2.2 Design Debt

DD is one of the instances of TD [4, 6, 8]. It is concerned with the design aspects of TD. DD includes design smells and violations of design rules [49]. Design smells are certain structures in the design that indicate violation of fundamental design principles and negatively impact design quality [49]. A possible symptom of DD is when code structures drift away from good OO-design principles [29]. There are various reasons for software projects to run into DD. For example, a common OO-design principle tells that a class should have a single purpose and should not implement many functions of the system. A class that tries to accomplish too many purposes is known as *God Class*

code smell. As we have mentioned earlier, code smells are an example of design flaws in OO-design, which may potentially lead to maintainability issues in future evolution of the software system [39].

Suryanarayana et al. [49] present six common causes of DD: *violation of design principles*, *inappropriate use of patterns*, *language limitations*, *procedural thinking in OO-paradigm*, *viscosity*, and *non-adherence to best practices and processes*. Design principles provide guidance to designers in creating high quality software. *Violation of design principles* are manifested as smells. For example, code smells are symptoms of design violations in the source code. *Inappropriate use of patterns* is related to the use of patterns without fully understanding the context. In some cases, misuse of a design pattern may potentially lead to an antipattern. *Language limitations* is related to deficiencies in programming languages. *Procedural thinking in OO-paradigm* is related to using procedural programming techniques in OO-paradigm. For example, using imperative names for classes, functional decomposition, and missing polymorphism with explicit type checks results in design smells in an OO-context. *Viscosity* is one of the reasons to use hacks instead of adopting a systematic process to achieve a particular requirement. Lastly, *Non-adherence to practices and processes* is related to best practices and processes that are not followed correctly or completely.

2.3 Software Quality

According to the IEEE Standard Glossary of Software Engineering Terminology [51], the quality of a software is defined as 1) *the degree to which a system, component, or process meets specified requirements*, and 2) *the degree to which a system, component, or process meets customer or user needs or expectations*. ISO/IEC 9126:2001 offer a valuable conceptual framework for software quality, where a distinction between *quality in use*, *external quality*, and *internal quality* is made [52]. *Quality in use* refers to the users view of system quality. *External quality* reflects the dynamic aspect of a software application, and is subdivided into six quality attributes. *Internal quality* refers to the criteria of each quality attribute. Software quality can be measured by using a product quality model. ISO/IEC 9126:2001 classifies software quality in a structured set of six quality attributes [52]. Bass et al. [53] describes these characteristics as software quality attributes. Table 2.3 present the ISO/IEC 9126:2001 quality model, with the quality attributes, and their criteria.

2.4 Object-Oriented Metrics

OO-metrics have been proposed as a quality indicator for OO-software systems. There are three traditional metrics that are widely used, and are well understood by researchers and practitioners [54]. These metrics are: *Cyclomatic Complexity* (CC), *Size*, and *Comment Percentage*. CC evaluates the complexity of an algorithm in a method. It is recommended that the CC for a method should be below 10 [54]. *Size* of a method is used

evaluate the understandability of the code. *Size* is measured in many different ways, including all physical lines of code, lines of statements, and number of blank lines. *Comment percentage* measures the number of comments in percent by counting total number of comments divided on total lines of code minus number of blank lines.

Chidamber and Kemerer [55] proposed a set of six software metrics to identify certain design traits of a software component. These metrics are: *Weighted Methods per Class (WMC)*, *Depth of Inheritance Tree (DIT)*, *Number of Children (NOC)*, *Lack of Cohesion in Methods (LCOM)*, *Coupling Between Objects (CBO)*, and *Response For a Class (RFC)*. The *WMC* is used to count the number of methods in a class, or to count the of sum of complexities of all methods in a class. The complexity of a method is measured by the CC. This metric measures understandability, maintainability, and reusability [54]. The *DIT* metric measures the maximum number of steps from a class node to the root in the inheritance hierarchy. The deeper a class is within the hierarchy, the greater number of methods it is likely to inherit, making it more complex to predict its behavior [54]. *DIT* metric is related to efficiency, reusability, understandability, and testability [54]. *NOC* metric measures the number of subclasses of a class in a hierarchy. The greater number of children may be an indication of misuse of subclasses, or improper parent abstraction. This metric evaluates efficiency, reusability, and testability [54]. The *LCOM* is used to measure the lack of cohesion in methods of a class. It measures the dissimilarity of methods in a class by looking at the instance variables or attributes used by the methods. High cohesion indicates good subdivision, while low cohesion increase the complexity of a class. This metric measures efficiency and reusability [54]. The *CBO* metric counts the number of other classes to which a class is coupled. Excessive coupling is detrimental to modular design and prevents reuse [54]. Larger number of coupled objects indicates higher sensitivity to changes in other parts of the design. *CBO* evaluates efficiency and reusability [54]. The *RFC* metric counts the total number of methods in a class that can be invoked in a response to a message sent to an object. This metric includes all methods accessible within the class hierarchy. *RFC* metric evaluates understandability, maintainability, and testability [54]. Basili et al. [56] investigated Chidamber and Kemerer's suite of metrics. Their result claim that several of Chidamber and Kemerer's OO-metrics appear to be useful to predict fault-prone classes during early phases of the software life-cycle. Similarly, Li et al. [57] conclude that OO-metrics are able to predict maintenance effort.

2.5 Software Evolution and Maintenance

Increasingly, more and more software developers are employed to maintain and evolve existing systems instead of developing new systems from scratch [58]. Lehman [59] introduced the study of software evolution. Software evolution is a process that usually takes place when the initial development of a software project is done and was successful [60]. The goal of software evolution is to incorporate new user requirements in the application, and adapt it to the existing application. Software evolution is important because it takes up to 85-90% of organizational software costs [58]. In addition, software

evolution is important because technology tend to change rapidly.

Software maintenance is defined as *modifications of a software after delivery to correct faults, to improve performance or other attributes, or to adapt the product to a modified environment* [61]. Maintenance can be classified into four types [60,61]: *adaptive maintenance*, *perfective maintenance*, *corrective maintenance*, and *preventive maintenance*. *Adaptive maintenance* is the modification of a software product performed after delivery to keep the computer program usable in a changed or changing environment. *Perfective maintenance* is the modification of a software product after delivery to improve performance and maintainability. *Corrective maintenance* is the reactive modification of a software product performed after delivery to correct discovered faults. Lastly, *Preventive maintenance* is the maintenance performed for the purpose of preventing problems before they occur.

Van Vliet [62] states that the real maintenance activity is corrective maintenance. 50% of the total software maintenance is spent on perfective, 25% on adaptive maintenance, and 4% on preventive maintenance. This leads to that 21% of the total maintenance activity is corrective maintenance, the 'real' maintenance [62]. This has not changed since the 1980s when Lientz and Swanson conducted a study on software maintenance [63]. The study points out that most severe maintenance problems were caused by poor documentation, demands from users for changes, poor meeting scheduled, and problems training new hires.

2.6 Software Reuse

Software reuse is the process of creating software systems from existing software rather than building software systems from scratch [64,65]. There are many various techniques that can be used to achieve software reuse: *abstraction*, *compositional reuse*, *generative reuse*, and *generation vs. composition* [66]. *Compositional reuse* is based on the idea of reusable components. It supports bottom-up development of systems where low-level components are available. *Compositional reuse* is also known as ad-hoc reuse or individual reuse. *Generative reuse* is based on the reuse of a generation process. It is often domain-specific, and is concerned with artifacts such as requirements, designs, and sub-systems [67]. *Generative reuse* is also known as systematic reuse [67]. *Generation vs. composition* is a combined approach.

There are several "assets" from a software project that can be reused, including *architectures*, *source code*, *data*, *designs*, *documentation*, *estimates*, *human interfaces*, *plans*, *requirements*, and *test cases* [65].

The most common benefits of software reuse are quality improvement and effort reduction. Software reuse includes lower costs, shorter development time, higher quality of reusable components, and higher productivity [68,69]. The quality of the software increases every time an asset is reused, because errors are discovered more frequently, thus making it easier to keep the artifact more stable [66]. However, software reuse may cause

problems as well. A case study on selected feature from self-driving miniature car development revealed that reuse of legacy, third party, or open-source code, was one of the root causes of TD accumulation [70]. Furthermore, Morisio et al. [71] identified three main causes of failures associated with software reuse: *not introducing reuse-specific processes*, *not modifying non-reuse processes*, and *not considering human factors*. These causes were combined with lack of commitment by top management. Organizations attempting to implement generative reuse face both technical and non-technical problems [72]. Technical problems includes challenges with tools, standards and technology. Non-technical problems include the human factor, cultural issues, economical issues, and organizational issues.

2.7 Refactoring

Fowler defined refactoring as a process of changing the software system in such way that it does not alter the external behavior of the code yet improves its internal structure [37]. Refactoring increases the code readability and maintainability by cleaning up code in a way such that changes of defects is reduced [37]. Refactoring is an act of improving the design and quality of an existing system [62]. It is believed that refactoring improves software quality and developer productivity by making it easier to understand and maintain software source code [73]. According to Opdyke [74], refactoring can be categorized into low-level and high-level refactoring. Low-level refactoring is related to changing a program entity, such as renaming a member variable. High-level refactoring is usually a sequence of low-level refactorings, such as defining an abstract superclass. In addition, Opdyke presents twenty-six low-level refactorings and three high-level refactorings. Fowler et al. [37] have summarized the work of Opdyke by cataloging his refactorings along with many other refactoring techniques that have been collected over the years. Refactoring can be applied to three types of software artifacts: *programs*, *designs*, and *software requirements*. *Programs* include refactoring at source code or program level. For example, *Extract Method* [37] is a refactoring technique that can be applied at this level. *Designs* include refactoring at design level, for example in the form of UML models. Design pattern is an example of an artifact at this level. Lastly, *software requirements* include refactoring at the level of requirements specification. For example, decomposing requirements into a structure of viewpoints.

2.8 Embedded Systems

The *IEEE Standard Glossary of Software Engineering Terminologies* define embedded system as *a computer system that is part of a larger system and performs some of the requirements of that system; for example, a computer system used in aircraft or rapid transit system* [51]. Embedded systems are special-purpose computer systems designed to perform certain dedicated functions under certain constraints. According to Koopman [75], there are four types of embedded systems: *general computing*, *control systems*,

signal processing, and *communication and networking*. The various types of embedded systems share common requirements such as: *real-time requirements*, *resource consumption*, *dependability*, and *life-cycle properties* [76]. Embedded systems are supposed to be failure-free [77], but these requirements may hinder embedded systems to deliver reliable service given a disturbance to its services [78].

Software constitutes only one part in embedded systems. Embedded software is defined as software part of a larger system which performs some of the requirements of that system [51]. Developing embedded software has proven to be difficult [79, 80]. Embedded software not only have to meet the functional requirements of "what", but also the non-functionality or "quality" requirements expected of them [80]. Non-functional requirements have come to be referred to as software quality attributes. Ebert et al. [79] state that quality of software in embedded systems is difficult to measure. Most behaviors for embedded systems are related to run-time quality attributes, depending on the system domain. For instance, reliability and safety qualities are important in the automotive domain [81], while performance and reliability is important in mission critical systems [82]. Other quality attributes such as maintainability and usability, are often compromised in favor of run-time quality attributes.

2.8.1 Safety-Critical Systems

Safety-critical systems are those systems whose failure could result in loss of life, significant property damage, or damage to the environment [83]. These systems put a huge demand on software reliability, because a minor error in safety-critical software can produce failure of a complete system [79, 81, 82]. Consider the role of safety-critical systems is often critical, the importance of managing software quality is therefore necessary to deliver software in a useful, safe, and reliable way. Given the lack of focus on non-functional requirements, the amount of TD in safety-critical systems grow continuously. Insufficient requirements and testing are two major cost drivers in embedded software development [1, 79]. 40% of all software defects results from insufficient requirements, especially non-functional requirements[6,7,45] [1, 79, 84]. Testing code consumes up to 40% of development resources. In addition, testing may potentially require 15-50% of total project duration [79].

Table 2.3: Software quality attributes, criteria, and description (ISO/IEC 9126)

Quality Attribute	Criteria	Description
Functionality		Systems ability to do work for which it was intended.
	Suitability	
	Accuracy	
	Interoperability	
	Security	
	Compliance	
Reliability		Systems ability to keep operating over time under certain conditions.
	Maturity	
	Fault Tolerance	
	Recoverability	
	Compliance	
Usability		Systems ability to be understood, learned, and used by others.
	Understandability	
	Learnability	
	Operability	
	Attractiveness	
	Compliance	
Efficiency		Systems ability to provide appropriate performance relative to the amount of resources used under stated conditions.
	Time behaviour	
	Resource utilization	
	Compliance	
Maintainability		Systems ability to be modified in the future.
	Analyzeability	
	Changeability	
	Stability	
	Testability	
	Compliance	
Portability		Systems ability to be transferred from one environment to another.
	Adaptability	
	Installability	
	Co-existence	
	Replaceability	
	Compliance	

CHAPTER 3

RESEARCH METHODOLOGY

The nature of this thesis makes it suitable as an empirical research. To answer the research questions that was stated in Section 1.3, Chapter 1, an empirical research needs to be carried out in order to collect some data. This chapter provides a brief introduction to research methods in software engineering, and describes the research conducted in the thesis. Section 3.1 describes the relevant research methods in software engineering. Section 3.2 describes the research method that was chosen for this study. Section 3.4 presents the research process we have followed throughout this thesis, which includes our research design.

3.1 Research Methods in Software Engineering

This section presents the relevant research methods applied in software engineering research.

Research is believed to be the most effective way of coming to know what is happening in the world [16]. Empirical software engineering is a field of research based on empirical studies to derive knowledge from an actual experience rather than from theory or belief [85]. Empirical studies can be explanatory, descriptive, or exploratory [15].

There are two types of research paradigms that have different approaches to empirical studies [15]; *the qualitative*, and *the quantitative* paradigm. *Qualitative* research is concerned with studying objects in their natural setting [15]. It is based on non-numeric data found in sources as interview tapes, documents, or developers' model. *Quantitative* research is concerned with quantifying a relationship or to compare two or more groups [15]. It is based on collecting numerical data.

To perform research in software, it might be to understand the different research strategies that are available in software engineering. Oates [14] presents six different research strategies; *survey*, *design and creation*, *case study*, *experimentation*, *action research*, and *ethnography*.

Survey focuses on collection data from a sample of individuals through their responses to questions. The primary means of gathering qualitative or quantitative data are interviews or questionnaires. The results are then analyzed using patterns to derive descriptive, exploratory, and explanatory conclusions.

Design and creation focuses on developing new IT products, or artifacts. It can be a computer-based system, new model, or a new method.

Case study focuses on monitoring one single 'thing'; an organization, a project, an information system, or a software developer. The goal is to obtain rich, and detailed data.

Experimentation are normally done in laboratory environment, which provides a high level of control. The goal is to investigate cause and effect relationships, testing hypotheses, and to prove or disprove the link between a factor and an observed outcome.

Action research focuses on solving a real-world problem while reflecting on the learning outcomes.

Ethnography is used to understand culture and ways of seeing of a particular group of people. The researcher spends time in the field by participating rather than observing.

3.2 Choice of Research Method in this Thesis

The main purpose of this research project is to gain an understanding about the nature of TD in software design, and its potential sources in safety-critical system software in order to improve the management of software evolution. Based on the research questions stated in Chapter 1, we have applied the case study approach to our research as it fits our desire to study an industrial system. Case studies provide both quantitative and qualitative information about the system [14], depending on the approach the case study is taking.

3.2.1 Case Study Method

Case study is an empirical method to investigate a single phenomenon within a specific time space in real-life context [15]. Case studies excels at bringing an understanding of why or how certain phenomena occur or to add strength to what is already known through previous research [15, 17]. Runeson et al. [86] suggests case study as the most appropriate research method to use when exploring how a problem behaves in a real life context. In addition, they conclude that case study is suitable for software engineering research. There has been suggested systematic approaches for organizing and conducting a research successfully [17, 86]. According to Yin [87], a research design is an action plan from getting "here" to "there", where "here" is defined as the initial set of questions

answered, and "there" is some set of conclusions about these questions. Moreover, a research design can be seen as a blueprint of research, dealing with at least four problems: *what questions to study, what data are relevant, what data to collect, and how to analyze the results* [87]. Soy [17] proposes six steps that can be used when carrying out a case study:

1. *Determine and Define the Research Questions*: The first step involves establishing a research focus by forming questions about the problem to be studied. The researcher can refer to the research focus and questions over the course of study.
2. *Select the Cases and Determine Data Gathering and Analysis Techniques*: The second step involves determining what approaches to use in selecting single or multiple real-life cases to examine, and which instruments and data gathering approaches to use. (whom we want to study, the case, cases, sample. and how we want to study it, design).
3. *Prepare to Collect Data*: The third step involves a systematic organization of the data to be analyzed. This is to prevent the researcher from being overwhelmed by the amount of data and to prevent the researcher from losing sight of the research focus and questions.
4. *Collect Data in the Field*: This step involves collecting, categorizing, and storing multiple sources of data systematically so it can be referenced and sorted. This makes the data readily available for subsequent reinterpretation.
5. *Evaluate and Analyze the Data*: The fifth step involves examining the raw data in order to find any connections between the research object and the outcomes with reference to the original research questions.
6. *Prepare the Report*: In the final step, the researcher reports the data by transforming the problem into one that can be understood. The goal of the written report is to allow the reader to understand, question, and examine the study.

