

Automated API Refactoring for Evolving Codebases

Daniel Ramos

Submitted in partial fulfillment of the requirements
for the dual degree of Doctor of Philosophy in

Software Engineering, and Computer Science

Committee

Co-Chairs:

| | | |
|---------------------|----------------------------|--------------|
| Dr. Claire Le Goues | Carnegie Mellon University | (Co-advisor) |
| Dr. Nuno Lopes | Instituto Superior Técnico | |

Members:

| | | |
|---------------------|-----------------------------|--------------|
| Dr. Ruben Martins | Carnegie Mellon University | (Co-advisor) |
| Dr. Joshua Sunshine | Carnegie Mellon University | |
| Dr. Vasco Manquinho | Instituto Superior Técnico | (Co-advisor) |
| Dr. Işıl Dillig | University of Texas, Austin | |

May 2024



This work is licensed under a Creative Commons Attribution 4.0 International License.

May 15th, 2024

Abstract

Modern software development heavily relies on third-party libraries and frameworks, which are known to yield significant productivity gains. Libraries expose functionality through Application Programming Interfaces (APIs). Although stable API selection is desirable, it is often not possible, as software must adapt to new technical requirements or shifts in stakeholder or market demands. Therefore, as libraries evolve, clients may need to migrate APIs to adapt to these changes.

The task of adapting APIs to accommodate non-functional changes is a form of *software refactoring*, a crucial practice in software engineering. Refactoring entails modifying code to improve its quality and reduce complexity. However, refactoring is typically labor-intensive and prone to errors. The complexity of API refactoring has spurred numerous research efforts towards automating this task. A widely-used method for automating API refactoring is to generate match-replace rules by mining vast amounts of data from client projects of the libraries, sourced from collaborative coding platforms like GitHub. However, a significant challenge with mining approaches is their limited effectiveness due to reliance on data from clients that have already undergone refactoring, which is often scarce.

In this thesis, we present novel methods for automated API refactoring that do not rely on extensive training data or specific refactoring examples from client projects. In particular, we explore three alternative data sources. First, we use API documentation to discover API mappings, which we use to both generate migration rules and as a heuristic to guide the migration process. Specifically, the API mappings are used as an heuristic to guide a program synthesis approach to migrate client code effectively and reliably. Second, we use the API development process, particularly library pull requests, to learn API migration rules for addressing breaking changes. Our core idea is that if a library changes functionality, its tests and internal usages likely change as well, providing a rich source of data for generating migration rules. Third, we exploit natural language, as software is enriched with an abundance of natural language data including commit messages, issue reports, and comments. We use this unstructured data to test equivalence between API usages by synthesizing pairs of code examples in the source and target libraries. Our goal is to then abstract these code examples to generate broadly applicable migration scripts.

So far, we have implemented our ideas in three automated refactoring tools. Our tools leverage state-of-the-art program synthesis and machine learning techniques for establishing API mappings, synthesizing migration scripts, and migrating client code directly. We evaluated the three tools on real datasets, by migrat-

ing client programs found on collaborative coding platforms. Our ongoing research aims to automatically generate training examples from natural language and documentation, which we will then use to generate migration scripts for libraries where migration pairs do not exist.

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 2 |
| 1.1 | Motivation | 2 |
| 1.2 | Goal of this work | 4 |
| 1.3 | Evaluation Methodology | 6 |
| 1.4 | Expected Contributions and Outline | 7 |
| 2 | Research Background and Related Work | 9 |
| 2.1 | What are Application Programming Interfaces? | 9 |
| 2.2 | What is Refactoring, and Why is it Important? | 10 |
| 2.3 | Automated Techniques for API Refactoring | 11 |
| 3 | Synthesis-Driven Refactoring with Documentation and Error Messages | 15 |
| 3.1 | Motivating Example | 16 |
| 3.2 | SOAR's Algorithm | 17 |
| 3.3 | Evaluation | 25 |
| 3.4 | Discussion and threats to validity | 29 |
| 3.5 | Key Takeaways and Contributions. | 31 |
| 4 | Mining API Refactoring Rules from the API Development Process | 33 |
| 4.1 | Motivating Example | 34 |
| 4.2 | MELT's Approach | 36 |
| 4.3 | Evaluation | 40 |
| 4.4 | Discussion | 46 |
| 4.5 | Key Takeaways and Contributions | 48 |
| 5 | A Language for Automating Logically Related Changes | 49 |
| 5.1 | Motivation | 50 |
| 5.2 | Preliminaries | 52 |
| 5.3 | Language Syntax and Overview | 54 |
| 5.4 | Language Runtime | 58 |
| 5.5 | Evaluation | 61 |

