

Bachelor's Thesis

# COMMIT-FEATURE INTERACTIONS: ANALYZING STRUCTURAL AND DATAFLOW RELATIONS BETWEEN COMMITTS AND FEATURES

SIMON STEUER

October 3, 2023

Advisor:

Sebastian Böhm Chair of Software Engineering

Examiners:

Prof. Dr. Sven Apel

Chair of Software Engineering

Andreas Zeller

CISPA Helmholtz Center for Information Security

Chair of Software Engineering  
Saarland Informatics Campus  
Saarland University



UNIVERSITÄT  
DES  
SAARLANDES



## **Erklärung**

Ich erkläre hiermit, dass ich die vorliegende Arbeit selbständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel verwendet habe.

## **Statement**

I hereby confirm that I have written this thesis on my own and that I have not used any other media or materials than the ones referred to in this thesis

## **Einverständniserklärung**

Ich bin damit einverstanden, dass meine (bestandene) Arbeit in beiden Versionen in die Bibliothek der Informatik aufgenommen und damit veröffentlicht wird.

## **Declaration of Consent**

I agree to make both versions of my thesis (with a passing grade) accessible to the public by having them added to the library of the Computer Science Department.

Saarbrücken, \_\_\_\_\_  
(Datum/Date)

\_\_\_\_\_  
(Unterschrift/Signature)



## ABSTRACT

---

Short summary of the contents in English...a great guide by Kent Beck how to write good abstracts can be found here:

<https://plg.uwaterloo.ca/~migod/research/beck00PSLA.html>



## CONTENTS

---

1	Introduction	1
1.1	Goal of this Thesis . . . . .	1
1.2	Overview . . . . .	1
2	Background	3
2.1	Code Regions . . . . .	3
2.2	Interaction Analysis . . . . .	4
3	Commit-Feature Interactions	7
3.1	Structural CFIs . . . . .	7
3.2	Dataflow-based CFIs . . . . .	8
3.3	Combination of CFIs . . . . .	8
3.4	Feature Size . . . . .	9
3.5	Implementation . . . . .	9
4	Methodology	11
4.1	Research Questions . . . . .	11
5	Evaluation	13
5.1	Results . . . . .	13
5.2	Discussion . . . . .	13
5.3	Threats to Validity . . . . .	13
6	Related Work	15
7	Concluding Remarks	17
7.1	Conclusion . . . . .	17
7.2	Future Work . . . . .	17
A	Appendix	19
	Bibliography	21

## LIST OF FIGURES

---

Figure 4.1	Kinds of commits in software projects investigated in this work. In the first two <b>RQs</b> we have discussed different kinds of commits and the ways in which they interact with features. Figure 1 showcases them in a venn diagram and illustrates the dependencies and divisions between them. . . . .	12
------------	---	----

## LIST OF TABLES

---

## LISTINGS

---

Listing 3.1	This code example contains both structural as well as dataflow-based commit feature interactions. Commit <code>fc3a17d</code> implements the functionality of <code>FeatureDouble</code> for this function. It follows that a structural commit-feature interaction can be found between them, as their respective commit and feature regions structurally interact. Commit <code>7edb283</code> introduces the variable <code>ret</code> that is later used inside the feature region of <code>FeatureDouble</code> . This accounts for a commit-feature interaction through dataflow, as data that was produced within a commit region is used as input by an instruction belonging to a feature region of <code>FeatureDouble</code> later on in the program. . . . .	8
-------------	--	---

## ACRONYMS

---



# INTRODUCTION

---

## 1.1 GOAL OF THIS THESIS

## 1.2 OVERVIEW



## BACKGROUND

---

In this section, we summarize previous research on the topic of code regions and interaction analysis. Thereby we give definitions of terms and discuss concepts that are fundamental to our work. In the next chapter, we use the introduced definitions and concepts to explain and investigate structural as well as dataflow-based commit-feature interactions.

### 2.1 CODE REGIONS

Software Programs consist of source code lines that are translated into an intermediate representation (IR) upon compilation. IR accurately represents the source-code information and is utilized in many code improvement and transformation techniques such as code-optimization. Furthermore, static program analyses are conducted on top of IR or some sort of representation using IR.