3.3 Case Context

We have chosen to study a commercial system by conducting a case study to study the consequences of DD in this thesis. The conducted case study took place at Autronica Fire and Security AS. Autronica is a leading innovator, manufacturer, and supplier of fire safety equipment and marine safety monitoring and surveillance equipment. Their headquarter is based in Trondheim, Norway. AutoSafe, a high-end distributed fire alarm system, is one of the products they offer. The product was first released around year 2000, and has been on sale since. The software of the product is mainly based on C/C++ source files. Project "Firmus" is the project name for the next generation AutoSafe. "Firmus" is a Latin word, which in English means: *solid, firm, strong, steadfast, steady, stable, reliable, and powerful*. The goal with Project "Firmus" is to adopt newer technologies and technology standards that are used today. We had the opportunity to conduct our case study on Project "Firmus". This project is still in the development phase. The goal

Table 3.1: System Metrics for Project "Firmus"

Project "Firmus"	
Lines	88465
Lines of Code	49287
Lines of Comments	23017
Components	13
Files	461
Number of Classes	339

of this research is to identify DD before it potentially creates more maintenance effort in the future. As we mentioned earlier, the product is mainly based of C/C++ source code files. The software architecture of the Project "Firmus" is component-based, where the different source files are divided into different components. The system consists of 13 components, and 461 source code files in total. Test files are found inside each component. To ensure reuse of source code, the company have developed a library that is used by this system and other systems as well as they can be seen as a case itself. However, we have decided not to include library files in our analysis. Table 3.1 summarizes the system metrics, which includes the test files.

3.4 Research Process

A research process provides a systematic approach on how to fulfill the goal of a research. We have chosen to follow the principles of the six steps defined by Soy [17] in this study: 1) *Determine and Define The Research Questions*, 2) *Select the Cases and Determine Data Gathering and Analysis Techniques*, 3) *Prepare to Collect Data*, 4) *Data Collection*, 5) *Evaluate and Analyze Data*, and 6) *Prepare the Report*.

3.4.1 Determine and Define the Research Questions

The first step of this study is to define the research goal, and the research questions. In our previous research [13], we stated that we are interested in getting a deeper insight into the field of TD by studying an industrial system. There are many subcategories of TD (See Table 2.1. With regards to that, we have chosen to investigate DD in safety-critical systems.

An analysis of the state-of-the-art was carried out to determine what prior studies have determine about the topic of DD. The goal with the analysis is to determine and define the research questions. *Google Scholar*, *ACM Digital Library*, *Scopus*, and *IEEE Xplore Digital Library* were tremendously used during the analysis to find research papers that are relevant for our thesis. The literature review showed that there have been very little prior research on DD in safety-critical systems. Research questions were defined based on the results from our literature review. We have defined four research questions,

Table 3.2: Research Questions

RQ1	How can DD be identified?
RQ2	What are the effects of DD?
RQ3	What kind of DD can be found in embedded systems?
RQ4	How to pay DD?

RQ1-4, which we have summarized in Table 3.2. These research questions will be our primarily driving force though this research.

3.4.2 Select the Cases and Determine Data Gathering and Analysis Techniques

A representative context has been chosen to analyze and investigate the research questions. We have chosen to conduct an exploratory and descriptive case study in real-life context to obtain knowledge about the problem to be studied. Several e-mails were sent out to organizations within the embedded industry. After a few weeks, we managed to get in touch with Autronica Fire and Security. The case study took place at their main office in Trondheim for approximately six weeks. A brief description of Autronica can be found in Section 3.3. We were provided with a workspace and multiple data sources, including access to the software's source code, issue for the project, system requirements, and documentation for system design and code. Data were mainly extracted from these sources.

The first step is to find the structural code and design attributes of the software system. OO-metrics can be used to assess the quality of the software. Therefore, design attributes can be identified by measuring the OO-metrics. Furthermore, as we mentioned in Section 2.2, code smell detection can be used to identify design violations in the source code. A part of the literature review was to get familiar with existing tools that have been used to address similar problems. There is a wide set of tools which enables software metrics measurement and automatic static analysis of code. Some of the tools are open-source, while others contain strict licensing. However, a few of the commercial tools provide license for academical purposes. We will only focus on the tools providing measurement for C/C++. We have listed the tools used in this thesis below.

Doxygen [88] is a free software for generating documentation from annotated C++ sources. *Doxygen* has the ability to generate documentation in HTML or in LaTeX. Since the documentation is extracted from the source code, it is easier to keep the documentation up to date. In addition, *Doxygen* can be configured to extract the code structure from undocumented sources files, which makes it possible to visualize the relations between various elements in the software. *Doxygen* is used by the company to keep the documentation up-to-date, and was therefore used in this project to learn about the system and to visualize the system.

ArgoUML [89] is an open source UML modeling tool that includes support for all stan-

dard UML 1.4 diagrams. It features reverse engineering of C++ projects by reading C++ source files and generate and UML model and diagrams. *ArgoUML* was primarily used to reverse-engineer the source code to extract class diagrams.

Enterprise Architect [90] is an UML modeling tool. It has the ability to produce UML diagrams from code. This tool was used to create class diagrams for each component, allowing us to identify possible code smells, such as *Large Class* and *Long Method* code smell. *Enterprise Architect* is a commercial software, but they offer a 30-day trial version of the software. The reason we used this was because ArgoUML was not fully able to reverse-engineer the source code.

SonarQube [50] is an open source platform for quality management of software code. It has the ability to monitor different types of TD. It supports multiple languages through plugins, including Java, C/C++, JavaScript, and PHP. A downside with *SonarQube* is that some of the plugins requires a commercial license, and that some features included are not applicable for C++. This tool was used to detect code smells in Project "Firmus".

Understand [91] is a code static analysis tool. It supports dozen of languages, including Java, C/C++, Fortran and Python. Understand can help developers analyze, measure, visualize, and maintain source code. It includes many features, including dependency graph visualization of code, and various metrics about the code (e.g., CBO).

CppDepend [92] is a commercial tool for analyzing source code in C/C++.

CppCheck [93] is an open-source static analysis tool for C/C++ source code.

CppClean [94] is another open-source static analysis tool for C/C++ source code. Among many functionalities, this tool supports finding unnecessary '#includes' in header files.

CCCC [95] is able to measure source code related metrics. It is able to measure some of the metrics defined by Chidamber and Kemerer [55]. We have used this tool to measure software metrics, and the results has been compared to the results from *Understand*.

3.4.3 Prepare To Collect Data

Preparing for data collection is the third step of this case study. *Doxygen* is used by the company to generate and keep system documentation up-to-date. We spent some time analyzing the system documentation to get familiar with the system. Moreover, *Doxygen* has the ability to generate various diagrams, including inheritance diagrams, and dependency graphs. However, a downside with *Doxygen* is that it does not allow us to interact with the diagrams. *Doxygen* allows us to specify depth of the graphs that are being generated, but output can be very large, hence we had some troubles understanding the graphs. Furthermore, *Doxygen* can generate a dependency graph for each file in a component, but it does not provide full dependency graphs for a component. There are many tools that offers reverse engineering of C/C++ source code, so we decided to try out a few of them, including *ArgoUML*, *Enterprise Architect*, and *Understand*. *Understand*, *CppClean*, *CppDepend*, and *CppCheck* will be used to extract design problems at code level by analyzing the source code.

A Word document and an Excel spreadsheet were created to keep track of the extracted data so we can review it later for analysis.

Metrics Selection

Various software metrics exist for system measurement. Unfortunately, there are no well known standardized set of software metrics aimed to measure safety-critical systems. Choosing right metrics for measuring safety-critical software is preferable because this kind of software is responsible for many things, such as human lives. To enhance quality in software systems, OO-metrics were established. They measure characteristics of OO-systems in a way to improve them. Many aspects have been defined to improve the quality of code using these metrics. We have mainly used the metric suite defined by Chidamber and Kemerer [55]. This suite of metrics is widely cited and has been used by many researchers, and their metrics suite is the deepest research in OO-metrics investigation. For example, Rosenberg et al. [96] have applied this set of metric in evaluation of many NASA projects. The projects were written in both C++ and Java. In addition to Chidamber and Kemerer's metrics, we have measured two additional metrics, one counting the number of instance methods, and the other counting the number of instance variables. The following metrics have been measured in this study by using *Understand*:

LCOM: A method is cohesive when it performs a single task. Low cohesion increases complexity, and will increase the likelihood for errors during development process. In general, the desirable value for LCOM is to be lower. The average percentage of class methods using a given class instance variable

DIT: Measures the longest path from a given class to the root class in the inheritance hierarchy.

CBO: Measures number of other classes that are coupled to a particular class. The desirable value is lower.

NOC: Measures number of subclasses from a given class.

RFC: It is the sum of number of methods that can be potentially executed in response to a message by an object of a class.

NIM: Measures number of instance methods in a class, which is methods that are accessible through an object of that class.

NIV: Measures number of instance variables in a class, that is, variables that are only accessible through an object of that class. Instance variables are used to measure LCOM of a class.

WMC: This metrics sums the CC of a class by counting the CC for each method.

3.4.4 Data Collection

The fourth step of the research process is to execute the plan that was created in step three. During the case study, data is collected from two different sources by using multiple tools to improve the reliability of the study. The first source of design flaws is to identify for code smells in the source code. Table 2.2 in Section 2.1.4, Chapter 2, summarizes the code smells that are presented by Fowler et al. [37]. By using automatic static analysis tools, we were able to identify multiple code smells in the system. Most of the code smells were manually verified by the researcher by inspecting the class and dependency diagrams for the class in which code smell exists. For instance, *Duplicated Code* code smell were identified using *SonarQube*. We inspected each file with duplicated code to verify the results. Another example of a code smell we identified is the *Long Method* code smell. *Long Method* code smell was identified using *CppDepend* and *Understand*. The results were verified by reverse engineering the source code to generate UML class diagrams. At first, we used the built-in functionality in *Doxygen* to generate the class diagrams. However, *Doxygen* was not able to provide full class diagrams. Therefore, we had to look for other options. We came across ArgoUML, an open source alternative to generate UML diagrams by reverse engineering C/C++ code. After comparing some of the results with snippets from *Doxygen* class diagrams, we noticed that *ArgoUML* failed to reverse engineer parts of code that contains classes, and their corresponding relations to other classes. This led us to look for commercial software. Using *Understand* and *Enterprise Architect*, we were able to extract the class diagrams for the system by using their built-in reverse engineering functionality. *Enterprise Architect* allow us to interact with the class diagram. In addition, it was able to create a class diagram for each component, which made the navigation much easier compared to *Understand*. UML class diagrams were used to verify *Large Class* code smell and *Long Method* code smell.

The second source of design flaw identification is to measure the OO-metrics for the system. OO-metrics are used to manage, predict, and improve the quality of a software product [97] by finding classes with large metric values. *Understand C++* was used for measuring all the metrics values of the system. The reason we chose to use this tool is because it provides wide set of metrics to be measured. In addition, *Scitools*, the developers behind *Understand C++*, offer a 15 days trial for academical purposes. The metrics for C++ are divided into four groups: file, class, project, and method. We have mainly focused on extracting class metrics, which includes *LCOM*, *DIT*, *CBO*, *NIV*, *NIM*, *NOC*, and *WMC*. In our measurements, we decided to exclude the test classes as they may potentially affect the metrics in both positive and negative way. After analyzing our project, we had the opportunity to extract the measured in HTML format. However, we wanted to export the tables to *Microsoft Excel 2016* to perform some calculations. We had to create a script to convert the HTML tables to Excel tables in order to compute descriptive statistics and graph creation. Descriptive statistics deal with presentation and numerical processing of a data set [15]. After collecting experimental data, descriptive statistics may be used to describe and graphically present interesting aspects of the data set. The goal with descriptive statistics is to get a feeling for how the data set is distributed, and may be used to understand the nature of the data and identify abnormal or false data points.

The HTML tables were imported to both Excel 2013 and Google Spreadsheet. After metrics extraction, descriptive statistics were computed for the whole project and for each component using formulas in *Google Spreadsheet*. These statistics aims to give a measure of the value of the metrics for all the classes, which we can use to identify classes with weak metric values. In addition, graphs were created using *Google Spreadsheet* and *Microsoft Excel 2016*.

The descriptive statistics that we have measured are:

- **Minimum:** The minimum value of a metric.
Google spreadsheet formula: *MIN(value1; [value2; ...])*
- **Maximum:** The maximum value of a metric.
Google spreadsheet formula: *MAX(value1; [value2; ...])*
- **Sample Mean:** The mean of the metric, that is, the average value of a metrics. It can be used to measure the center of the data.
Google spreadsheet formula: *AVERAGE(value1; [value2; ...])*
- **Mean:** Represents the middle value of a data set. It is calculated by sorting the samples in ascending order and picking the middle sample [15].
Google spreadsheet formula: *MEDIAN(value1; [value2; ...])*
- **Standard Deviation:** A measure of how spread out the numbers are. Higher values indicates greater spread.
Google spreadsheet formula: *STDEV(value1; [value2; ...])*
- **Kurtosis:** A measure of whether the data set have a peak or not. A positive kurtosis value indicates peaked data distribution, while negative value indicate flat data distribution.
Google spreadsheet formula: *KURT(value1; [value2; ...])*
- **Skewness:** A measure which gives information about the symmetry of data distribution. Normal distribution has a skewness equal zero. Positive skewness value indicate longer right tail with sample mean placed on the right side of the peak value. Negative skewness value indicate longer left tail with sample mean placed on the left side of the peak value.
Google spreadsheet formula: *SKEW(value1; [value2; ...])*

Code smell data were collected using automatic static analysis tools. *CppDepend*, *Cp-pClean*, *CppCheck*, *SonarQube*, and *Understand* were mainly used to detect the code smells. In addition, we have manually verified some of the detected code smells in order to determine whether the hits are false negative or false positive. For example, to verify Long Method code smell, we manually inspected the class in which the method is located.

3.4.5 Evaluate and Analyze the Data

In the fifth step of the case study, we examined the data collected that was in the fourth step. Although OO-metrics are well-known, there is still a discussion about identifying their threshold and usage [98]. In the previous step, we collected data from multiple OO-metrics. The data can be used to measure the software quality from different points of views, search for potential problems in the source code, and identify possible candidates in which inspection may potentially be needed. However, the usage of OO-metrics alone do not tell us which classes have acceptable metric values, and which classes need an inspection. In addition to OO-metric measurement, we also need to determine when the metric values represent a positive characteristic, and when the values becomes some sort of warning sign for possible unwanted aspect of the software. The boundaries between two sets of values are known as thresholds [99]. Our goal with the analysis part is to identify the classes with metric values above a defined threshold.

Threshold values were identified using statistical information method [100]. For each metric, we used descriptive statistics, the sample mean and the standard deviation, to derive its threshold values. There are three different threshold values we identified. The first value corresponds with the sample mean and represents the most typical value in the data set. The second value is calculated as the sum of the sample mean and the standard deviation. It represents high values of the metric, but acceptable in some cases. The last value is the second value multiplied with 1.5 [100]. It represents extreme values and should not be present in the data set.

3.4.6 Prepare the Report

Lastly, the methods conducted in this research will be reported through this thesis. This involves all the steps that we have gone through this research. The report also includes findings from the literature review, and how they are related to our findings from the case study. Lastly, the report conclusion summarizes our contributions, and points out suggestions for further research, so that other researchers may apply these techniques in some other context to determine whether findings are similar to our research or not.

3.5 Summary of the Research Design

The research questions described in Subsection 3.4.1 were answered using the case study methodology and a literature review. The literature review helped us getting familiar with the topic of area, and to identify the tools needed to mine necessary data. A case study was conducted in collaboration with Autronica Fire and Security AS, where we analyzed a safety-critical system to investigate DD.

This chapter presents the results of this study. We have evaluated the design and software quality of the system by measuring OO-metrics related to classes, and by identifying code smells in the system. Section 4.1 presents the OO-metrics we have measured for Project "Firmus". In Subsection 4.1.1, we present the OO-metrics for each component in Project "Firmus". Finally, Section 4.2 summarizes the code smells we have detected using automatic static analysis tools.

4.1 Object-Oriented Metrics in Project "Firmus"

The metrics that have been used to measure the quality of the code is mostly based on the work of Chidamber and Kemerer. [55]. The following metrics were measured: *WMC*, *DIT*, *NOC*, *LCOM*, *RFC*, and *CBO*. In addition to these metrics, we have measured *NIV*, and *NIM* in each class. A short description of each metric is provided in Section 3.4. We present their descriptive statistics which includes the *minimum*, *maximum*, *median*, *sample mean*, *standard deviation*, *kurtosis*, and *skewness* values for the whole system. In addition to metrics for Project "Firmus", we present descriptive statistics for each component in the system.

A total of 317 files were analyzed. These files contains 226 classes, and 31204 lines of code, excluding the test files. Table 4.1 presents descriptive statistics for class level metrics for the whole project.

LCOM: A class is cohesive if its LCOM value is low. In this analysis, LCOM is measured in percent. Our data revealed that LCOM value lies between a range from 0 to 100, indicating that there are classes with high and low cohesion. Figure 4.1 shows the frequency distribution of LCOM values. There are 74 classes with LCOM value of 0. These

Table 4.1: OO-metrics and descriptive statistics for Project "Firmus"

Metric	Min	Max	Median	Sample Mean	Standard Deviation	Kurtosis	Skewness
LCOM	0	100	55	42.711	32.961	-1.501	-0.274
DIT	0	4	1	1.066	1.069	-0.676	0.647
CBO	0	30	5	6.178	5.162	2.093	1.214
NOC	0	20	0	0.467	1.866	61.606	7.064
RFC	0	115	10	15.991	8.769	9.279	2.936
NIM	0	48	7	8.458	7.006	8.117	2.515
NIV	0	18	1	2.195	2.828	5.131	2.003
WMC	0	325	10	20.293	31.934	41.063	5.314

classes have high cohesion. However, 118 classes have a value of LCOM larger than 50. Among them, 7 classes have a value of LCOM larger than 90, whereas 2 classes have a LCOM value of 100. Kurtosis and skewness values tell us that the normal distribution is almost perfectly symmetrical. Classes with low cohesion increase the complexity of the software, and may therefore increase the likelihood of errors during development. It is necessary to split one class to two or more classes to make them more cohesive to improve the class design.