| | |
|---|-----------|
| 5.6 Discussion | 69 |
| 5.7 Key Takeaways and Contributions | 71 |
| 6 Automating Logically Related Changes | 73 |
| 6.1 Proposed Research Overview | 74 |
| 6.2 Generating Synthetic Pairwise Examples | 74 |
| 6.3 Generating Rules Graphs in the POLYGLOTPIRANHA language | 75 |
| 6.4 Evaluation Strategy | 76 |
| 7 Conclusions | 77 |
| Bibliography | 77 |

1

Introduction

Contents

| | |
|--|---|
| 1.1 Motivation | 2 |
| 1.2 Goal of this work | 4 |
| 1.3 Evaluation Methodology | 6 |
| 1.4 Expected Contributions and Outline | 7 |

1.1 Motivation

Modern software development relies heavily on third-party libraries and frameworks. These libraries facilitate software reuse [1], allow developers to leverage quality implementations for a desired functionality, and yield significant productivity benefits [2]. Libraries expose their functionality through *Application Programming Interfaces (APIs)*. APIs serve as contracts between the library developers and its clients [3], providing functionality through a set of functions and methods, and hide concrete implementation details [4]. [5] Although stable API selection is desirable, the dynamic nature of software often renders it impractical. This dynamism in software [6] is driven by changing technical requirements and shifts in stakeholder or market needs [7]. As requirements and libraries evolve, clients may need to migrate APIs to accommodate these

shifts [7]. Furthermore, APIs themselves might break, become deprecated [8], or become exploitable, posing severe security risks, thus forcing downstream clients to update their usage accordingly.

The task of changing APIs to accommodate non-functional changes is a specific instance of *software refactoring*, a crucial software engineering practice. Refactoring involves modifying code with the goal of enhancing its quality and reducing its overall complexity [9]. Broadly, refactoring facilitates new feature development [10], assists managing technical debt [11], and prevents codebase decay [12]. Neglecting to refactor can escalate future costs; for example, the Consortium for Information & Software Quality (CISQ) estimates the cost to address the accumulated technical debt in the U.S. at approximately \$1.31 trillion [13].

However, refactoring is generally labor-intensive and error-prone [14]. Even seemingly straightforward tasks, like moving between two libraries that provide similar functionality, can be challenging and tricky. For instance, *PyTorch* [15] and *Tensorflow* [16], two of the most popular deep learning libraries, have different conventions regarding the order of values in tensor dimensions, which can lead to subtle bugs during a migration. Migrating APIs requires significant domain-specific expertise in both the source and target libraries [17]. Furthermore, achieving the desired functionality often extends beyond simple method mappings. Developers may need to write additional boilerplate code, figure out the correct argument combination in the replacement API, and manage cascading changes like type migrations and removing/adding import statements. Migration tasks are also difficult as APIs continuously evolve, often rendering prior knowledge obsolete.

The complexity of API refactoring has inspired numerous research efforts towards the automation of this task. At a high level, the goal is to automatically infer and generalize API refactorings from minimal user/developer input in order to reapply them on a large scale. These refactorings include migrations, updates from breaking changes, or handling deprecations. Like in any automation task, the goal is to orchestrate refactoring in such a way that user intervention is minimal. However, full automation is hard, and existing tools typically provide partial support rather than a fully integrated refactoring experience (e.g., [18–22]).

The most common approach to automated API refactoring is to use heuristics or learning algorithms to establish mappings between APIs and create match-replace rules [23]. The data for these algorithms is typically sourced from large-scale collaborative coding platforms like GitHub [24]. The data is obtained either by mining commits from library client projects that have undergone migrations [19, 21, 22], and can also be supplemented by information from new client projects in the most up-to-date APIs [18].

One significant challenge with these mining approaches is that the effectiveness of the tools is limited by their reliance on mining data from client projects that have already refactored their APIs. Unfortunately, this data is scarce. For example, a recent study found that 81.5% of projects maintain outdated dependencies [25]. Additionally, the mining process can only occur after clients begin transitioning between versions, precluding use shortly after a new version of the library is released [26, 27]. On the other hand, unsupervised learning methods [28] circumvent the issue of pair-wise data scarcity. However, they require extensive training data. This would in theory primarily benefit well-established libraries and APIs but, in practice unsupervised

methods are not as effective. Moreover, they cannot tackle lesser-known libraries and proprietary code.