Code regions are comprised of IR instructions that are consecutive in their control-flow. They are used to represent abstract entities of a software project, such as commits and features. For this, each code region carries some kind of variability information, detailing its meaning. It should be noted that there can be several code regions with the same meaning scattered across the program. We discuss two kinds of variability information, namely commit and feature variability information, allowing us to define two separate code regions.

Commits are used within a version control system to represent the latest source code changes in its respective repository. Inside a repository revision, a commit encompasses all source code lines that were added or last changed by it.

**Definition 1** (commit regions). The set of all consecutive instructions, that stem from source code lines belonging to a commit, is called a *commit region*.

As the source code lines changed by a commit are not necessarily contiguous, a commit can have many commit regions inside a program. To properly address the entire set of these commit regions within a program, we introduce the term regions of a commit. We say that the *regions of a commit* are the set of all commit regions that stem from source code lines of the commit.

In general, features are program parts implementing specific functionality. In this work we focus on features that are modelled with the help of configuration variables that specify whether the functionality of a feature should be active or not. Inside a program, configuration variables decide whether instructions, performing a feature's intended functionality, get executed. Detecting these control-flow dependencies is achieved by deploying extended static taint analyses, such as Lotrack[5].

**Definition 2** (feature regions). The set of all consecutive instructions, whose execution depends on a configuration variable belonging to a feature, is called a *feature region*.

Similarly to commits, features can have many feature regions inside a program as the instructions implementing their functionalities can be scattered across the program. Here, we also say that the *regions of a feature* are the set of all feature regions that stem from a feature’s configuration variable.

## 2.2 INTERACTION ANALYSIS

In the previous chapter, we have already shown how different abstract entities of a software project, such as features and commits, can be assigned concrete representations within a program. We can use interaction analyses to infer interactions between them by computing interactions between their concrete representations. This can improve our understanding of them and give us insights on how they are used inside software projects. An important component of computing interactions between code regions is the concept of interaction relations between them introduced by Sattler [6]. As we investigate structural and dataflow-based interactions, we make use of structural interaction ( $\odot$ ) and dataflow interaction ( $\rightsquigarrow$ ) relations in this work.

The structural interaction relation between two code regions  $r_1$  and  $r_2$  is defined as follows:

**Definition 3** (structural interaction relation). The interaction relation  $r_1 \odot r_2$  evaluates to true if at least one instruction that is part of  $r_1$  is also part of  $r_2$ .

The dataflow interaction relation between two code regions  $r_1$  and  $r_2$  is defined as:

**Definition 4** (dataflow interaction relation). The interaction relation  $r_1 \rightsquigarrow r_2$  evaluates to true if the data produced by at least one instruction  $i$  that is part of  $r_1$  flows as input to an instruction  $i'$  that is part of  $r_2$ .

Essential to the research conducted in this paper is the interaction analysis created by Sattler [6]. Their interaction analysis tool VaRA is implemented on top of LLVM and PhASAR. In their novel approach SEAL[8], VaRA is used to determine dataflow interactions between commits. This is accomplished by computing dataflow interactions between the respective regions of commits. Their approach for this will be shortly discussed here, however we advice their paper for a more thorough explanation.

The first step of their approach is to annotate code by mapping information about its regions to the compiler’s intermediate representation. This information is added to the LLVM-IR instructions during the construction of the IR. In their paper, Sattler et al. [8] focus on commit regions, which contain information about the commit’s hash and its respective repository. The commit region of an instruction is extracted from the commit that last changed the source code line the instruction stems from. Determining said commit is accomplished by accessing repository meta-data.

The second analysis step involves the actual computation of the interactions. For this, Sattler et al. implemented a special, inter-procedural taint analysis. It’s able to track data flows between the commit regions of a given target program. For this, information about which commit region affected it, is mapped onto data and tracked along program flow. In this context, we would like to introduce the concept of taints, specifically commit taints. Taints are used to give information about which code regions have affected an instruction through dataflow. For example, an instruction is tainted by a commit if its taint stems from a commit region. It

follows, that two commits, with their respective commit regions  $r_1$  and  $r_2$ , interact through dataflow, when an instruction is part of one's commit region,  $r_2$ , while being tainted by  $r_1$ . Consequently data produced by the commit region  $r_1$  flows as input to an instruction within the commit region  $r_2$ , matching the [dataflow interaction relation](#).