CBO: In general, higher values of CBO indicate fault prone classes. Our analysis show that 194 classes have a value of CBO less than 10, and that 4 classes have a value larger than 20. The maximum value of CBO is 30. This class is an example of a class that is hard to understand, harder to reuse, and more difficult to maintain.

DIT: DIT value appears to be generally low in the captured statistics. A class with DIT value of 0 is the root in a class hierarchy. Figure 4.3 shows that 89 classes have a DIT value of 0, and 65 classes have a DIT value of 1. The median value suggests that most classes tend to be close to the root in the inheritance hierarchy. DIT metric can be used to determine whether the design is top heavy or bottom heavy [55]. A design is top heavy if there are too many classes near the root, and bottom heavy if most classes are near the bottom of the hierarchy. The system appears to be top heavy by observing the empirical data of DIT, which indicates that there may be lack of reuse through inheritance. Classes with a DIT value of 2 and 3 indicates higher degree of reuse through inheritance. Widentified 70 classes in total with DIT value of 2 and 3. The maximum value of DIT is 4. These values show that inheritance is used in most of the classes to an optimal level. However, there may be some possibilities for improvements for classes with DIT value of 0.

NOC: NOC metric measures the number of subclasses of a class. The median value of NOC show that most half of the classes has no immediate children, which indicates that inheritance may not be used enough. According to the Figure 4.4, approximately 86% of the classes seem to have no subclasses, hence affecting the normal distribution due to a high value of kurtosis. In addition, the skewness value indicate that long tail is on the positive side of the peak (i.e., on the right side). Furthermore, the max value of NOC is 20, which may indicate a misuse of subclassing. Classes with high NOC value are difficult to modify, and they usually require more testing because of the effects on changes on all the children.

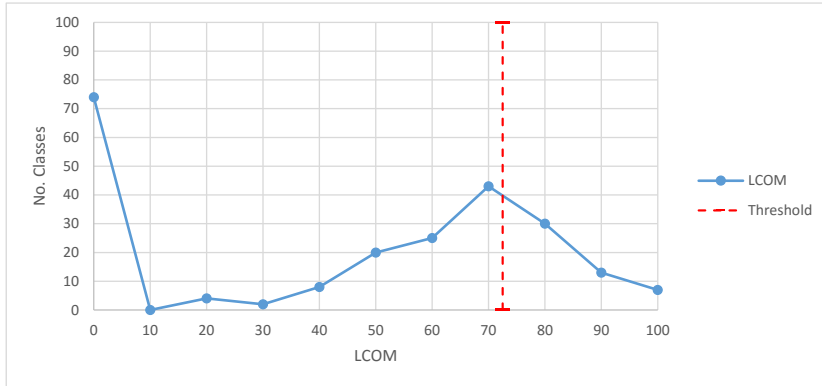


Figure 4.1: Frequency chart of the LCOM metric

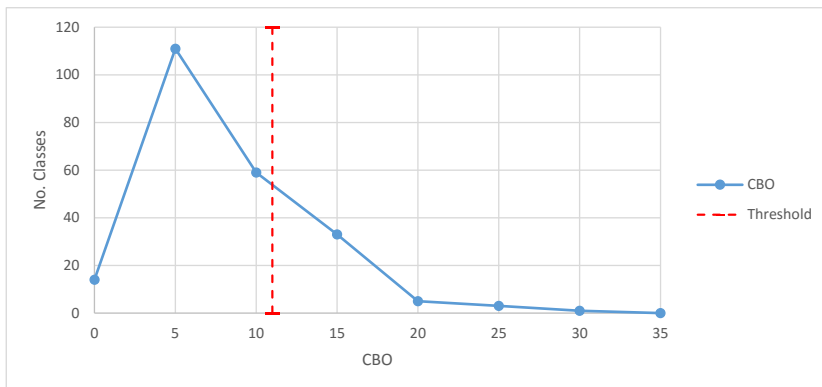


Figure 4.2: Frequency chart of the CBO metric

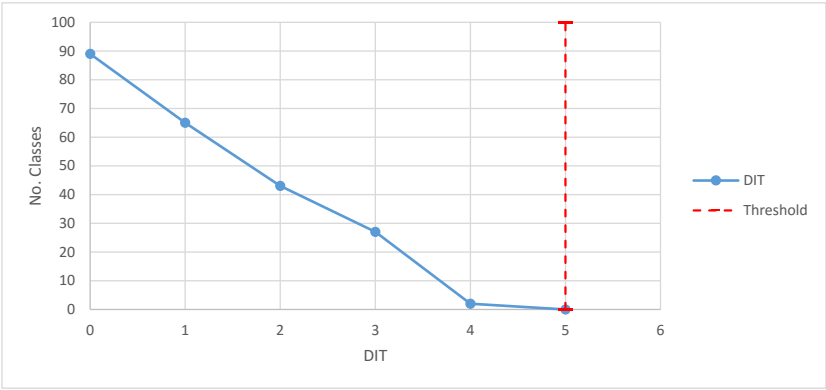


Figure 4.3: Frequency chart of the DIT metric

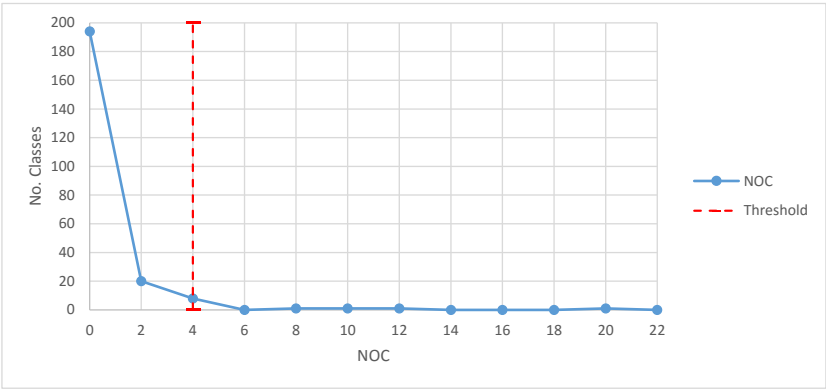


Figure 4.4: Frequency chart of the NOC metric

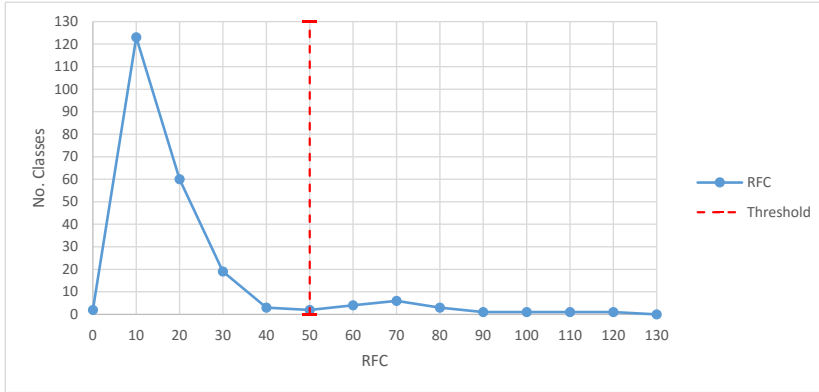


Figure 4.5: Frequency chart of the RFC metric

RFC: Classes with large RFC tend to be complex and have decreased understandability. Testing classes with large RFC is more complicated. The RFC statistics reveal that majority of the classes have a RFC of less than 20. There are only 22 classes with a value of RFC larger than 30, where 2 of them have a value larger than 100. The maximum value of RFC in this system is 115. A suggestion would be to inspect classes with higher values of RFC.

WMC: We observe that majority of the classes have a value of WMC less than 10 by examining Figure 4.6. More precisely, 123 classes has a value of WMC less than 10. Moreover, 5 classes have a value of WMC2 larger than 100. The maximum value of WMC is set to 325.

NIM and NIV: The NIM and NIV metrics report the number instance methods and instance variables in a class. Our analysis show that most classes are small. The sample mean of NIV tells us that each class has an average of 2 instance variables. The maximum value of NIV is 18, indicating that there is at least one class that contains 18 instance variables. The sample mean of NIM show us that each class has an average of 8 instance methods. More precisely, there are 170 classes with value of NIM less than 10. The maximum value of NIM is 48, which indicates that there is at least one class contains 48 methods.

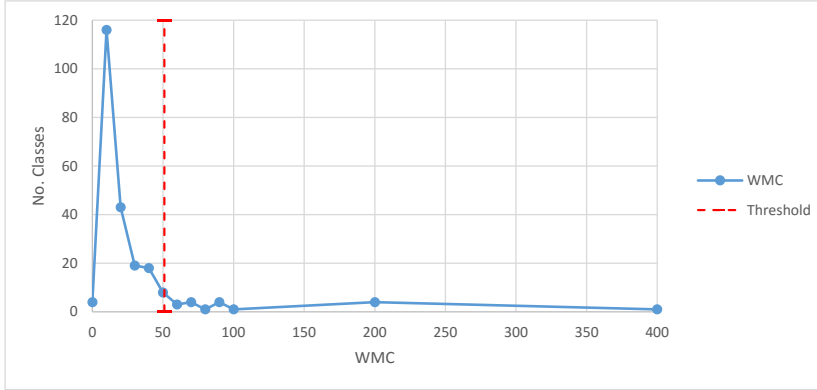


Figure 4.6: Frequency chart of the WMC metric

4.1.1 Object-Oriented Metrics for the Components

Descriptive statistics in Table 4.1 reveal statistics for class level metrics for the whole project. However, the statistics do not say anything about class level metrics in the different components. Some of the components may have good OO-metric values, while other components have bad statistics. To identify components that are more fault-prone, we calculated descriptive statistics for each component.

Component A

Component A contains 53 files. Among these files, we identified 37 classes and 5427 lines of code. Figure 4.7 presents the frequency distribution of the measured OO-metrics. Table 4.2 presents common descriptive statistics of the metric distribution.

The DIT values indicate that inheritance hierarchies is somehow flat. Classes with flat inheritance hierarchy usually hints that reuse through inheritance is not used. There are 8 classes with flat inheritance hierarchy. Rest of the classes inherit for at least one class. The max value captured show that some classes have deep hierarchy. Higher values for DIT indicates higher degree of reuse, but as trade off, it may potentially increase the complexity of the class. Moreover, the results indicate that most classes only have a few subclasses. We identified 29 classes with no children. Moreover, we identified a class with

Table 4.2: OO-metrics and descriptive statistics for Component A

Metric	Min	Max	Median	Sample Mean	Standard Deviation	Kurtosis	Skewness
LCOM	0	94	62	46.405	34.321	-1.477	-0.472
DIT	0	4	2	1.567	1.167	-0.707	0.270
CBO	0	30	6	6.758	6.639	4.00	1.766
NOC	0	8	0	0.758	1.706	8.508	2.743
RFC	6	115	38	43.649	31.515	-0.708	0.559
NIM	4	40	10	13.135	8.891	3.674	1.991
NIV	0	12	1	2.270	2.969	3.694	1.917
WMC	3	194	19	31.081	36.554	10.362	2.821

a NOC value of 8.

The results show that 32.4% (i.e., 12 classes) of all classes in Component A are strongly cohesive, which implies that more than half of the classes show lack of cohesion. We notice that two classes has LCOM values larger than 90 by examining Figure 4.7. These classes indicate loose class structures. The kurtosis value of LCOM is negative, so we can consider LCOM having a flat distribution, which can be seen Figure 4.7.

Most classes have small CBO values, indicating that they are self-contained. However, the frequency distribution shows that few of the classes are strongly coupled. There is one class in component B with a CBO value of 29. This class may potentially be fault-prone class. Additionally, this class may affect its reusability and maintainability. This particular class has LCOM value of 62, WMC value of 95, and RFC value of 25, hence being a possible Large Class code smell.

The results show that each class have at least two methods. More than half of the classes have low RFC values, which indicates greater polymorphism. However, there are few classes in this component that has larger values of RFC. The maximum RFC value is 115, and classes with high RFC are usually difficult to maintain and test.

The values of WMC ranges from 3 to 194. The median value indicate that at least 50% of the classes have a cyclomatic complexity of 17 or less. However, the sample mean is revealed to be larger than the median value. This implies that there are few classes with large values of WMC, which is evident by inspecting the standard deviation value.

There are few classes with large values of NIV. At least 50% of the classes have one instance variable or less. The largest number of instance variables in a class is 12. This may potentially imply that the system does not apply information hiding principle appropriately for this class. The same class have its LCOM value at 90, which indicates that most variables are not shared across the member functions. Furthermore, each class is revealed to have at least four instance methods. At least 50% of the classes have 10 instance methods or less. This means that rest of the classes have more than 10 instance methods, which indicates that classes may potentially provide several services to other classes. This could be the reason behind large values of LCOM. The maximum value of NIM captured is 40. There are two classes with NIM value of 40, one having a WMC value of 194 and the other having a WMC value of 72. Both classes have a LCOM value larger than 90. This class may potentially be a *Large Class*.

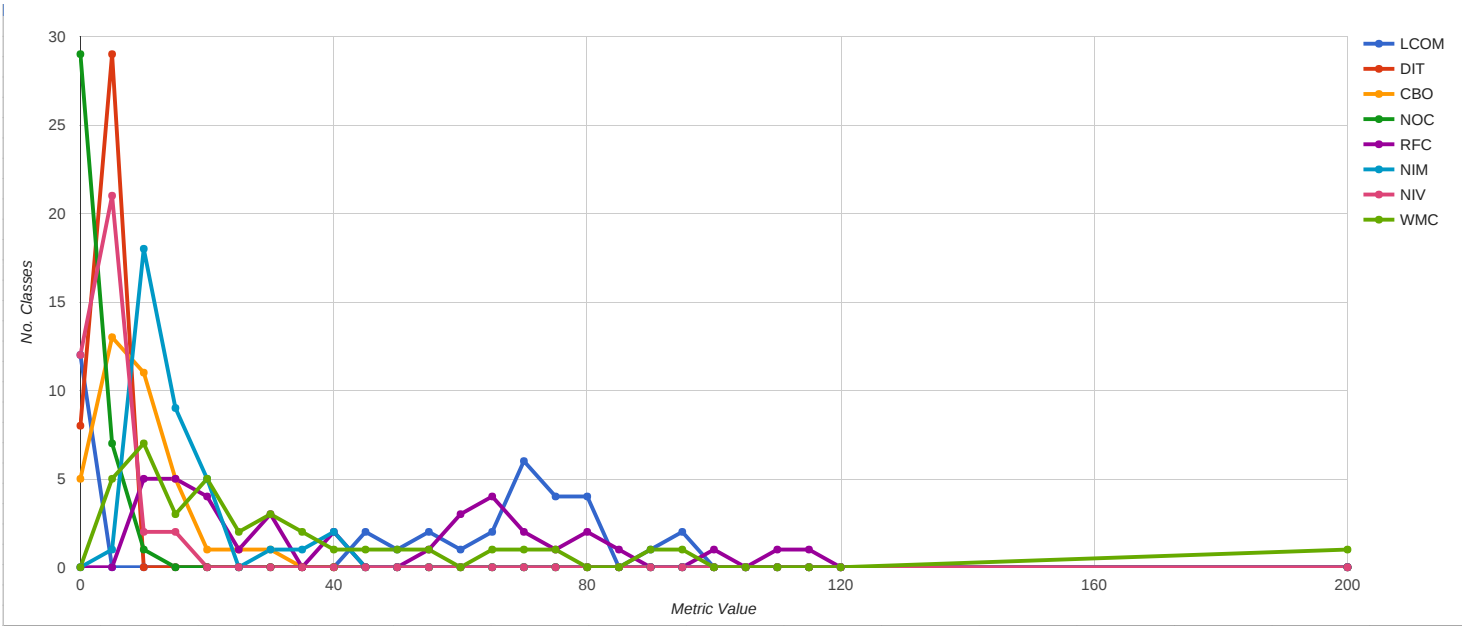


Figure 4.7: Frequency distribution of OO-metrics in Component A

Table 4.3: OO-metrics and descriptive statistics for Component B

Metric	Min	Max	Median	Sample Mean	Standard Deviation	Kurtosis	Skewness
LCOM	22	91	66	63.348	15.916	1.179	-0.761
DIT	0	1	1	0.609	0.499	-1.950	-0.477
CBO	1	21	5	6.13	4.799	2.653	1.350
NOC	0	1	0	0.087	0.288	8.605	3.140
RFC	6	48	9	14.217	11.977	2.867	1.831
NIM	6	48	9	13.434	10.974	3.822	1.986
NIV	1	10	2	3	2.504	2.51	1.711
WMC	3	325	19	35.826	66.796	17.288	3.964

Component B

We identified 23 classes in Component B. These classes are spread across 42 files, which in total contains 3905 lines of code. Figure 4.8 presents a frequency chart of the OO-metric results for Component B. Table 4.3 presents descriptive statistics of the analyzed metrics.