1.2 Goal of this work

The goal of this work is to develop novel methods for automated API refactoring that do not rely on extensive training data or pairwise migration examples from downstream client projects of libraries. In particular, we focus on techniques for two API refactoring tasks: 1. **cross-library migrations**, and 2. **same-library migration to cope with breaking changes**. Our techniques are based on two different approaches:

1. inferring match-replace rules for updating APIs;
2. refactoring the APIs directly using a tool.

We observe that a wealth of high-quality information for automated API refactoring can be sourced from API documentation, the development processes of the APIs, and other self-contained information in libraries, along with natural language associated with the API development process.

Thesis Statement

Automated API refactoring can be effectively achieved without relying on extensive pair-wise training data. By leveraging self-contained data in libraries and leveraging synthesis techniques, we can generate rules and migrate code to facilitate API refactoring.

Next, we explain the rationale and hypotheses that underlie our proposed approaches.

1.2.1 API Documentation

APIs intended for widespread reuse are often reasonably well-documented [29]. The quality, quantity, and structure of this documentation can vary widely [30]. However, as code meant to be called and reused by unrelated client applications, documentation is often key to successful API uptake [30]. High-quality API documentation usually includes detailed descriptions of each function, method, and class, along with their expected inputs, outputs, and error conditions [31]. This information can be leveraged for classifying APIs, finding alternatives, or adapting usage when breaking changes are introduced.

In general, documentation may also contain examples and best practices [31, 32], which can be analyzed to infer latent properties of the APIs. These properties can then inform automated refactoring tools, ensuring that the transformation results in a use case aligned with the API's purpose. Furthermore, modern API documentation is increasingly enriched with metadata, such as annotations in dynamically typed programming languages [33]. This metadata can also be leveraged to guide the automated refactoring process.

In the case of breaking API changes, documentation evolves alongside the API itself [34], reflecting changes and deprecations in API behavior. If the documentation of evolving APIs provides transition examples from

old usage to new usages, or even natural language descriptions of how to adapt to the changes, then these can also be leveraged for automated refactoring.

In short, our hypothesis is that we can: *1. leverage API documentation to establish mappings between closely related APIs and libraries to facilitate API refactoring, 2. use examples and structured information from the documentation to guide the refactoring process, and 3. leverage metadata to increase the accuracy of the refactoring.*

1.2.2 API Development Process

Pull requests have become the *de facto* standard for software development on collaborative platforms like GitHub [35, 36]. In this method, a developer first clones the project (i.e., makes a personal copy of the project), makes changes in their copy, and finally submits these changes to the central repository for review. Pull requests typically include a title, a natural language description of the proposed changes, how they relate to project milestones or issues, and a set of commits (i.e., code file changes). These are reviewed by a core group of maintainers who decide whether to accept, request revisions, or reject the changes.

Pull requests involving API changes are also a rich source of data for mining transformation rules for API refactoring. Our hypothesis is that by examining pull requests (PRs) submitted to libraries where APIs are broken, we can identify and extract rules governing these changes. First, we can use tags in PRs to identify API changes by searching for labels such as "deprecated", "breaking change", and "API change". If the PR corresponds to such an API change, we can use self-contained commits in the PRs to learn code transformation rules for updating client code. The internal updates to the library source code resulting from the API change (such as test case changes) serve as the ground truth for mining rules and adapting client code.

1.2.3 Natural Language

Software development extends beyond coding and is enriched by a plethora of natural language data, such as commit messages, issue reports, discussions, and code comments [37]. This unstructured data can also be crucial in informing and guiding the development of automated API migration tools.

Developers frequently use natural language to describe intended changes or updates to APIs, whether in documentation or issue trackers [38]. Leveraging general-purpose large language models [39] is a promising avenue for transforming this unstructured data into structured data for our techniques. We seek to explore natural language-to-code techniques to create API refactoring rules directly from language (rule synthesis). We also aim to use natural language comments and descriptions in pull requests to generate code examples, which can be used to mine and generalized into refactoring rules.

During the automated refactoring process, error messages from the compiler can also be used in pinpointing issues, and suggesting how to fix them. These messages often contain details about issues in API calls or usage. By interpreting these error messages, we can better align with the correct API usage.