## COMMIT-FEATURE INTERACTIONS

---

In this section, we define structural and dataflow-based commit-feature interactions as well as properties related to them. Furthermore, their meaning and relationship inside a software project is explained here. In the [Overview](#) chapter, we discussed what purpose commits and features serve in a software project. Commits are used to add new changes, whereas features are cohesive entities in a program implementing a specific functionality. In this work, structural interactions are used to investigate how commits implement features and their functionality. In addition, dataflow interactions are examined to gain additional knowledge on how new changes to a program, in the form of commits, affect features.

### 3.1 STRUCTURAL CFIS

In the background chapter, we discussed the concept of [Code Regions](#), especially commit and feature regions. Logically, we speak of structural interactions between features and commits when their respective regions structurally interact. This structural interaction between code regions occurs when at least one instruction is part of both regions. This is the case when code regions structurally interact through the [structural interaction relation](#) ( $\odot$ ).

**Definition 5.** A commit  $C$  with its commit regions  $r_{1C}, r_{2C}, \dots$  and a feature  $F$  with its feature regions  $r_{1F}, r_{2F}, \dots$  structurally interact, if at least one commit region  $r_{iC}$  and one feature region  $r_{jF}$  structurally interact with each other, i.e.  $r_{iC} \odot r_{jF}$ .

Structural CFIs carry an important meaning, namely that the commit of the interaction was used to implement or change functionality of the feature of said interaction. This can be seen when looking at an instruction accounting for a structural interaction, as it is both part of a commit as well as a feature region. From the definition of [commit regions](#), it follows that the instruction stems from a source-code line that was last changed by the region's respective commit. From the definition of [feature regions](#), we also know that the instruction implements functionality of the feature region's respective feature. Thus, the commit of a structural interaction was used to extend or change the code implementing the feature of that interaction. Following this, we can say that the commits, a feature structurally interacts with, implement the entire functionality of the feature. That is because each source-code line of a git-repository was introduced by a commit and can only belong to a single commit. Thus every instruction, including those part of feature regions, is annotated by exactly one commit region.

Knowing the commits used to implement a feature allows us to determine the authors that developed it. This is made possible by simply linking the commits, a feature structurally interacts with, to their respective authors. By determining the authors of a feature, we can achieve a deeper insight into its development than solely focusing on the commits that implemented it.

### 3.2 DATAFLOW-BASED CFIS

Determining which commits affect a feature through dataflow can reveal additional interactions between commits and features that cannot be discovered with a structural analysis. Especially dataflows that span over multiple files and many lines of code might be difficult for a programmer to be aware of. Employing VaRA's dataflow analysis, that is discussed in section 3.5, facilitates the detection of these dataflow interactions.

Commit interactions based on dataflow were explained in the [Interaction Analysis](#) section and can be considered as precursors to dataflow-based commit-feature interactions. Similarly to commits interacting with other commits through dataflow, commits interact with features through dataflow, when there exists dataflow from a commit to a feature region. This means that data allocated or changed within a commit region flows as input to an instruction located inside a feature region. This pattern can also be matched to the [dataflow interaction relation](#) ( $\rightsquigarrow$ ) when defining dataflow-based commit-feature interactions.

**Definition 6.** A commit  $C$  with its commit regions  $r_{1C}, r_{2C}, \dots$  and a feature  $F$  with its feature regions  $r_{1F}, r_{2F}, \dots$  interact through dataflow, if at least one commit region  $r_{iC}$  interacts with a feature region  $r_{jF}$  through dataflow, i.e.  $r_{iC} \rightsquigarrow r_{jF}$ .