The values of LCOM in Component B range from 22 to 91, which indicates that none of the classes in Component B are strongly cohesive. There are only two classes with LCOM values below 50, both having low values in terms of complexity, method count, coupled objects, and instance variables. The WMC metric values range from 3 to 325. We decided to examine the class with max value of WMC. This class has a LCOM value of 74, and a CBO value of 10. Its RFC value is 44, while NIM is 36. This class has no subclasses, but it does inherit methods and variables from one superclass. This class is revealed to be the most complex class in the system in terms of cyclomatic complexity. Moreover, there are 8 classes in this component with a LCOM value of 66. We did notice that 4 of these classes has identical DIT, CBO, RFC, NIM, NIV, and WMC values by examining the metrics of the metrics for these classes. Manual inspection of the classes did not reveal any duplicated code. In terms of refactoring, we do think that all these classes need the same effort.

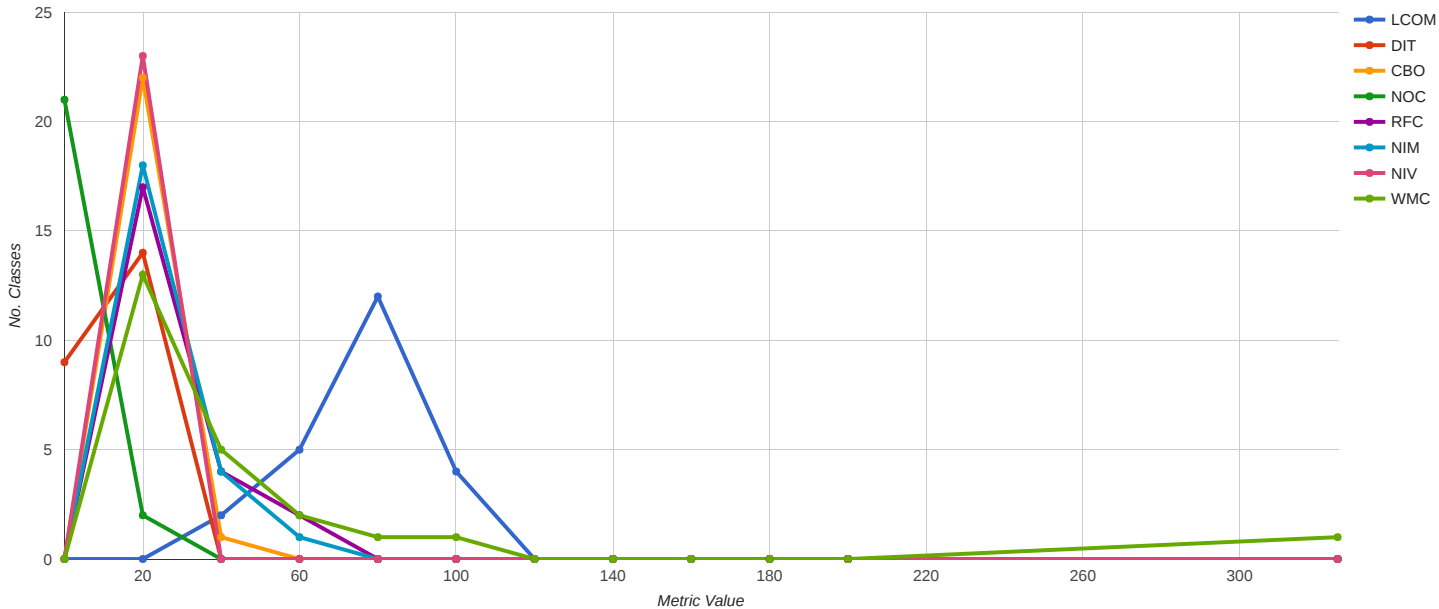


Figure 4.8: Frequency distribution of OO-metrics in Component B

Table 4.4: OO-metrics and descriptive statistics for Component C

Metric	Min	Max	Median	Sample Mean	Standard Deviation	Kurtosis	Skewness
LCOM	0	89	61	54.8	22.139	1.015	-1.089
DIT	0	1	0	0.1	0.308	7.037	2.888
CBO	1	19	4.5	6.5	5.587	-0.045	1.030
NOC	0	1	0	0.05	0.223	20	4.472
RFC	3	26	8.5	10.3	5.741	1.825	1.363
NIM	3	26	8.5	9.85	5.153	3.967	1.626
NIV	0	9	2	3.15	3.013	-0.217	1.084
WMC	3	106	12.5	27.65	30.975	1.472	1.544

Component C

Our analysis shows that component C contains 29 files, 20 classes, and 4792 lines of code. We present the descriptive statistics for component C in Table 4.4, and the frequency distribution of the metrics in Figure 4.9.

LCOM values of Component C range from 0 to 89. The median value reveal that more than half of the classes have LCOM value of 60 or more, indicating possibilities for design improvements by splitting up the classes. DIT and NOC metric values is very low, implying that inheritance may potentially not be used enough. Moreover, CBO values range from 1 to 19. Each class has a CBO value of 5 in average. One of the classes have a CBO value of 19. This class prevents reuse due to its modular design. Strong coupling complicates a system since a class is harder to understand and modify. The WMC metric values range from 3 to 106. The median value suggests that complexity in this component is well managed for most classes. However, we identified one class with a WMC value of 106. This particular class has a LCOM value of 63, and is coupled to 18 other objects. RFC, and NIM values of this class is set to 26, which are the maximum values that we analyzed for the corresponding metrics. All these values are an indication of a possible fault-prone and complex class. In addition, this class may potentially be affected by the Large Class code smell.

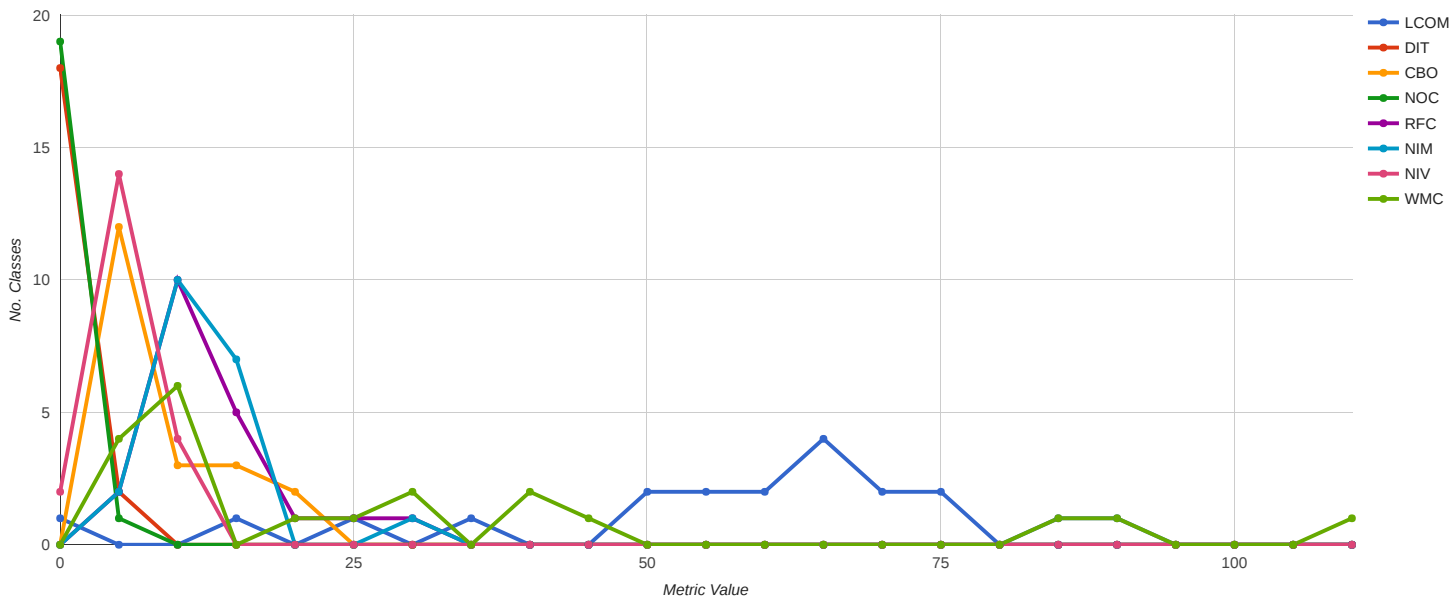


Figure 4.9: Frequency distribution of OO-metrics in Component C

Component D

Our analysis show that Component D consists of 13 files and 1647 lines of code. Among these files, we were only able to identify one class. This class has LCOM value of 68. Furthermore, our results show a DIT value of 1 for this class, and a NOC value of 0. The CBO value is set to 10. The values of RFC, and NIM have a value of 8. Moreover, the class has one instance variable. The sum of complexity in this class is 17.

Component En

Similar to Component D, our analysis identified only one class among 3 files and 367 lines of code in Component En. The metrics are very similar to Component D metrics. The class has a LCOM value of 62, indicating that the class may potentially have low cohesion in some of the methods. However, compared to Component D, this class is only coupled to one object. The sum of complexity of methods in this class is 15.

Component Ex

Our analysis found 48 files in Component Ex. These files consist of 4089 lines of code, and among these, we identified 86 classes. 38% of the number of classes of the system are located in this component. Table 4.5 presents the descriptive statistics for Component Ex, while Figure 4.10 presents the frequency distribution of the measured metrics.

Despite the fact that at least 50% of the classes have high cohesion, there are still many classes that show lack of cohesion. Our results show that there are 22 classes with LCOM value larger than 50, whereas 6 classes have a LCOM value of 80 or more. There are 2 classes with LCOM value of 100. Moreover, the average cyclomatic complexity of the classes have a value of 8.744, while the median has a value of 7. These values imply that most classes may have more polymorphism and less complexity. Figure 4.10 shows that there are only one 8 classes with value of WMC larger than 20, whereas the maximum value is 41. These values indicate that the complexity is well managed to this point. We observe that only 7 classes are self-contained by examining the CBO values. Rest of the classes are coupled to other objects. The majority of the classes has a CBO value of 5 or less. There is only one class with CBO value of 16, indicating that this class may potentially be difficult to understand and maintain.

Table 4.5: OO-metrics and descriptive statistics for Component Ex

Metric	Min	Max	Median	Sample Mean	Standard Deviation	Kurtosis	Skewness
LCOM	0	100	0	26.302	33.082	-0.971	0.756
DIT	0	3	2	1.581	1.121	-1.313	-0.234
CBO	0	16	4	5.314	4.459	-0.511	0.783
NOC	0	20	0	0.744	2.736	32.072	5.338
RFC	0	28	8	10.279	6.030	0.800	0.87
NIM	0	22	3	4.907	3.846	4.082	1.785
NIV	0	10	0	1.209	2.098	5.887	2.415
WMC	0	41	7	8.744	8.117	3.91	1.882

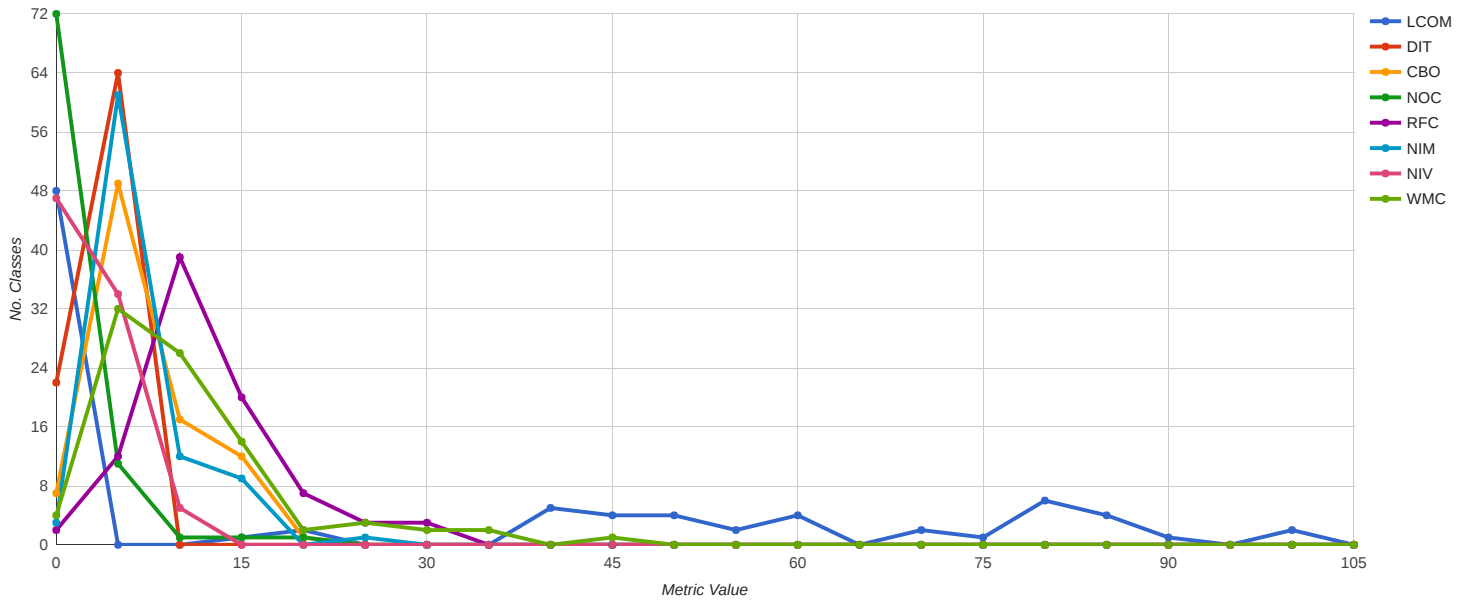


Figure 4.10: Frequency distribution of OO-metrics in Component Ex

Table 4.6: OO-metrics and descriptive statistics for Component G

Metric	Min	Max	Median	Sample Mean	Standard Deviation	Kurtosis	Skewness
LCOM	0	94	60	50.25	31.236	-0.767	-0.794
DIT	0	2	1	0.625	0.609	-0.582	0.399
CBO	0	23	5.5	6.312	4.987	1.898	1.152
NOC	0	2	0	0.25	0.622	4.231	2.357
RFC	2	30	9	10.187	6.382	2.178	1.356
NIM	0	29	7	8.437	5.459	5.669	1.745
NIV	0	18	2	3.062	3.926	5.983	2.225
WMC	1	123	12	19.437	24.794	9.572	2.824

Component G

Component G consists of 59 files. These files include 3701 lines of code and 32 classes. Descriptive statistics for Component G are summarized in Table 4.6, and its frequency distribution of metrics are presented in Figure 4.11.

Overall, the statistics indicate that there are some accumulated DD in this component. The values of LCOM range from 0 to 94, whereas 8 classes have a LCOM value of 0. However, values of LCOM for rest of the classes are larger than 50. Moreover, the statistics show that at least 16 classes have a DIT value of larger than 0, which indicates a higher degree of reuse. There are 5 classes with subclasses in this component. The WMC metric values in this component range from 1 to 123, where only class has a value of WMC larger than 100. We decided to examine the class with WMC value of 123. The metric of that class reveal a LCOM value of 92, CBO value of 22, RFC value of 30, NIM value of 9, and NIV value of 18. These values show that this class is probably influenced by the Large Class code smell, and is a candidate for inspection and possible refactoring.

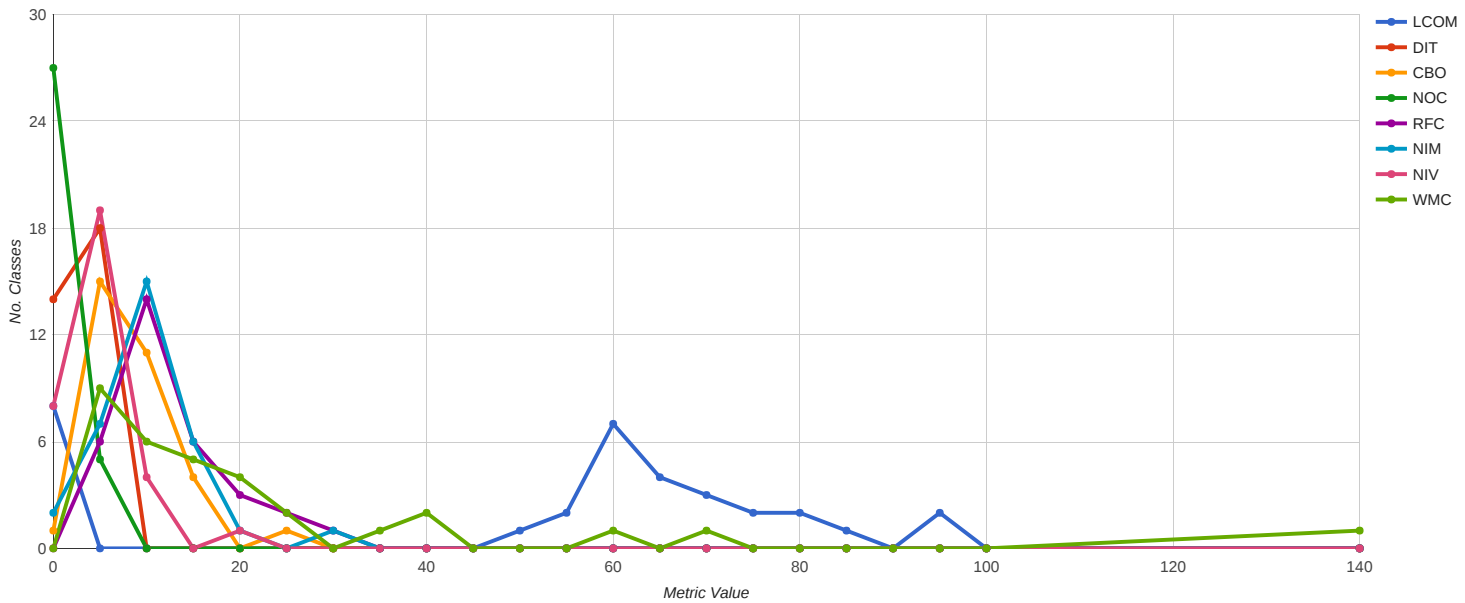


Figure 4.11: Frequency distribution of OO-metrics in Component G

Table 4.7: OO-metrics and descriptive statistics for Component L

Metric	Min	Max	Median	Sample Mean	Standard Deviation	Kurtosis	Skewness
LCOM	0	80	68	50.857	35.130	-0.908	-1.135
DIT	0	1	1	0.571	0.534	-2.8	-0.374
CBO	1	14	6	6.714	4.609	-0.712	0.327
NOC	0	0	0	0	0	N/A	N/A
RFC	5	12	9	8.571	2.936	-2.012	-0.239
NIM	3	12	9	8.286	3.402	-1.221	-0.555
NIV	0	5	1	1.571	2.070	-0.535	1.120
WMC	4	42	11	19.571	14.524	-1.521	0.577

Component L

Component L consists of 16 files, 849 lines of code. Among these, we identified 7 classes. Figure 4.12 presents the frequency distribution of the metrics, while Table 4.7 presents the descriptive statistics for this component. The values of LCOM range from 0 to 80. There are only 2 classes with LCOM value at 0. However, rest of the classes are showing lack of cohesion. These classes are candidates for inspection, and might eventually be split up into multiple classes. Moreover, 4 classes have a DIT value of 1. WMC metric values range from 4 to 42. There are only 2 classes with WMC values larger than 30.