1.2.4 Program Synthesis

So far, our observations have highlighted alternative data sources that can be leveraged for automated API migration. In this section, we discuss how we will leverage them towards automated API refactoring. Specifically, in our proposed work, we frame automated API refactoring as a program synthesis problem.

Program synthesis is a research area focused on automatically generating programs that comply with a given specification, such as natural language or input-output examples [40]. In our work, we frame API refactoring as a synthesis problem in two distinct manners:

1. **Direct Migration using Synthesis:** We approach the API refactoring problem by synthesizing programs directly. Specifically, given a program that uses library A , our goal is to generate an equivalent program that replaces all usages of A with an alternative library B . Unlike the general setup in synthesis (such as programming-by-example [40]), our specification is complete. That is, our source program fully describes the intended behavior of the program to be synthesized.
2. **Script-based Migration Synthesis:** At times, it is desirable to synthesize a script for migration, rather than directly migrate the code. In these cases, our goal is to generate a script that can automatically migrate programs from A to B . In our work, we represent our scripts using declarative languages for matching and replacing code.

1.3 Evaluation Methodology

We test our hypotheses by developing automated API refactoring methods based upon them. The evaluation of each method is as follows.

1.3.1 Metrics

To measure the effectiveness and efficiency of our techniques, we use the following metrics:

1. **Refactoring Accuracy:** This metric evaluates the quality of the refactorings produced by our tools. Specifically, for fixing breaking changes, we upgrade client projects to a new library version (e.g., v_1 to v_2) and run their tests post-upgrade both before and after refactoring. In the context of library migrations, we migrate the code and run the same test suite where feasible. If this is not possible, we use automatically generated tests.
2. **Runtime Performance:** We assess the runtime taken by our tools during refactoring. We also measure the impact of each component of our proposed approaches in ablation studies. Finally, we compare our tools to state-of-the-art alternatives where such comparison is possible and fair.
3. **Refactoring Rule Accuracy:** Validating refactoring rules is hard, as ground truth for the rule typically does not exist, and there are infinitely many programs that could spur different behaviors from the rule.

Thus, instead, we manually validate the quality of the match-replace rules generated by our proposed approach. Our manual analyses always take into consideration the context and the intent of the rule. For rules to fix breaking changes, we examine documentation, developer comments, and online discussions to determine if our rule for addressing the breaking change preserves the original behavior. We use inter-rater agreement [41] as outlined by best practices for all our manual validations.

1.3.2 Benchmark Dataset

The benchmarks for each technique are chosen based on the task at hand. For evaluating refactoring accuracy and efficiency, our benchmarks consist of programs written in popular languages like python or java, each accompanied by a set of test cases. The benchmarks feature outdated API usages that require refactoring. In some cases, they are supplemented with automatically generated test cases. Tests play a crucial role in ensuring that the program's behavior remains consistent before and after the refactoring.

For assessing rule accuracy, we select relevant migrations and breaking changes from open-source libraries, particularly those where the author possesses the most expertise. We infer and validate behavior-preserving refactoring rules. The rule accuracy validation is done by manually checking if match-replace rules are generally applicable, considering existing documentation, developer discussions, and other publicly available information.

1.4 Expected Contributions and Outline

We propose several techniques and prototype tools based on our ideas:

- **Synthesis-based Refactoring Using Documentation and Error Messages (Chapter 3):** A novel technique, named SOAR, that uses readily available API documentation to learn API representations and migrate code between libraries (Chapter 3). SOAR uses program synthesis to automatically compute the correct configuration of arguments and necessary glue code for API invocation. It also integrates the interpreter's error messages to refine the search space during refactoring. This work was published in the proceedings of the *International Conference on Software Engineering 2021* [42].
- **Mining API Migration Rules from Pull Requests (Chapter 4):** A novel technique, named MELT, that generates lightweight API refactoring rules for fixing breaking changes by using data from the pull requests that broke the API. The data is used in two ways: first, natural language descriptions and code changes in pull requests are used to generate adaptation examples. The examples are then tested and generalized into API transformation rules. Secondly, code changes to test cases within the library that were adapted to cope with the breaking change are also used to mine rules. This work was published in the proceedings of the *International Conference on Automated Software Engineering 2023* [43].