---

1. <code>int calc(int val) {</code>	▷ d93df4a	
2. <code>int ret = val + 5;</code>	▷ 7edb283	
3. <code>if (FeatureDouble) {</code>	▷ fc3a17d	▷ FeatureDouble
4. <code>ret = ret * 2;</code>	▷ fc3a17d	▷ FeatureDouble
5. <code>}</code>	▷ fc3a17d	▷ FeatureDouble
6. <code>return ret;</code>	▷ d93df4a	
7. <code>}</code>	▷ d93df4a	

---

Listing 3.1: This code example contains both structural as well as dataflow-based commit feature interactions. Commit `fc3a17d` implements the functionality of `FeatureDouble` for this function. It follows that a structural commit-feature interaction can be found between them, as their respective commit and feature regions structurally interact. Commit `7edb283` introduces the variable `ret` that is later used inside the feature region of `FeatureDouble`. This accounts for a commit-feature interaction through dataflow, as data that was produced within a commit region is used as input by an instruction belonging to a feature region of `FeatureDouble` later on in the program.

### 3.3 COMBINATION OF CFIS

When investigating dataflow-based CFIs, it is important to be aware of the fact that structural CFIs heavily coincide with them. This means that whenever commits and features structurally interact, they are likely to interact through dataflow as well. As our structural analysis has already discovered that these commits and features interact with each other, we are more interested in commit-feature interactions that can only be detected by a dataflow analysis. That this relationship between structural and dataflow interactions exists becomes clear when looking at an instruction accounting for a structural CFI. From definition 5, we know that the instruction belongs to a commit region of the interaction's respective commit. It follows that data changed inside the instruction produces commit taints for instructions that use the data



as input. Now, if instructions that use the data as input are also part of a feature region of the interaction's respective feature, the commit and feature of the structural interaction will also interact through dataflow. However, such dataflow is very likely to occur, as features are functional units, whose instructions build and depend upon each other. Knowing this we can differentiate between dataflow-based CFIs that occur within the regions of a feature and those where data flows from outside the regions of a feature into it. From prior explanations, it follows that this differentiation can be accomplished by simply checking whether a commit that influences a feature through dataflow, also structurally interacts with it.

### 3.4 FEATURE SIZE

When examining commit-feature interactions in a project, it is helpful to have a measure that can estimate the size of a feature. We can use such a measure to compare features with each other and, thus, put the number of their interactions into perspective. Considering our implementation, it makes most sense to define the size of a feature as the number of instructions implementing its functionality inside a program. As the instructions inside the regions of a feature implement its functionality, we can define the size of a feature as follows:

**Definition 7.** The *size* of a feature is the number of instructions that are part of its feature regions.

It's possible to calculate the defined size of a feature by calculating the number of instructions in which structural CFIs occur. That is, because every instruction that is part of a feature region accounts for a structural commit-feature interaction, as every instruction is part of exactly one commit region as shown in the beginning of this section. It follows, that we do not miss any instructions that are part of feature regions and do not count any such instruction more than once.

### 3.5 IMPLEMENTATION

The detection of structural as well as dataflow-based commit-feature interactions is implemented in VaRA [7]. Additionally to commit regions, VaRA maps information about its feature regions onto the compiler's IR during its construction. Commit regions contain the hash and repository of their respective commits, whereas feature regions contain the name of the feature they originated from. VaRA also gives us access to every llvm-IR instruction of a program and its attached information. Thus, structural CFIs of a program can be collected by examining its compiled instructions. According to definition 5, we can store a structural interaction between a commit and a feature, if an instruction is part of a respective commit and feature region. For each such interaction we also save the number of instructions it occurs in. This is accomplished by incrementing its instruction counter if we happen to encounter a duplicate.

With this, we are also able to calculate the size of a feature based on our collected structural CFIs and their respective instruction counters. Following the explanations from section 3.4, we can compute the size of a feature as the sum over the instruction counters of all found structural CFIs the feature is part of.

In the [Interaction Analysis](#) section, we discussed the taint analysis deployed by VaRA. There, VaRA computes information about which code regions have affected an instruction through dataflow. Checking whether a taint stems from a commit region allows us to extract information about which commits have tainted an instruction. Thus, dataflow-based commit-feature interactions can also be collected on instruction level. According to definition 6, we can store a dataflow-based interaction between a commit and a feature, if an instruction has a respective commit taint while belonging to a respective feature region. Consequently said instruction uses data, that was changed by a commit region earlier in the program, as its input, while belonging to a feature region.