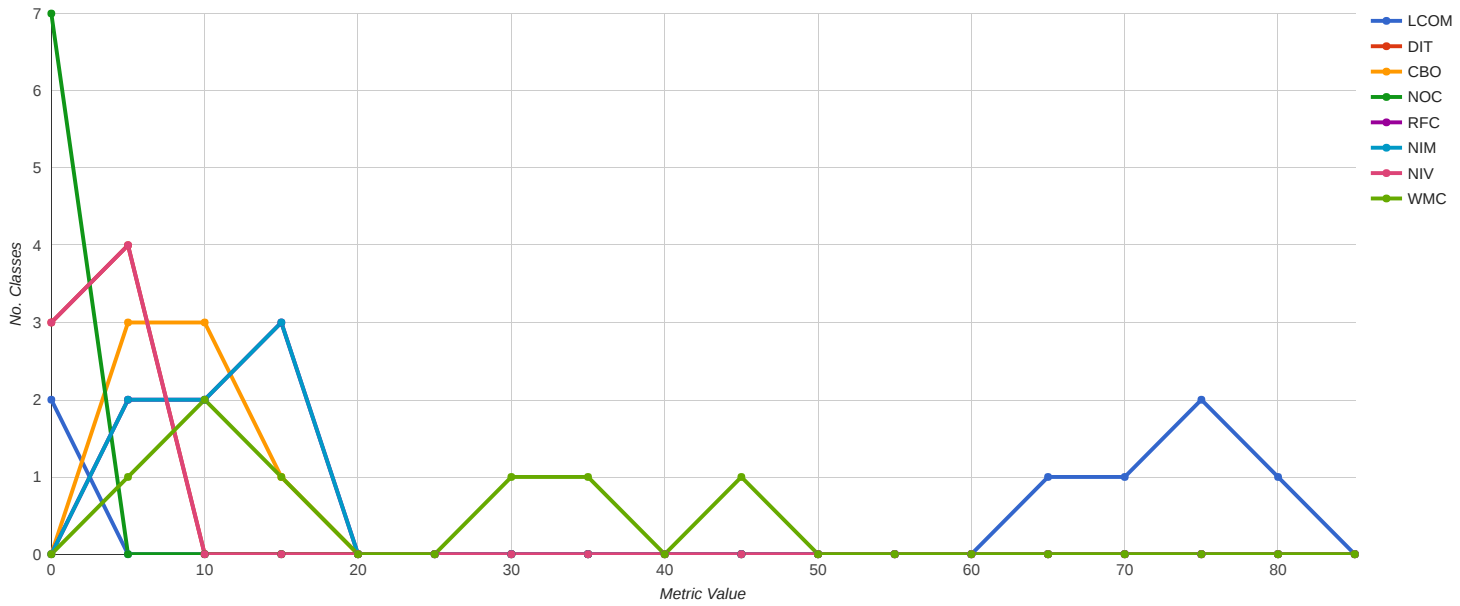


Figure 4.12: Frequency distribution of OO-metrics in Component L

Table 4.8: OO-metrics and descriptive statistics for Component N

Metric	Min	Max	Median	Sample Mean	Standard Deviation	Kurtosis	Skewness
LCOM	0	79	70.5	54.5	33.899	-0.083	-1.378
DIT	0	1	0	0.25	0.463	0	1.440
CBO	3	19	12.5	11.5	4.536	1.941	-0.410
NOC	0	1	0	0.125	0.353	8	2.828
RFC	6	32	9	11.625	8.568	6.227	2.413
NIM	6	21	8.5	9.75	5.036	3.976	1.896
NIV	0	8	5.5	4.375	3.068	-1.142	-0.630
WMC	4	125	32	40.375	36.707	5.172	2.092

Component N

Component N consists of 17 files. There are 1839 lines of code spread across these files. We were able to identify 8 classes in this component. Descriptive statistics for Component N are presented in Table 4.8, while the frequency distribution of the metrics are presented in Figure 4.13.

LCOM metric values range from 0 to 79. The median value shows that more than half of the classes has LCOM value of 70 or more, indicating that classes in this component can potentially be split into more classes to increase the cohesion of each class. By taking a closer look at Figure 4.13, we identify two classes with LCOM values interval from 75 to 80. More precisely, one class has LCOM value of 78 while the other class has LCOM value of 79. We decided to examine the class with LCOM value of 79, and identified that this class has WMC value of 125, RFC value of 32, CBO value of 19, and NIM value of 21. These values tell us that this class may be affected by Large Class code smell.

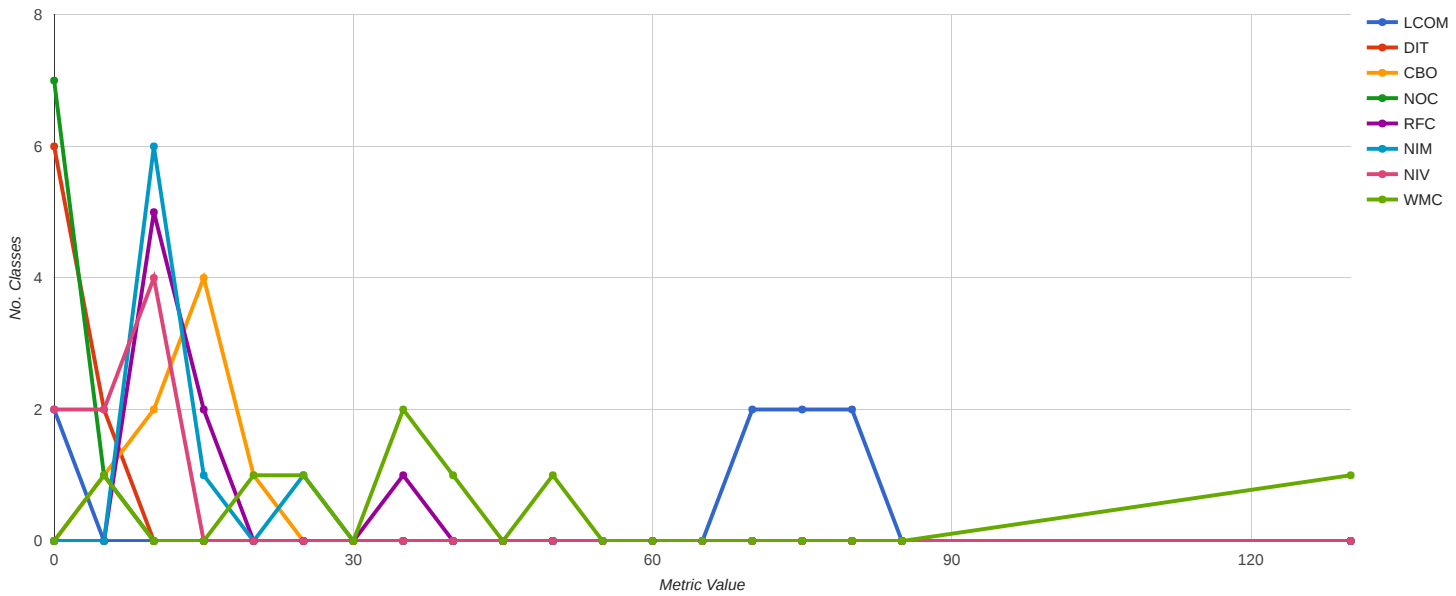


Figure 4.13: Frequency distribution of OO-metrics in Component N

Table 4.9: OO-metrics and descriptive statistics for Component P

Metric	Min	Max	Median	Sample Mean	Standard Deviation	Kurtosis	Skewness
LCOM	0	81	63	57.25	25.409	4.625	-1.987
DIT	0	1	0	0.25	0.463	0	1.440
CBO	0	13	6	6	4.472	-0.645	0.255
NOC	0	1	0	0.125	0.353	8	2.828
RFC	2	14	6.5	7.375	3.502	1.523	0.647
NIM	2	13	6.5	6.875	3.044	2.986	0.765
NIV	0	6	3.5	3.125	2.031	-0.886	-0.223
WMC	1	47	8	15.75	15.809	-1.037	1.350

Component P

The descriptive statistics calculated for Component P are presented in Table 4.9. The frequency distribution of the metrics are presented in Figure 4.14. We identified 12 files in Component P, consisting of 12 files, 722 lines of code, and 8 classes.

LCOM values in Component P ranges from 0 to 81. A median value shows the level of cohesiveness in the system. In this context, the median value shows that more than half of the classes have large LCOM values, implying that these classes are improperly designed and should be split up to make them more cohesive. We identified three classes with values of LCOM larger than 70. The DIT results range from 0 to 1, implying that most classes have flat inheritance hierarchies. There are only 2 classes having a DIT value of 1. Despite the fact that 6 of the classes have DIT metric of 0, they alone may not tell us if classes are part of an inheritance tree or if they are root classes. By examining the NOC results, we see that only one class has a NOC value of 1. Six classes have both NOC and DIT values of zero, indicating that they are not part of an inheritance hierarchy. CBO values ranges from 0 to 12, with a mean and median of 5.5. WMC values range from 1 to 47.

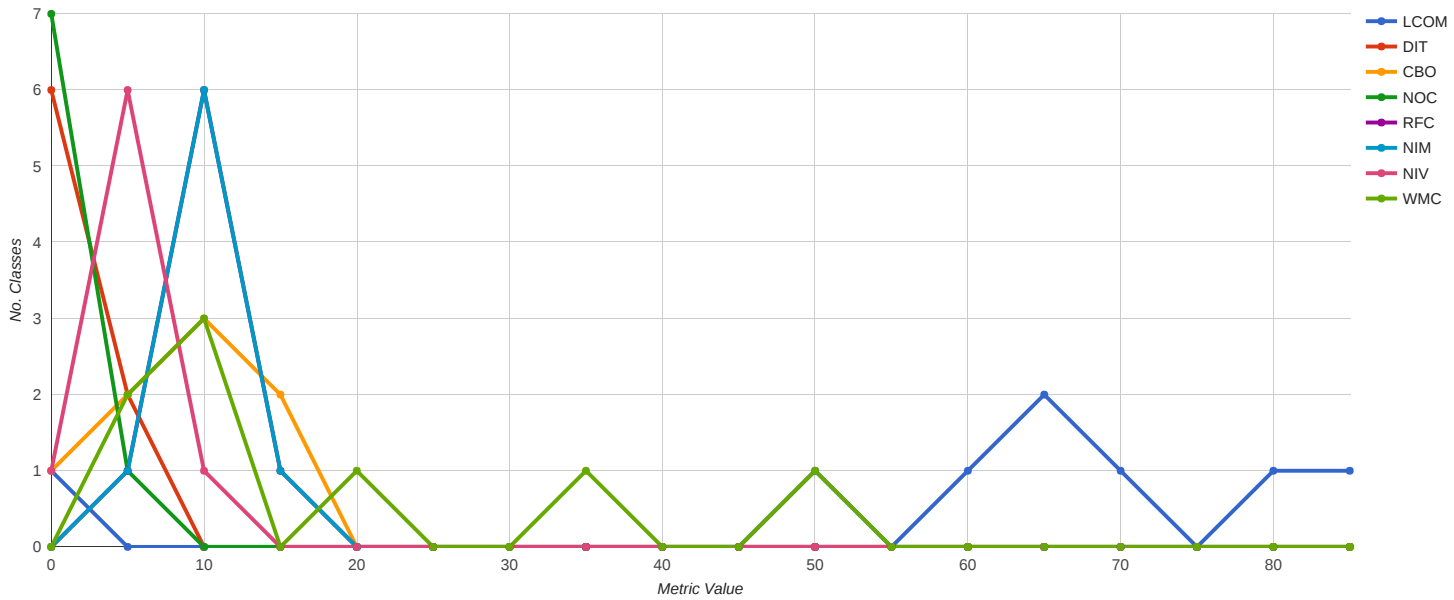


Figure 4.14: Frequency distribution of OO-metrics in Component P

Table 4.10: OO-metrics and descriptive statistics for Component S

Metric	Min	Max	Median	Sample Mean	Standard Deviation	Kurtosis	Skewness
LCOM	33	50	41.5	41.5	12.021	N/A	N/A
DIT	0	1	0.5	0.5	0.707	N/A	N/A
CBO	3	7	5	5	2.829	N/A	N/A
NOC	0	0	0	0	0	N/A	N/A
RFC	6	9	7.5	7.5	2.121	N/A	N/A
NOM	6	9	7.5	7.5	2.121	N/A	N/A
NIM	6	9	7.5	7.5	2.121	N/A	N/A
NIV	1	1	1	1	0	N/A	N/A
WMC	6	19	12.5	12.5	9.192	N/A	N/A

Component S

Table 4.10 presents the descriptive statistics for Component S. Component S contains 4 files, which consists of 223 lines of code and 2 classes. The LCOM values show that there are possibilities to improve the design of this component by splitting up methods that does not share fields with each other to separate classes. Moreover, the results point out that none of the classes has any subclasses. However, the results reveal that one of the classes inherits methods and variables from a superclass. Moreover, WMC metric values indicate that the classes have less complexity and greater polymorphism.

Component W

Component W contains only one file. This file consists of 69 lines of code. This component has one class. The class have a LCOM value of 58, indicating that some instance variables are not shared across the member functions. Moreover, DIT and NOC values of 0 indicates that this class does not inherit from a superclass, or has any subclasses. The results reports a CBO value of 4, implying that this class is not tightly coupled with other objects. There are only 2 instance variables in this class, and 6 methods, in which all are instance methods. WMC value of this class is set to 8.

4.2 Identifying Code Smells using Automatic Static Analysis Tools

Code smells are used to find problematic classes. As we explained in Chapter 2, one of the ways to identify DD is to look at the number of code smell in the source code. Code smells are an indicator of software design flaws in OO-systems that can decrease software maintainability which may lead to issues in further evolution of the system [39]. The analysis of code smells in the system is based on ASA tools. We investigated each tool that we could use in this project, and evaluated the smells they are able to detect.

Table 4.11 describes the number of code smells that were identified using automatic static analysis tools.

Table 4.11: Number of Code Smells detected

Code Smell	Detected
Long Method	10
Long Parameter List	15
Duplicated Code	Approximately 5% of the source code.
Speculative Generality	1153
Dead Code	151

Duplicated Code

Duplicated Code is found by looking for pieces of code that appears at multiple places in the source code, both internally in a file or in another file. A piece of code is considered duplicated if the piece of code contains at least 10 lines of code and occurs at multiple places in the source code. Table 4.11 shows the number of duplicated code found by *SonarQube*, expressed as a percentage value. The results show that roughly 5% of the source code contains duplicated code including the test files. This corresponds to approximately 4395 lines of code, affecting 39 files across the system. We identified that roughly 54% of the duplicated code is located in Component A. The other duplicated lines are spread across Component B, N, P, C, D, L, Ex, S, and G. Table 4.12 summarizes duplicated lines in the different components. The first column contains the name of the components, while the second column summarizes number of files affected by duplicated code and the number of duplicated lines among these files.

Table 4.12: Duplicated Code in Project "Firmus"

Component	Information
Component A	12 files, 2400 LOC
Component B	6 files, 366 LOC
Component C	3 files, 284 LOC
Component D	2 files, 80 LOC
Component Ex	4 files, 305 LOC
Component G	4 files, 311 LOC
Component L	1 file, 30 LOC
Component N	3 files, 301 LOC
Component P	2 files, 124 LOC
Component S	2 files, 194 LOC

Long Method

Understand for C++ considers a *Long Method* as code smell if lines of code in method exceeds 200 lines. We identified 10 long methods, spread across six different files. 7 of 10 long methods are located in the test files.