- **A Multi-Language Code Transformation Tool (Chapter 5):** We develop a new declarative domain-specific language (DSL) for expressing interdependent multi-language code transformations . The language aims to make lightweight transformation tools more expressive for complex refactorings. The language takes inspiration from other lightweight match-replace languages and enhances their design. We demonstrate the language and toolset effectiveness and expressiveness in an industrial setting, as well as its ease of use. This work was just accepted at the *International Conference on Programming Language Design and Implementation 2024* [5].
- **Generating Refactoring Programs from Natural Language (Planned – Chapter 6):** Our proposed work aims to generate refactoring scripts to migrate complex API usages. Our goal is extend the MELT approach to be able to tackle the types of migrations discussed in Chapter 5.

This work is a compilation of the papers discussed above.

2

Research Background and Related Work

Contents

| | |
|---|----|
| 2.1 What are Application Programming Interfaces? | 9 |
| 2.2 What is Refactoring, and Why is it Important? | 10 |
| 2.3 Automated Techniques for API Refactoring | 11 |

In this chapter, we delve into two fundamental concepts to this work and to modern software development in general: Application Programming Interfaces (APIs) and refactoring. We start by offering a concise overview of APIs, discussing their role in creating re-usable software systems (Section 2.1). Following this, we overview the practice of refactoring in software engineering, a key process for enhancing code quality (Section 2.2). The chapter concludes with an overview of how researchers have developed automated approaches for API refactoring, highlighting trends and gaps in this area (Section 2.3).

2.1 What are Application Programming Interfaces?

Developers oftentimes build software meant for use by other software, rather than directly by end-users. An important consideration in such cases is to decide which functionality is intended to be provided and the abstractions to be exposed to developer clients. According to design best practices [44], there must be a

clear separation between the software's functional requirements (what it does) and its concrete implementation (how it does it). This separation is critical to ensure that clients are not burdened with irrelevant implementation details. The functionalities made available to clients are collectively known as an Application Programming Interface, or API.

The term API has been loosely used to describe a wide array of concepts, interpreted differently across technical domains [45]. Indeed, there is considerable debate about its meaning and usage [46]. However, in this work, we use the term 'API' to denote a software API, typically exposed and implemented as a library with a collection of classes, functions, or methods (depending on the programming language) that expose certain functionalities for other programmers to use. A single API can have multiple implementations (or none, if abstract) in the form of bindings that share the same programming interface. For example, because Scala and Java compile to compatible bytecode, Scala programs may use APIs implemented in Java.

Since APIs are intended for reuse, they often come with documentation, describing classes, methods, and sometimes typical usage cases, design rationales, and performance discussions [47, 48]. Regardless of what constitutes a particular API, the important underlying concept is that an API is a well-defined interface providing specific services to other software components.

APIs are the cornerstone of modern software engineering. Modern applications are often built on top of many APIs, which are also built upon other APIs. This approach yields significant productivity gains, allowing developers to focus on their specific tasks. Various research studies have investigated API design [49], evolution [50], and usability [51].

2.2 What is Refactoring, and Why is it Important?

Real-world software is constantly evolving, often driven by the need for enhancements and changes to meet new requirements [52]. As software adapts to additional functionality, code complexity tends to increase, and the software drifts from its original design. This issue is further exacerbated when developers prioritize short-term fixes over comprehensive, long-term solutions. While these immediate solutions might initially appear cost-effective, they frequently result in significant, hidden long-term maintenance costs. Such escalating complexity eventually leads to a gradual deterioration, a phenomenon known as *technical debt* [11]. The longer technical debt remains unaddressed, the more 'interest' it accrues, making it increasingly challenging to develop new features and maintain existing code.

Providing accurate estimates for software maintenance spending is challenging, but it is widely accepted that the expenses related to software maintenance — including bug fixes, design enhancements, and code restructuring — vastly exceed the cost of actual new feature development. Complex code often results in extended development periods, subtle bugs, and increased cost of change [53]. Additionally, according to Minelli et al. [54], developers allocate approximately 70% of their time to understanding code rather than

coding. Therefore, a well-thought-out design and adherence to best practices are crucial to reduce this effort [55].

One way to tackle this spiral of complexity is through a software engineering practice known as *refactoring*. Refactoring is defined as *the process of changing a software system in such a way that it does not alter the external behavior of the code, yet improves its internal structure* [9]. Sometimes, refactoring is colloquially used in a more generic and less rigorous manner than this definition suggests, and indeed, some authors argue that refactoring is not always behavior-preserving [56]. However, in this document, we are primarily concerned with behavior-preserving refactorings, aiming to improve maintainability without external changes, as refactorings that do not preserve behavior have limited potential for full automation.