## METHODOLOGY

---

The purpose of this chapter is to first formulate the research questions that we examine in our work and then propose our method of answering them.

### 4.1 RESEARCH QUESTIONS

In the [Interaction Analysis](#) chapter we discussed the different meanings of structural and dataflow-based interactions. With this knowledge, we can answer many interesting research topics. These topics include patterns in feature development and usage of commits therein as well as findings about how likely seemingly unrelated commits are to affect features inside a program.

#### *RQ1: How do commits and features structurally interact with each other?*

We intend to research two main properties which already provide a lot of insight into the development process of features and best practices of commits therein. Firstly, we examine the amount of commits features interact with structurally. This gives us a direct estimate on how many commits were used in the development of a feature. Our analysis also allows us to measure the size of feature, which can put the amount of commits used to implement a feature into perspective. Secondly, we want to examine how many features a commit interacts with structurally, e.g. how many features a commit usually changes. This is especially interesting when considering best practices surrounding the usage of commits. It is preferred to keep commits atomic[2] meaning they should only deal with a single concern. As different features implement separate functionalities, it's unlikely for a commit to change several features while dealing with the same concern. Transferring this to our work, high quality commits should mostly change a single feature. Acquiring data on this issue might show how strictly this policy is enforced in the development of features across different projects.

#### *RQ2: How do commits interact with features through dataflow?*

Investigating dataflow can unveil interactions between parts of a program that were previously hidden from programmers. This can help a programmer understand the extent to which new changes affect other parts of a program. Deploying the introduced analysis in a direct manner could even aid a programmer when fixing bugs of features. Bugs occurring in certain features could be traced back to the authors responsible for them by factoring in recent commits affecting said features through dataflow.

Previous research has laid the groundwork for researching dataflow interactions between different parts of a program. However, it has focused solely on dataflow interactions between commits. That's why we aim to provide first insights on the properties of dataflow-based CFIs. Specifically, we investigate how connected commits and features are by analyzing the

amount of features a commit usually affects through dataflow. Knowing what fraction of all commits contributing code to a project are part of dataflow-based interactions can show how often new commits affect the data of a feature. Regarding this, it is worth considering that commits constituting code of a feature are very likely to influence said feature through dataflow, as discussed in section 3.3. Since structural interactions coinciding with dataflow interactions are so obvious, programmers are also much more aware of them. Depending on the prevalence of feature regions in a project's code space, this could heavily skew the data in one direction, as a large portion of all dataflow interactions would stem from these obvious interactions. Therefore we want to especially focus on commits that aren't part of a feature, because programmers might not be aware or intend that changes introduced with these commits also affect features through dataflow. With the gathered information, another interesting aspect of dataflow-based CFIs to examine, is the relationship between the size of a feature and the number of outside commits interacting with it through dataflow. Determining to what extent feature size is the driving factor in this, could tell us whether it is worth considering other possible properties of features responsible for the number of commits affecting its data.

**RQ3: How do authors interact with features?**

Usually there are many programmers working on the same software project, implementing different features, sometimes alone, sometimes with the help of colleagues. We want to shine some light on the exact statistics of this by combining structural commit-feature interactions with high-level repository information. One major question we want to answer is how many authors implement a feature on average, where considering the size of a feature could help put this data into perspective. The collected results could serve as advice for software companies on how to allocate programmers on to-be implemented features.