Long Parameter List

Long Parameter List code smell is detected by comparing the total number of parameters in a method against a fixed threshold. The maximum number of parameters allowed in a method using *CppDepend* is set to 5. This means that 6 or more parameters in a method are considered as code smell. The results from *CppDepend* report 15 hits of *Long Parameter List* code smell, where 3 hits are considered as critical. A *Long Parameter List* is considered as critical when total of parameters in a method is higher than 8. The largest number of parameters in a method we identified was 12. These results were verified manually by examining the class diagrams for the corresponding methods.

Speculative Generality

Speculative Generality is detected by detecting unused classes, methods, fields, or parameters. Table 4.13 summarizes *Speculative Generality* code smell that were identified through a code analysis using *Understand for C++*. The results are divided into the categories unused methods, unused local variables, and unused static globals.

Table 4.13: Speculative Generality in Project "Firmus"

Category	Hits
Unused Methods	794
Unused Local Variables	346
Unused Static Globals	13

Dead Code

Fowler and Beck [101] do not classify *Dead Code* as code smell. However, *Dead Code* should be classified as a code smell, as it is a quite common problem as it hinders code comprehension and makes the current program structure less obvious [38]. We examined three types of *Dead Code* code smell in Project "Firmus": "Commented Out" Code, Unreachable Code, and Unnecessary Includes in Header Files. In total, we found 151 hits of *Dead Code* code smell, which we have summarized in Table 4.14.

Table 4.14: Dead Code in Project "Firmus"

Category	Hits
"Commented Out" Code	67
Unreachable Code	10
Unnecessary Includes in Header Files	74

This thesis has investigated DD in safety-critical systems. Section 5.1 analyses OO-metrics that have been measured by applying a set of threshold values to uncover possible insecure parts of the code. Section 5.2 evaluates the research by answering the research questions.

5.1 Analysis of Object-Oriented Metrics by Applying Threshold Values

Software metrics measurement in OO-software is important in terms of quality management [98, 99], as software metrics can be used as predictors of fault-prone classes in OO-systems [56]. A study by Basili et. al [56] assessed Chidamber and Kemerers [55] suite of OO-metrics as predictors of fault-prone classes. Their results implied that OO-metrics appear to be useful to predict class fault-proneness during early phases in the software life-cycle. However, it is hard to assess the quality of a software with single OO-metrics.

A metric is meaningless without its threshold values. According to Tarcisio [98], software metrics are not effectively used in software industry due to the fact that for the majority of metrics, thresholds are not defined. Threshold is defined as values used to set ranges of desirable and undesirable metric values for measured software [99]. Knowing thresholds for metrics allow us to assess the quality of a software, and we may be able to identify where in a design errors are likely to occur. Lanza et al. [100] present two ways to identify major sources for threshold values; *statistical information* (i.e., thresholds that are based on statistical measurements), and *general accepted semantics* (i.e., thresholds that are based on information which is considered common).

We identified the following thresholds for the individual metrics using their descriptive statistics from Table 4.1. The goal with applying threshold values to the metrics is to identify the classes with major design violations. Threshold values have been identified using statistical information [100]. For each metric, we have used the sample mean and standard deviation from Table 4.1. There are three different threshold values: *low/good*, *high/typical*, and *extreme/bad*. **Low/good** value corresponds with the mean, and represents the most typical value in the data set [102]. Another way to calculate *low/good* is by subtracting the standard deviation from the sample mean [100]. However, the standard deviation in the data set may be larger than the sample mean, which leads to negative threshold values. This is why we have chosen to use the median value. **High/typical** value is the sum of the sample mean and the standard deviation. Classes with *high/typical* values are acceptable, but they should be inspected if possible. **Extreme/bad** value is simply the second threshold value multiplied with 1.5 [100]. It is considered as extreme, and should not be included in the data set.

Table 5.1 presents the metrics and their threshold values. The table includes threshold values that has been proposed by other researchers for the selected metrics. Despite the fact we were able to derive threshold values, our study should still be viewed as evaluation of one software system. Although our techniques can be applied to other kinds of safety-critical software, it is necessary to review that our threshold represent only local data. For general usage, it is recommended to compare our threshold with the threshold defined by other researchers.

Table 5.1: Thresholds for OO-software metrics

Metric	Observed Value	Low	High	Extreme	Recommended Value
LCOM	100	42	75	112.5 (100+)	72.5 [98]
DIT	4	1	2	3+	4 [98], 5 [96, 103, 104]
CBO	30	6	11	17+	5 [96, 104], 14 [105, 106]
NOC	20	0	2	4+	3 [98], 5 [104], 10 [103]
RFC	115	15	34	51+	50 [96, 103]
NIM	48	8	15	23+	N/A
NIV	18	2	5	7+	N/A
WMC	325	19	51	76+	34 [98], 40 [96], 50 [106]

Depth in Inheritance Tree

DIT indicates how deep a class is in the inheritance tree. Classes with higher value of *DIT* are associated with higher defects [107]. It is evident that a deep inheritance makes software maintenance more difficult [108]. Moreover, higher degree of *DIT* indicate a trade-off between increased complexity and increase reuseability. A low median value indicates that at least 50% of the classes tend to be close to the root in the inheritance

hierarchy by following Chidamber and Kemerer's [55] guide to interpreting *DIT* metric using descriptive statistics. A low median value had a typical value of 1 and 3 in their study. However, if the majority of *DIT* values are below 2, it may represent a poor exploitation of the advantages of OO-design and inheritance, because a *DIT* value of 2 and 3 indicates higher degree of reuse.

Figure 4.3 in Chapter 4 show that approximately 39% of the classes have a *DIT* value of 0, while 26.7% of the classes have a *DIT* value of 1. These classes are considered to be close to the root in the inheritance tree, and there may be a probability of not exploiting the advantages of OO-methodologies. Classes with high values of *DIT* have shown to be very significant in identifying fault-prone classes [56,109]. We have derived the following threshold values for the *DIT* metric: a *low/good* value of 1, a *high/typical* value of 2, and an *extreme/bad* value of 3. We observe that at least 26.7% of the classes satisfies good value by applying these thresholds. However, the *extreme/bad* value of 3 does not comply with the median value of 3, which is considered as a good value [55]. Therefore, we have chosen to apply the recommended max value of 5. As shown in Figure 4.3, there are no classes with a *DIT* value of 5 or more. However, there are two classes with *DIT* value of 4. These results do indicate that reuse opportunities through inheritance is limited and perhaps compromised in favor of comprehensibility of the overall architecture. On the other hand, low values of *DIT* suggest that appropriate design preferences are being followed by the company [107].

Number Of Children

NOC measures the number of children a class has. The greater number of children a class has, the greater potential of reuse, since inheritance is a form of reuse. However, large values of *NOC* indicate a misuse of subclassing [56]. Moreover, classes may require more testing. Our results reveal a *NOC* median value is 0 (See Table 4.1). The distribution of *NOC* metrics in Figure 4.4 shows that approximately 86.5% of all classes have no children. In addition, Figure 4.4 show that a small number of classes have many immediate subclasses. Both Chidamber and Kemerer [55] and Basili et al. [56] have observed the similar median values for *NOC* metric in their respective studies. These values suggest that designer may not be fully exploiting the advantages of inheritance as a basis for designing classes. Lack of communication could be another reason between class designers, which may hinder developers to reuse.

The following threshold values for *NOC* metric were derived from its descriptive statistics in Table 4.1: a *low/good* value of 0, a *high/typical* value of 2, and an *extreme/bad* value 4. These values can be compared to threshold values identified by Filo et al. [98]. The difference is that they classify *extreme/bad* value as 4. According to Figure 4.4, 91.2% of the classes are classified as *low/good*, 5.2% of the classes are classified as *high/typical*, and 3.6% of the classes are classified as *extreme/bad*. Classes classified as *extreme/bad* are an indication of classes that may be hard to understand and maintain, and are potential candidates for inspection. However, similar to *DIT* metric values, *extreme/bad* values indicate limitation in reuse opportunities in favor of comprehensibility of the overall architecture. This suggest that appropriate design preferences for the

system are probably being followed.

Lack of Cohesion in Methods

LCOM is related to the counting of methods using common attributes in a class. As mentioned, smaller values of *LCOM* represent cohesive and independent classes, which is desirable since it promotes encapsulation. Larger values of *LCOM* increase the complexity of the class, hence increasing the likelihood of errors during the development process. Our derived threshold values for *LCOM* metric reveal a *low/good* value of 42, a *high/typical* value of 75, and an *extreme/bad* value of 112.5. However, as we see in Figure 4.1, there are no classes with *LCOM* value larger than 100. This result is influenced by the special shape of the data from *LCOM* metric. The recommended threshold value for *LCOM* is 72.5. This value is similar to our definition of *high/typical*. We chose to apply the recommended threshold value for *LCOM* metric. We identified 41 classes with larger *LCOM* than 72.5, which we have illustrated in Figure 4.1. Our data reveal that *LCOM* values seem to increase with the size of classes. Most classes with a high value of *LCOM* revealed to have large number of methods. These methods indicate higher disparateness in the functionality provided by the class. Chidamber and Kemerer [55] and Basili et al. [56] state that classes with large values of *LCOM* could be more error prone, and more difficult to test. A refactoring option would be to split the classes into two or more classes that are more well defined in terms of behavior. Moreover, *LCOM* metric can be used by the developers to keep track of whether the cohesion principle is adhered.

Coupling Between Object classes

CBO refers to the number couplings between object classes. Higher values of *CBO* indicate the extent of lack of reuse potential of a class, and that more effort may be required to maintain and test the class. Classes with higher values of *CBO* are associated with higher defects [107]. According to Chidamber and Kemerer [55], a low median had a typical value of 0, while a high median had a value of 9. A median value of 0 suggests that at least half of the classes are self-contained and do not refer to other classes. Furthermore, Basili et al. [56] state that *CBO* appears to be useful to predict class fault-proneness.

Table 4.1 reveal a *CBO* median value of 5, indicating that at least 50% of the classes refer to 5 or less object classes. The *CBO* value is generally less for most classes. Therefore, these classes are easy to understand, reuse, and maintain. However, we did notice that some of the classes with low *CBO* value had higher values in the other metrics, such as *WMC* and *LCOM*.

Threshold values for the *CBO* metric were derived, and the following values were identified: a *low/good* value of 6, a *high/typical* value of 11, and an *extreme/bad* value of 17 or more. Chidamber and Kemerer [55] classified a median value of 9 as high, which is similar to the recommended values in Table 5.1. Therefore, we have chosen to apply our derived *high/typical* threshold value of 11 on the *CBO* metric results. Figure 4.2 reveal

that 11.3% of all classes have a coupling of 11 or more. *CBO* should be kept to a minimum in order to promote encapsulation and modularity for a class. Classes with large values of *CBO* are more sensitive to changes in other parts of the design, which affects the systems maintainability, reusability, and testability. Furthermore, the *CBO* metric can be used as a way to track whether the class hierarchy is losing its integrity, and whether different parts of the system are developing unnecessary relations between classes.

Response For Class

RFC is defined as the total number of methods that can be executed in response to a message to a class. This count includes all methods available in the class hierarchy. Large values of *RFC* makes testing and debugging more complicated since it requires a greater level of understanding on the part of the tester [55]. In addition, large values of *RFC* may increase the complexity of the class. It can be hard to predict the behavior of the class since it requires a deep understanding of the potential interactions that the objects of the class have with the rest of the system, hence affecting the classes' testability. Our analysis reveal that most classes are able to invoke a small number of methods. According to Figure 4.5, 74.2% of the classes have a *RFC* value of less than 15.

The following threshold values were derived for the *RFC* metric: a *low/good value* of 15, a *high/typical* value of 34, and an *extreme/bad* value of 51 or more. The *extreme/bad* value is very similar to the recommended max value, although it is acceptable to have a *RFC* up to 100 [103]. However, it is recommended that a class does not have *RFC* value larger than 50. We applied the threshold value on Project "Firmus", and identified 17 classes with a *RFC* value larger than 50. One of the class have its *RFC* at 115. Research point out that *RFC* has been found to be highly correlated with *WMC* and *CBO* [55]. Our results reveal that *RFC* is not correlated with either *WMC* or *CBO*. For example, the class with *RFC* value of 115 had a *CBO* value of 2 and *WMC* value of 22. However, its *DIT* value is 4 and may explain that the large values of *RFC* comes from inherited methods. We identified the same pattern with the class with *RFC* of 109. This particular class has a *CBO* value of 3 and a *WMC* value of 10. However, its *DIT* value is 4. This observation indicates that *RFC* is associated with *DIT*. Moreover, *RFC* values tend to be low because there are a number of classes that have no parents.

Weighted Methods per Class

WMC refers to the sum of complexity in each method. A greater value of *WMC* indicates a complex class. *WMC* is usually affected by the number of methods in a class, but there may be some cases where some methods may have low complexity while other methods have high complexity. Table 4.1 reports a median value of 10 and sample mean value of 19.707 on *WMC* metric. The results are quite similar to what Chidamber and Kemerer [55] and Basili et al. [56] measured in their respective studies. Another aspect of *WMC* data is the similarity of the frequency distribution of the metric values in both of their and our study. In Figure 4.6, we notice that approximately 60% of all classes

have a *WMC* value of less than 10. This suggest that most classes have a small number of methods. The following threshold values were derived for *WMC* metric: a *low/good* value of 19, *high/typical* value of 51, and an *extreme/bad* value of 76. According to Table 5.1, the recommended max value are ranging between 30-50, which is close to our *high/typical* value. Therefore, it seems like this value is more appropriate to use to identify classes with high values of *WMC*. We apply the threshold, and identify 7.4% (i.e., 17 classes) with a *WMC* value larger than 51. One of the classes has a *WMC* value of 325, which indicates that time and effort to develop and maintain this class may be high. The same class have 44 methods and a *CBO* value of 10. A class with many methods are most likely to be application specific, hence reducing its reuse potential. Moreover, an increase in number of methods is associated with an increase in defects [107]. Another class with a *WMC* value of 194 reveal to have 40 methods. This suggest that classes with more methods tend to have higher complexity, which indicates an increase in defects. In such cases, maintenance effort increase drastically. The 17 classes with large value of *WMC* are primary candidate classes in which code inspection and potentially refactoring is needed.

5.2 Research Evaluation

This thesis is focused on identifying DD in safety-critical systems. In particular, we were able to identify DD by measuring OO-metrics and detecting code smells. The software design was analyzed using a suite of OO-metrics proposed by Chidamber and Kemerer, one of the most used OO-measures. The result does indicate that embedded software developers accumulate TD, despite the fact that they cannot contain any errors [79, 81, 82]. However, it seems like the accumulated TD is manageable. An important aspect when delivering a product is to make sure the product is stable and reliable. According to our results, we can now answer the research questions that were stated in Chapter 1. Table 5.2 show the connection between our contributions, and research questions.

Table 5.2: Connection between contributions, and research questions

Contributions	Description	Research Questions
C1	Empirical knowledge about DD identification in safety-critical systems by object-oriented metric analysis and code smell detection.	RQ1
C1.1	A set of threshold values for the measured OO-metrics. See Table 5.1.	RQ1
C2	Empirical knowledge about the effects of having DD in safety-critical systems.	RQ2
C3	Empirical knowledge about the different types of DD in safety-critical systems.	RQ3
C4	Empirical knowledge about paying DD.	RQ4

RQ1: How can DD be identified?

The first research question is related to the techniques we have used to identify DD in this research. Automatic static code analysis, and OO-metrics measurements using Chidamber and Kemerer's suite of metrics have proven to be useful in the context of DD identification. We were able to collect a large amount of measurements that characterize the software by using various tools. Tools that have been used through this research are *Doxygen*, *Understand*, *SonarQube*, *CCCC*, *CppCheck*, *CppDepend*, and *Enterprise Architect*.

Code smells are an example of design flaws that can degrade maintainability of source code, which implies that code smells can be used as an indicator to identify fault-prone files in the system. Moreover, code smells are an indication of refactoring possibilities in code base. We have been able to identify *Duplicated Code*, *Long Method*, *Long Parameter List*, *Speculative Generality*, *Dead Code*, and *Large Classes* code smells in our analysis. *Duplicated Code* was identified using *SonarQube*. *SonarQube* analyzed the entire source code base, and identified similar code blocks that had an appearance in multiple places. *Long Method*, *Speculative Generality*, and *Dead Code* were identified using *CppDepend*, *Understand*, and *CCCC*. *Long Parameter List* code smell were detected using *CppDepend* and confirmed using *Enterprise Architect*.

Another approach to assessing the software quality is based on OO-metrics [48]. OO-are able to predict maintenance effort more than traditional metrics can [57] by identifying design flaws, and defect-prone, change-prone, and fault-prone classes [56]. In addition, OO-metrics may potentially affect the quality attributes of a system. For example, large values of *WMC* will affect a system's maintainability and reusability [54]. OO-metrics are calculated over data that are extracted from the systems source code. Chidamber and Kemerer's suite of metrics have been applied in empirical investigations of OO-systems by multiple researchers, including Basili et al. [56], Chidamber and Kemerer [55], Okike [110], and Bakar et al. [111]. We have applied this suite of OO-metrics to discover potential fault-prone classes.

Understand C++ is an integrated development environment that enables static code analysis through visualization, documentation, and metric tools. The software is capable of analyzing projects with multiple lines. An academical license tool was provided to us, and the tool was used in our case study to compute software metrics. Every file in the system were analyzed, and metrics are then extracted from these files. The tool has been used by researchers. *Understand* have been proven to be useful for code analysis. Malhotra et al. [112] calculated threshold values of OO-metrics by using statistical models. *Understand* was used to extract relevant metric data from one of the systems. Furthermore, Codabux et al. [48] extracted class-level metrics for defect- and change-prone classes using *Understand*.

Threshold values were derived in this study to assist us on identifying classes with design flaws. Threshold can be defined as the upper bound value for a metric. A metric value with a greater than its upper bound threshold value can be considered as problematic, while a metric value lower than its upper bound threshold value can be considered

as acceptable. Both metric values and threshold values can be compared in design phase of software development to identify the metrics whose value is bigger than its threshold. Results from our study suggests that refactoring can be applied to classes with larger metric value than its threshold. In other words, we were able to predict possible fault-prone classes by applying metric threshold. For instance, a class with a *WMC* value larger than its threshold value indicate high complexity. Fault-prone classes should be used as early quality indicators, and actions should be take based on extent of the problem. For example, the project team may choose to redesign the entire class to achieve the desired metric value.

RQ2: What are the effects of DD?

The second research question is related to the consequences of having design flaws in safety-critical software. OO-metrics have proven to be indicators of problems in software design. Classes with larger metric values than its threshold values may affect the quality attributes of a system. Our analysis reported multiple classes in which following metric values were larger than its threshold values; *CBO*, *RFC*, *WMC*, *LCOM*, *CBO*, and *NOC*.

Classes with large values of *WMC* are likely to be more application specific, hence affecting the software's understandability, reusability, flexibility, and maintainability quality attributes [113, 114]. Moreover, classes with large values of *RFC* may be harder to understand and test, which also affects the software's understandability, testability, maintainability [113]. In addition, *RFC* may affect system's functionality and reusability as objects communicates by message passing [114]. *LCOM* measures the cohesion of a system. Lack of cohesion in a class increases the classe's complexity, which may lead to errors during development. This metric affects the systems efficiency, reusability, and understandability [113, 114]. Large values *CBO* complicates a system, since a module is harder to understand, change, reuse, and maintain due to its excessive coupling with other classes. *CBO* evaluates the systems understandability, extendability, efficiency, reusability, testability, and maintainability of a class [113, 114]. *DIT* and *NOC* metric are related to inheritance which enables reuse. Large values of *DIT* indicates deep hierarchy which constitutes greater design complexity. Deep hierarchy enhances the potential reuse of inherited methods but in trade-off, complexity will increase which affects other quality attributes. *DIT* metric evaluates efficiency, reusability, understandability, and testability [113] of a software. *DIT* metric may also be related to flexibility, extendability, effectiveness, and functionality [114]. Furthermore, *NOC* primarily evaluates efficiency, testability, and resuability of a system [113], but it may also influence flexibility, understandability, extendability, and effectiveness of a system [114].

Code smells are manifestations of design flaws that can have negative influence on software quality. Although the results did not investigate the effects of the identified code smells on the system studied, we have reviewed the consequences of code smell that other researchers have discovered. Sjoberg et al. [115] investigated the relationship between 12 different code smells and maintenance effort. Their result show that none of these code smells were associated with more maintenance effort. Similarly, Hall et al. [116] state that some smells indicate fault-prone code in some circumstances but that the effect

these smells have on faults and software maintainability is small. Lindsay et al. [117] state that not all *Large Class* code smell are bad, some of them could be explained by decisions that perhaps are not entirely under the control of the developer. For example, a choice of particular design pattern may lead to this smell. On the other hand, both Li et al. [118] and Dhillon et al. [119] state that bad smells are positively associated with increased error rate in software projects. Furthermore, Olbrich et al. [120] proved that *God Class* code smell have a negative effect on software quality in terms of change frequency, change size, and number of weighted defects. Khomh et al. [40] provided evidence that classes with specific code smells are more subject to change than others.

These articles state that certain types of code smells can have minimal effects, while other types of code smells can have greater effects on the quality attributes of a system. For example, *God Class* code smell will create more maintenance effort than *Dead Code* code smell will. To determine the actual effects of having code smells, the file and class metric in which the code smells are located should be measured and an action should be based on those results.

RQ3: What kind of DD can be found in embedded systems?

The third research question is related to our results from the case study. We were able to identify fault-prone classes by applying threshold on the derived OO-metrics, and code smells using automatic static analysis code tools during our case study. We have been able to identify *Duplicated Code*, *Dead Code*, *Speculative Generality*, and *Long Method* as the most common code smells in general. Furthermore, OO-metrics were useful for investigating each individual class and the software in general. For example, we identified possible *Large Class* code smell by examining the metrics. These smells indicate possible fault-prone classes.

RQ4: How to pay DD?

The last research question is related to the management of DD. A common approach to keep DD from growing, or to pay back DD, is to conduct refactoring and re-engineering. Codabux et al. [5] mention that refactoring is a common way to manage and ultimately get rid of TD. In addition, refactoring seems to improve important internal measures for reusability of OO-classes [121]. Without refactoring, the design of the program will decay over time.

Our results show that 5% of the source code, including the test files, contain duplicated code. Removing duplicated code will reduce the number of lines of code by default. Some duplicated code were found in the same class, while others were spread across multiple classes. Therefore, they must be handled differently. In most cases during our data collection, the same block of code occurred in the same file or class. If the same code of block exists in two different methods of the same class, *Extract Method* refactoring technique should be applied. This technique creates a new method, which can be invoked from both methods that contained duplicated code [37]. In addition, we identified some

cases where the same block of code occurred in different classes in the system. For example, two identical files were identified in both Component S and B. If that is the case, then *Extract Class* refactoring technique should be applied. This technique creates a new class, superclass, or a subclass, which can be reused by both of the components with duplicated code. However, if the method is a critical part of one class, then the method should be invoked by the other class.

Moreover, we identified 10 *Long Method* code smells. *Extract Method* can be applied to shorten a method. This technique finds parts the method that seem to go nicely together, and creates a new method [37]. In some cases, a method may have lots of parameters and eventually temporary variables. If the method has lots of temporary variables to hold the result of an expression, then *Replace Temp with Query* refactoring technique should be applied [37]. *Replace Temp with Query* extracts the expression into a method, and all references to the temporary variables will be replaced with the expression. Moreover, this will allow us to reuse the new method in other methods. A method with long parameter list can be slimmed down with *Introduce Parameter Object* and *Preserve Whole Object* [37].

Furthermore, we were able to identify 15 methods with the *Long Parameter List* code smell. Three of the methods were listed as critical. As mentioned in the paragraph above, Long Parameter List can be slimmed down with *Introduced Parameter Object*, and *Preserve Whole Object*. *Introduce Parameter Object* should be used when a group of parameters naturally go together, while *Presever Whole Object* should be applied when a whole object can be sent as parameter in a method call rather than several values from an object. In addition, *Replace Parameter with Method* refactoring technique can be applied when you can get the data in one parameter by making a request of an object you already know about [37].

Our results reveal 1153 hits of *Speculative Generality* code smell, including unused methods, local variables, and static globals. Fowler et al. [37] suggest that unused variables and static globals should simply be deleted. Refactoring unused methods can be done by either removing them, or by applying *Inline Method* refactoring technique if a method body is more obvious than the method itself. *Inline Method* will replace calls to the method with the method content and delete the method.

The last code smell we identified is *Dead Code*. This smell includes "commented out" code, unreachable code, and unnecessary "#includes in header files". The quickest way to conduct refactoring on Dead Code code smell is to delete unused code and unneeded files. Notice that this will reduce the number of lines of code.

To summarize, DD can be paid by conducting necessary refactoring. We believe that refactoring will keep the software quality stable, which ultimately mitigates DD issues. However, it is worth noticing that refactoring not necessary is the solution to pay DD. There may be some cases where refactoring is unlikely to reduce fault-proneness in classes, and may increase fault-proneness in a class instead [116]. This phenomenon affects some of the OO-metrics. A consideration should be taken where both metrics and other classes affected by the refactoring are involved. For example, refactoring code smells like *Long Method* and *Duplicated Code* may increase the number of methods and

coupling. Code smells should be managed by prioritizing the most critical smells. If the goal is to maintain and improve a certain metric value, then remove the smell that allows to improve this metric. Furthermore, we do believe that improving the software metrics and other work practices may be better more beneficial than refactoring code smells to reduce the maintenance effort.

5.3 Threats To Validity

Validity is related to how much the results can be trusted [15]. We consider threats to the external, internal, and conclusion validity of this study.

5.3.1 Internal Validity

Interval validity is the degree to which we can conclude that the dependent variable is accounted for by the independent variable [15]. In other words, it concerns any factors that could have influenced our study results without the researcher's knowledge. The internal validity threats can be sorted into three categories [15]: *single group threats*, *multiple group threats*, and *social threats*. *Single group threats* apply to experiments with single groups, which is relevant in our case. One factor is *instrumentation*, that is, the effect caused by the artifacts used for experiment execution [15]. We have used multiple tools to analyze the system. In some cases, these tools may result false positives, which may affect our results.

Another factor that may reduce the internal validity is the selection of subjects from a larger group [15]. We found our subject for the case study after being recommended by others. The internal validity would have been stronger if the subjects had been selected using random sampling.

5.3.2 External Validity

External validity refers to the degree to which the results from our study can be generalized to industrial practice [15]. This study has used one commercial system within a single application domain, so the threat to external validity increases. That said, we cannot generalize the results to other safety-critical systems, both open-source and commercial. More studies need to be conducted on different systems, both commercial and open-source, to compare with the results obtained in this thesis.

5.3.3 Conclusion Validity

Conclusion validity refers to the degree in which correct conclusion can be drawn from the relationship between treatment and the outcome [15]. One important issue here was

the sample size of the experiment. This study has only studied one commercial system, so the statistical power is very low. This is something we are aware of, and therefore more deeper studies need to be performed to confirm if our results have a more general applicability. Another threat may potentially be "*Fishing and error rate*". Searching or "*fishing*" for a specific result is a threat [15], since the researchers may influence the result by looking for a specific outcome. We do not think that we have influenced the results. Reliability of measures may be another threat. Running the tools twice gave us the same results. However, it is not guaranteed that the outcome would be same if another tool had been applied. Another potential problem can be the lack of experience on conducting case studies.

CHAPTER 6

CONCLUSION

DD is an instance of TD which describes the problem of increasing complexity of software design, and deterioration of its maintainability. DD is recognized to be a major problem for many software projects today. As the debt increases, more time will be spent on maintaining the system which means that the software development process and software evolution will become less effective. This will delay other development processes, such as releasing new features. Being able to measure and monitor the quality of software design will help the development team to understand the current situation of having DD by making the debt visible. Once the problems are visible, suitable actions may potentially be performed, which includes refactoring and eventually re-engineering.

This thesis presents the result of a case study that has been conducted in real-life context in collaboration with Autronica Fire and Security AS. The purpose of the study was to investigate DD in safety-critical systems. OO-metrics were measured, and descriptive statistics were computed to analyze and interpret the data. In addition, a set of threshold value were derived to identify classes that are most likely to pose problems for a system. Moreover, automatic static analysis tools were applied to detect code smells in the system. We do believe that these data provide enough information for the project team as a basis for further inspection. It is up to the team to determine the criticality of these classes to make the final determination.

The work contributes mainly to improvement in software metrics and software quality. The stated contributions of this work are:

- **C1:** Empirical knowledge about DD identification in safety-critical systems by OO-metric analysis and code smell detection.
 - **C1.1:** A set of threshold values for the measured OO-metrics.

- **C2:** Empirical knowledge about the effects of having DD in safety-critical systems.
- **C3:** Empirical knowledge about the different types of DD in safety-critical systems.
- **C4:** Empirical knowledge about paying DD.

6.1 Future Work

Based on the results from this thesis, we outline some possibilities for future research.

- We have mainly focused on measuring Chidamber and Kemerer's suite of OO-metrics. A possibility for future research would be to measure other software metrics. For example, cyclomatic complexity and concurrency.
- Implementation of a tool which can detect code smells and suggest possible refactoring options.
- A deeper study on a system by analyzing multiple releases of the system. The evolution of metrics and patterns over time can then be compared to see if the design is decaying or not.
- A study on the use of design patterns, design decay, design grime, and design rot in safety-critical systems.
- A study on differences between metrics thresholds measured in this study and data from other similar systems. For example, it may be possible to compare our results with the measurement of thresholds of open-source systems.
- Naming convention on source files.

BIBLIOGRAPHY

- [1] B. Graaf, M. Lormans, and H. Toetenel, “Embedded software engineering: the state of the practice,” *Software, IEEE*, vol. 20, no. 6, pp. 61–69, 2003.
- [2] P. Ray, R. Laupers, and G. Ascheid, “Compose: A composite embedded software synthesis approach,” in *Innovations in Information Technology (IIT), 2015 11th International Conference on*, pp. 29–34, Nov 2015.
- [3] A. Kyte, “Gartner estimates global it debt to be 500 billion dollars this year, with potential to grow to 1 trillion dollars by 2015,” 2010.
- [4] Z. Li, P. Avgeriou, and P. Liang, “A systematic mapping study on technical debt and its management,” *Journal of Systems and Software*, vol. 101, pp. 193–220, 2015.
- [5] Z. Codabux and B. Williams, “Managing technical debt: An industrial case study,” in *Proceedings of the 4th International Workshop on Managing Technical Debt, MTD ’13*, (Piscataway, NJ, USA), pp. 8–15, IEEE Press, 2013.
- [6] N. Brown, Y. Cai, Y. Guo, R. Kazman, M. Kim, P. Kruchten, E. Lim, A. MacCormack, R. Nord, I. Ozkaya, R. Sangwan, C. Seaman, K. Sullivan, and N. Zazworka, “Managing technical debt in software-reliant systems,” in *Proceedings of the FSE/SDP Workshop on Future of Software Engineering Research, FoSER ’10*, (New York, NY, USA), pp. 47–52, ACM, 2010.
- [7] E. Tom, A. Aurum, and R. Vidgen, “An exploration of technical debt,” *Journal of Systems and Software*, vol. 86, no. 6, pp. 1498–1516, 2013.
- [8] N. Zazworka, C. Seaman, and F. Shull, “Prioritizing design debt investment opportunities,” in *Proceedings of the 2Nd Workshop on Managing Technical Debt, MTD ’11*, (New York, NY, USA), pp. 39–42, ACM, 2011.

- [9] N. Zazworka, R. O. Spínola, A. Vetro', F. Shull, and C. Seaman, "A case study on effectively identifying technical debt," in *Proceedings of the 17th International Conference on Evaluation and Assessment in Software Engineering, EASE '13*, (New York, NY, USA), pp. 42–47, ACM, 2013.
- [10] P. Kruchten, R. Nord, and I. Ozkaya, "Technical debt: From metaphor to theory and practice," *Software, IEEE*, vol. 29, pp. 18–21, Nov 2012.
- [11] R. Mo, J. Garcia, Y. Cai, and N. Medvidovic, "Mapping architectural decay instances to dependency models," in *Proceedings of the 4th International Workshop on Managing Technical Debt*, pp. 39–46, IEEE Press, 2013.
- [12] C. Izurieta and J. M. Bieman, "How software designs decay: A pilot study of pattern evolution," in *Empirical Software Engineering and Measurement, 2007. ESEM 2007. First International Symposium on*, pp. 449–451, IEEE, 2007.
- [13] S. K. Bhuiyan, "Managing technical debt in embedded systems," 2015.
- [14] B. J. Oates, *Researching Information Systems and Computing*. Sage Publications Ltd., 2006.
- [15] C. Wohlin, P. Runeson, M. Höst, M. C. Ohlsson, B. Regnell, and A. Wesslén, *Experimentation in Software Engineering: An Introduction*. Norwell, MA, USA: Kluwer Academic Publishers, 2000.
- [16] M. Bassey, "Case study research," *Educational Research in Practice*, pp. 111–123, 2003.
- [17] S. K. Soy, "The case study as a research method," 1997.
- [18] W. Cunningham, "The wycash portfolio management system," *SIGPLAN OOPS Mess.*, vol. 4, pp. 29–30, Dec. 1992.
- [19] E. Allman, "Managing technical debt," *Commun. ACM*, vol. 55, pp. 50–55, May 2012.
- [20] Y. Guo and C. Seaman, "A portfolio approach to technical debt management," in *Proceedings of the 2Nd Workshop on Managing Technical Debt, MTD '11*, (New York, NY, USA), pp. 31–34, ACM, 2011.
- [21] T. Klinger, P. Tarr, P. Wagstrom, and C. Williams, "An enterprise perspective on technical debt," in *Proceedings of the 2Nd Workshop on Managing Technical Debt, MTD '11*, (New York, NY, USA), pp. 35–38, ACM, 2011.
- [22] R. L. Nord, I. Ozkaya, P. Kruchten, and M. Gonzalez-Rojas, "In search of a metric for managing architectural technical debt," in *Software Architecture (WICSA) and European Conference on Software Architecture (ECSA), 2012 Joint Working IEEE/IFIP Conference on*, pp. 91–100, IEEE, 2012.
- [23] M. Fowler, "Technical Debt Quadrant," 2009.
- [24] S. McConnell, "Technical debt," 2007.