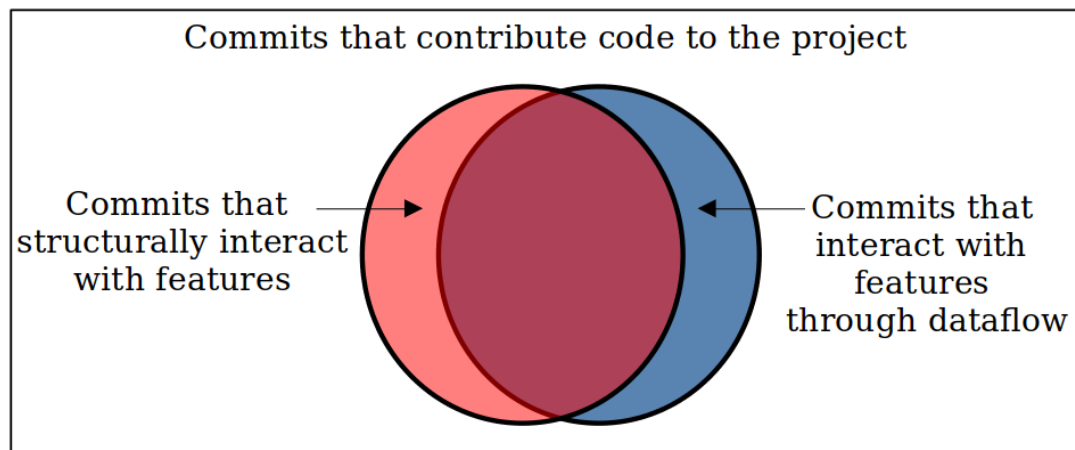


Figure 4.1: Kinds of commits in software projects investigated in this work. In the first two **RQs** we have discussed different kinds of commits and the ways in which they interact with features. Figure 1 showcases them in a venn diagram and illustrates the dependencies and divisions between them.

## EVALUATION

---

This chapter evaluates the thesis core claims.

### 5.1 RESULTS

In this section, present the results of your thesis.

### 5.2 DISCUSSION

In this section, discuss your results.

### 5.3 THREATS TO VALIDITY

In this section, discuss the threats to internal and external validity.



## RELATED WORK

---

This chapter presents related work.

For example, Kapser and Godfrey [3] investigated ...

Apel et al. [1] analyzed ...

In earlier work [1, 4], they have shown ...





## CONCLUDING REMARKS

---

### 7.1 CONCLUSION

### 7.2 FUTURE WORK





## APPENDIX

---

This is the Appendix. Add further sections for your appendices here.



## BIBLIOGRAPHY

---

- [1] Sven Apel, Don Batory, Christian Kästner, and Gunter Saake. *Feature-Oriented Software Product Lines: Concepts and Implementation*. Springer, 2013.
- [2] Christopher Hundhausen, Adam Carter, Phillip Conrad, Ahsun Tariq, and Olusola Adesope. “Evaluating Commit, Issue and Product Quality in Team Software Development Projects.” In: SIGCSE ’21. Virtual Event, USA: Association for Computing Machinery, 2021, pp. 108–114. ISBN: 9781450380621. DOI: [10.1145/3408877.3432362](https://doi.org/10.1145/3408877.3432362). URL: <https://doi.org/10.1145/3408877.3432362>.
- [3] Cory J. Kapser and Michael W. Godfrey. “Supporting the Analysis of Clones in Software Systems: A Case Study.” In: *Journal of Software Maintenance and Evolution: Research and Practice (SME)* 18.2 (2006), pp. 61–82.
- [4] Christian Kästner, Sven Apel, and Martin Kuhlemann. “A Model of Refactoring Physically and Virtually Separated Features.” In: *Proc. Int. Conf. Generative Programming and Component Engineering (GPCE)*. ACM, 2009, pp. 157–166.
- [5] Max Lillack, Christian Kästner, and Eric Bodden. “Tracking load-time configuration options.” In: *Proceedings of the 29th ACM/IEEE international conference on Automated software engineering*. 2014, pp. 445–456.
- [6] Florian Sattler. “Understanding Variability in Space and Time.” To appear. dissertation. Saarland University, 2023.
- [7] Florian Sattler. *VaRA is an analysis framework that enables users to build static and dynamic analyses for analyzing high-level concepts using advanced compiler and analysis technology in the background*. <https://vara.readthedocs.io/en/vara-dev/> [Accessed: (24.05.2023)]. 2023.
- [8] Florian Sattler, Sebastian Böhm, Philipp Dominik Schubert, Norbert Siegmund, and Sven Apel. *SEAL: Integrating Program Analysis and Repository Mining*. 2023.