- [25] S. McConnell, "Managing technical debt." <http://2013.icse-conferences.org/documents/publicity/MTD-WS-McConnell-slides.pdf>, 2007. [Online; accessed 03-May-2016].
- [26] E. Lim, N. Taksande, and C. Seaman, "A balancing act: What software practitioners have to say about technical debt," *Software, IEEE*, vol. 29, pp. 22–27, Nov 2012.
- [27] Y. Guo, C. Seaman, R. Gomes, A. Cavalcanti, G. Tonin, F. Q. Da Silva, A. L. Santos, and C. Siebra, "Tracking technical debt—an exploratory case study," in *Software Maintenance (ICSM), 2011 27th IEEE International Conference on*, pp. 528–531, IEEE, 2011.
- [28] C. S. A. Siebra, G. S. Tonin, F. Q. Silva, R. G. Oliveira, A. L. Junior, R. C. Miranda, and A. L. Santos, "Managing technical debt in practice: An industrial report," in *Proceedings of the ACM-IEEE International Symposium on Empirical Software Engineering and Measurement, ESEM '12*, (New York, NY, USA), pp. 247–250, ACM, 2012.
- [29] N. Zazworka, M. A. Shaw, F. Shull, and C. Seaman, "Investigating the impact of design debt on software quality," in *Proceedings of the 2nd Workshop on Managing Technical Debt*, pp. 17–23, ACM, 2011.
- [30] F. Buschmann, "To pay or not to pay technical debt," *Software, IEEE*, vol. 28, no. 6, pp. 29–31, 2011.
- [31] N. Zazworka, C. Izurieta, S. Wong, Y. Cai, C. Seaman, F. Shull, *et al.*, "Comparing four approaches for technical debt identification," *Software Quality Journal*, vol. 22, no. 3, pp. 403–426, 2014.
- [32] S. Huynh and Y. Cai, "An evolutionary approach to software modularity analysis," in *Proceedings of the First International Workshop on Assessment of Contemporary Modularization Techniques*, p. 6, IEEE Computer Society, 2007.
- [33] S. Wong, Y. Cai, M. Kim, and M. Dalton, "Detecting software modularity violations," in *Proceedings of the 33rd International Conference on Software Engineering*, pp. 411–420, ACM, 2011.
- [34] SourceMaking, "Design patterns." https://sourcemaking.com/design_patterns. (Accessed on 06/07/2016).
- [35] C. Izurieta and J. M. Bieman, "A multiple case study of design pattern decay, grime, and rot in evolving software systems," *Software Quality Journal*, vol. 21, no. 2, pp. 289–323, 2013.
- [36] C. Izurieta and J. M. Bieman, "Testing consequences of grime buildup in object oriented design patterns," in *Software Testing, Verification, and Validation, 2008 1st International Conference on*, pp. 171–179, IEEE, 2008.

- [37] M. Fowler, K. Beck, J. Brant, W. Opdyke, and D. Roberts, “Refactoring: Improving the design of existing programs,” 1999.
- [38] M. Mäntylä, J. Vanhanen, and C. Lassenius, “A taxonomy and an initial empirical study of bad smells in code,” in *Software Maintenance, 2003. ICSM 2003. Proceedings. International Conference on*, pp. 381–384, IEEE, 2003.
- [39] S. Olbrich, D. S. Cruzes, V. Basili, and N. Zazworka, “The evolution and impact of code smells: A case study of two open source systems,” in *Proceedings of the 2009 3rd international symposium on empirical software engineering and measurement*, pp. 390–400, IEEE Computer Society, 2009.
- [40] F. Khomh, M. D. Penta, and Y.-G. Gueheneuc, “An exploratory study of the impact of code smells on software change-proneness,” in *Reverse Engineering, 2009. WCRE’09. 16th Working Conference on*, pp. 75–84, IEEE, 2009.
- [41] G. Travassos, F. Shull, M. Fredericks, and V. R. Basili, “Detecting defects in object-oriented designs: using reading techniques to increase software quality,” in *ACM Sigplan Notices*, vol. 34, pp. 47–56, ACM, 1999.
- [42] R. Marinescu, “Detecting design flaws via metrics in object-oriented systems,” in *Technology of Object-Oriented Languages and Systems, 2001. TOOLS 39. 39th International Conference and Exhibition on*, pp. 173–182, IEEE, 2001.
- [43] O. Ciupke, “Automatic detection of design problems in object-oriented reengineering,” in *Technology of Object-Oriented Languages and Systems, 1999. TOOLS 30 Proceedings*, pp. 18–32, IEEE, 1999.
- [44] R. Marinescu, “Detection strategies: Metrics-based rules for detecting design flaws,” in *Software Maintenance, 2004. Proceedings. 20th IEEE International Conference on*, pp. 350–359, IEEE, 2004.
- [45] J. Schumacher, N. Zazworka, F. Shull, C. Seaman, and M. Shaw, “Building empirical support for automated code smell detection,” in *Proceedings of the 2010 ACM-IEEE International Symposium on Empirical Software Engineering and Measurement*, p. 8, ACM, 2010.
- [46] N. Rutar, C. B. Almazan, and J. S. Foster, “A comparison of bug finding tools for java,” in *Software Reliability Engineering, 2004. ISSRE 2004. 15th International Symposium on*, pp. 245–256, IEEE, 2004.
- [47] T. Copeland, “Pmd applied,” 2005.
- [48] Z. Codabux and B. J. Williams, “Technical debt prioritization using predictive analytics,” in *Proceedings of the 38th International Conference on Software Engineering Companion*, pp. 704–706, ACM, 2016.
- [49] G. Suryanarayana, G. Samarthayam, and T. Sharma, *Refactoring for Software Design Smells: Managing Technical Debt*. Morgan Kaufmann, 2014.
- [50] S. SonarSource, “Sonarqube,” tech. rep., Technical report, last update: June, 2013.

- [51] J. Radatz, A. Geraci, and F. Katki, “IEEE Standard Glossary of Software Engineering Terminology,” *IEEE Std*, vol. 610121990, no. 121990, p. 3, 1990.
- [52] ISO/IEC, *ISO/IEC 9126. Software engineering – Product quality*. ISO/IEC, 2001.
- [53] L. Bass, P. Clements, and R. Kazman, *Software Architecture in Practice*. Addison-Wesley Professional, 3rd ed., 2012.
- [54] G. Quenel and H. Lövdahl, “Object oriented software quality models,” *membres.multimania.fr/gquenel/site/files/doc/OOSoftQualMod.pdf*.
- [55] S. R. Chidamber and C. F. Kemerer, “A metrics suite for object oriented design,” *Software Engineering, IEEE Transactions on*, vol. 20, no. 6, pp. 476–493, 1994.
- [56] V. R. Basili, L. C. Briand, and W. L. Melo, “A validation of object-oriented design metrics as quality indicators,” *Software Engineering, IEEE Transactions on*, vol. 22, no. 10, pp. 751–761, 1996.
- [57] W. Li and S. Henry, “Object-oriented metrics that predict maintainability,” *Journal of systems and software*, vol. 23, no. 2, pp. 111–122, 1993.
- [58] I. Sommerville, *Software Engineering (9th Edition)*. Pearson Addison Wesley, 2011.
- [59] M. M. Lehman, “Programs, life cycles, and laws of software evolution,” *Proceedings of the IEEE*, vol. 68, no. 9, pp. 1060–1076, 1980.
- [60] K. H. Bennett and V. T. Rajlich, “Software maintenance and evolution: A roadmap,” in *Proceedings of the Conference on The Future of Software Engineering*, ICSE ’00, (New York, NY, USA), pp. 73–87, ACM, 2000.
- [61] “IEEE Standard for Software Maintenance,” *IEEE Std 1219-1998*, 1998.
- [62] H. v. Vliet, *Software Engineering: Principles and Practice*. Wiley Publishing, 3rd ed., 2008.
- [63] B. P. Lientz and E. B. Swanson, “Software maintenance management: a study of the maintenance of computer application software in 487 data processing organizations,” 1980.
- [64] C. W. Krueger, “Software reuse,” *ACM Computing Surveys (CSUR)*, vol. 24, no. 2, pp. 131–183, 1992.
- [65] W. Frakes and C. Terry, “Software reuse: metrics and models,” *ACM Computing Surveys (CSUR)*, vol. 28, no. 2, pp. 415–435, 1996.
- [66] J. Sametinger, *Software engineering with reusable components*. Springer Science & Business Media, 1997.
- [67] W. Frakes, “Systematic software reuse: a paradigm shift,” in *Software Reuse: Advances in Software Reusability, 1994. Proceedings., Third International Conference on*, pp. 2–3, IEEE, 1994.

- [68] O. P. N. Slyngstad, A. Gupta, R. Conradi, P. Mohagheghi, H. Rønneberg, and E. Landre, “An empirical study of developers views on software reuse in statoil asa,” in *Proceedings of the 2006 ACM/IEEE International Symposium on Empirical Software Engineering*, ISESE '06, (New York, NY, USA), pp. 242–251, ACM, 2006.
- [69] W. C. Lim, “Effects of reuse on quality, productivity, and economics,” *Software, IEEE*, vol. 11, no. 5, pp. 23–30, 1994.
- [70] M. Al Mamun, C. Berger, and J. Hansson, “Explicating, understanding, and managing technical debt from self-driving miniature car projects,” in *Managing Technical Debt (MTD)*, 2014 Sixth International Workshop on, pp. 11–18, Sept 2014.
- [71] M. Morisio, M. Ezran, and C. Tully, “Success and failure factors in software reuse,” *Software Engineering, IEEE Transactions on*, vol. 28, pp. 340–357, Apr 2002.
- [72] W. B. Frakes and S. Isoda, “Success factors of systematic reuse,” *Software, IEEE*, vol. 11, no. 5, pp. 14–19, 1994.
- [73] M. Kim, T. Zimmermann, and N. Nagappan, “A field study of refactoring challenges and benefits,” in *Proceedings of the ACM SIGSOFT 20th International Symposium on the Foundations of Software Engineering*, FSE '12, (New York, NY, USA), pp. 50:1–50:11, ACM, 2012.
- [74] W. F. Opdyke, *Refactoring object-oriented frameworks*. PhD thesis, University of Illinois at Urbana-Champaign, 1992.
- [75] P. Koopman, “Embedded systems in the real world,” *Course Material: Dependable Embedded Systems*, 1999.
- [76] I. Crnkovic, “Component-based approach for embedded systems,” in *Ninth International Workshop on Component-Oriented Programming (WCOP)*, 2004.
- [77] Z. You, “The reliability analysis of embedded systems,” in *Information Science and Cloud Computing Companion (ISCC-C)*, 2013 International Conference on, pp. 458–462, IEEE, 2013.
- [78] S. Patil and L. Kapaleshwari, “Embedded software-issues and challenges,” tech. rep., SAE Technical Paper, 2009.
- [79] C. Ebert and C. Jones, “Embedded software: Facts, figures, and future,” *Computer*, no. 4, pp. 42–52, 2009.
- [80] T. Sherman, “Quality attributes for embedded systems,” in *Advances in Computer and Information Sciences and Engineering*, pp. 536–539, Springer, 2008.
- [81] A. Pretschner, M. Broy, I. H. Kruger, and T. Stauner, “Software engineering for automotive systems: A roadmap,” in *2007 Future of Software Engineering*, pp. 55–71, IEEE Computer Society, 2007.

- [82] J. J. Trienekens, R. J. Kusters, and D. C. Brussel, "Quality specification and metrikation, results from a case-study in a mission-critical software domain," *Software Quality Journal*, vol. 18, no. 4, pp. 469–490, 2010.
- [83] J. C. Knight, "Safety critical systems: challenges and directions," in *Software Engineering, 2002. ICSE 2002. Proceedings of the 24rd International Conference on*, pp. 547–550, IEEE, 2002.
- [84] H. Washizaki, Y. Kobayashi, H. Watanabe, E. Nakajima, Y. Hagiwara, K. Hiranabe, and K. Fukuda, "Quality evaluation of embedded software in robot software design contest," *Progress in Informatics*, vol. 5, pp. 35–47, 2007.
- [85] M. A. B. Sarfraz Nawaz Brohi, "Empirical research methods for software engineering," 2001.
- [86] P. Runeson and M. Höst, "Guidelines for conducting and reporting case study research in software engineering," *Empirical Softw. Engg.*, vol. 14, pp. 131–164, Apr. 2009.
- [87] R. K. Yin, "Case study research: Design and methods. volume 5," 2003.
- [88] Doxygen, "Doxygen." <http://www.stack.nl/~dimitri/doxygen/>. (Accessed on 06/04/2016).
- [89] Tigris.org, "Welcome to argouml," 2001.
- [90] S. Systems, "The ultimate design & build platform." <http://www.sparxsystems.com.au/>. (Accessed on 06/03/2016).
- [91] Scitools, "Understand." <https://scitools.com/>. (Accessed on 06/03/2016).
- [92] CppDepend, "Cppdepend :: C/c++ static analysis and code quality tool." <http://www.cppdepend.com/>. (Accessed on 06/03/2016).
- [93] CppCheck, "Cppcheck - a tool for static c/c++ code analysis." <http://cppcheck.sourceforge.net/>. (Accessed on 06/03/2016).
- [94] CppClean, "cppclean." <https://github.com/myint/cppclean>. (Accessed on 06/03/2016).
- [95] T. Littlefair, "Cccc - c and c++ code counter." <https://scitools.com/>. (Accessed on 06/03/2016).
- [96] L. Rosenberg, R. Stapko, and A. Gallo, "Risk-based object oriented testing," *24th SWE*, 1999.
- [97] D. Rodriguez and R. Harrison, "An overview of object-oriented design metrics," 2001.
- [98] T. G. Filo, M. A. Bigonha, and K. A. Ferrira, "A catalogue of thresholds for object-oriented software metrics," 2015.

- [99] K. A. Ferreira, M. A. Bigonha, R. S. Bigonha, L. F. Mendes, and H. C. Almeida, "Identifying thresholds for object-oriented software metrics," *Journal of Systems and Software*, vol. 85, no. 2, pp. 244–257, 2012.
- [100] M. Lanza and R. Marinescu, *Object-oriented metrics in practice: using software metrics to characterize, evaluate, and improve the design of object-oriented systems*. Springer Science & Business Media, 2007.
- [101] *Refactoring: Improving the Design of Existing Code*. Boston, MA, USA: Addison-Wesley Longman Publishing Co., Inc., 1999.
- [102] Š. Cais and P. Pícha, "Identifying software metrics thresholds for safety critical system," 2014.
- [103] A. Brooks, "Metrics overview."
- [104] Objectteering, "Objectteering metrics user guide."
- [105] H. A. Sahraoui, R. Godin, and T. Miceli, "Can metrics help to bridge the gap between the improvement of oo design quality and its automation?," in *Software Maintenance, 2000. Proceedings. International Conference on*, pp. 154–162, IEEE, 2000.
- [106] PHPDepend, "Software metrics supported by phpdepend."
- [107] R. Subramanyam and M. S. Krishnan, "Empirical analysis of ck metrics for object-oriented design complexity: Implications for software defects," *Software Engineering, IEEE Transactions on*, vol. 29, no. 4, pp. 297–310, 2003.
- [108] J. Daly, A. Brooks, J. Miller, M. Roper, and M. Wood, "Evaluating inheritance depth on the maintainability of object-oriented software," *Empirical Software Engineering*, vol. 1, no. 2, pp. 109–132, 1996.
- [109] K. El Emam, W. Melo, and J. C. Machado, "The prediction of faulty classes using object-oriented design metrics," *Journal of Systems and Software*, vol. 56, no. 1, pp. 63–75, 2001.
- [110] E. Okike, "A pedagogical evaluation and discussion about the lack of cohesion in method (lcom) metric using field experiment," *arXiv preprint arXiv:1004.3277*, 2010.
- [111] A. A. Bakar and N. Sham, "The analysis of object-oriented metrics in c++ programs," *Lecture Notes on Software Engineering*, vol. 4, no. 1, pp. 48–52, 2014.
- [112] R. Malhotra and A. J. Bansal, "Fault prediction considering threshold effects of object-oriented metrics," *Expert Systems*, vol. 32, no. 2, pp. 203–219, 2015.
- [113] L. H. Rosenberg, "Applying and interpreting object-oriented metrics," 1998.
- [114] J. Bansiya and C. G. Davis, "A hierarchical model for object-oriented design quality assessment," *Software Engineering, IEEE Transactions on*, vol. 28, no. 1, pp. 4–17, 2002.

- [115] D. I. Sjöberg, A. Yamashita, B. C. D. Anda, A. Mockus, and T. Dyba, “Quantifying the effect of code smells on maintenance effort,” *Software Engineering, IEEE Transactions on*, vol. 39, no. 8, pp. 1144–1156, 2013.
- [116] T. Hall, M. Zhang, D. Bowes, and Y. Sun, “Some code smells have a significant but small effect on faults,” *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 23, no. 4, p. 33, 2014.
- [117] J. Lindsay, J. Noble, and E. Tempero, “Does size matter?: a preliminary investigation of the consequences of powerlaws in software,” in *Proceedings of the 2010 ICSE Workshop on Emerging Trends in Software Metrics*, pp. 16–23, ACM, 2010.
- [118] W. Li and R. Shatnawi, “An empirical study of the bad smells and class error probability in the post-release object-oriented system evolution,” *Journal of systems and software*, vol. 80, no. 7, pp. 1120–1128, 2007.
- [119] P. K. Dhillon and G. Sidhu, “Can software faults be analyzed using bad code smells?: An empirical study,” *International Journal of Scientific and Research Publications*, vol. 2, no. 10, pp. 1–7, 2012.
- [120] S. M. Olbrich, D. S. Cruze, and D. I. Sjøberg, “Are all code smells harmful? a study of god classes and brain classes in the evolution of three open source systems,” in *Software Maintenance (ICSM), 2010 IEEE International Conference on*, pp. 1–10, IEEE, 2010.
- [121] R. Moser, A. Sillitti, P. Abrahamsson, and G. Succi, “Does refactoring improve reusability?,” in *Reuse of Off-the-Shelf Components*, pp. 287–297, Springer, 2006.