In general, a refactoring is parameterized as a sequence of program transformations applied either manually or automatically with a tool, along with a specification that guarantees behavior preservation if satisfied. In modern software development, tests are the most widely adopted form of specification. Formally, we run a given program on a set of predefined inputs and then check if the observed output matches the expected output. Although refactoring can be applied in different paradigms, we focus on imperative programming languages. Many forms of refactoring also exist, like renaming variables, moving methods, migrating across languages [57], or refactoring APIs [58] (the primary focus of this work).

There are two primary API refactoring tasks of interest: library migration and library updates. A *library migration* involves replacing a third-party library (the source library) in a project with an alternative one (the target library). In the case of a *library update*, the task is to replace an old version with its newer iteration. Numerous factors can drive such migrations or updates, including (but not limited to) the introduction of new features, performance enhancements, improved community support and documentation, compatibility issues, license changes or restrictions, and deprecation, among others.

2.3 Automated Techniques for API Refactoring

The challenges associated with refactoring have motivated multiple research efforts in automation. In this section, we provide a brief overview of key ideas and trends in automated API refactoring.

2.3.1 Code Transformation Languages

As discussed in Section 2.2, a refactoring is parameterized as a sequence of program transformations. For large-scale, automated refactoring, it is necessary to represent this sequence of transformations in a language. Here, we give a brief summary of code transformation languages proposed over the years.

Code transformation languages and tool sets are often declarative. In general, these tools work with find-replace rules; each rule has two parts: (1) a template for selecting the source code to be transformed, (2) and the replacement template shows how the matched code should be transformed. Declarative match-

replace tool sets can either be language-specific or language-agnostic. Language-specific tool sets include `CocciNelLe` [59], widely adopted in the Linux community for C, and `Refaster` [60], used at Google for Java. Variants of these tools include `CocciNelLe4J` [61] for Java and `GoPatch` [62] for Go. Since these tools are language-specific, they often leverage existing compiler infrastructure and use semantic information, such as control flow, to enable precise transformations. While these tools provide similar languages, they leverage different aspects of each language. However, significant efforts are required to introduce and maintain new language front-ends for these tools.

In contrast, lightweight tools do not depend on compiler infrastructure and are thus language-agnostic, akin to regular expressions. Some examples include `Comby` [23] and `ast-grep` [63]. However, lightweight tools generally have a limited understanding of code context and semantics, making it challenging to express complex code transformations.

2.3.2 Learning API Transformations

2.3.2.A Mining Migration Scripts

For automating API refactoring, it is necessary to source refactoring rules automatically from a reliable origin. Prior research has primarily focused on mining API clients to learn these rules. The key idea behind these approaches is to locate API clients on collaborative online coding platforms such as GitHub and analyze their version history (i.e., commits). The goal is then to identify commits where clients have either migrated to a newer API version or updated their API usage. This data is aggregated and used to mine refactoring rules using various algorithms.

The majority of approaches to mining rules have primarily targeted APIs in object-oriented languages (specifically, Java and C#). Examples include `A3` [19] and `Mediator` [21]. While these tools use varied techniques for the mining process and adopt different internal representations for refactoring, their approaches are based on similar ideas. Some tools represent refactorings as sequences of edits on structure rather than as match-replace rules. For example, `APIFix` [18] mines transition examples from both previously-migrated and new client repositories to learn refactoring actions. `APIFix` represents code transformation as a sequence of tree edits using `Refazer's` [64] program synthesis engine, rather than using regular match-replace rules. Abstract Syntax Tree (AST) edits are typically more challenging to understand [65].

`APIMigrator` [22] and `AppEvoLve` [20] also mine client repositories for commit data to generate refactoring rules and apply them to clients. A key difference in these tools compared to prior work is that they also leverage differential testing techniques [66] to validate edits on clients, rather than only checking the syntactical validity of transformed code.

A major issue with existing mining approaches is that data for mining is often very scarce, which heavily limits the applicability and usefulness of these approaches. A recent study by Kula *et al.* found out that 81.5% of projects keep outdated dependencies. Additionally, the mining process can only occur after clients begin

transitioning between versions, precluding use shortly after a new version of the library is released [26, 27].

2.3.2.B Mining API Mappings

A different line of work focuses on learning API mappings, rather than refactoring rules, to offer partial automation and assist developers with manual refactoring. For example, `Semdiff` [68] recommends API changes to developers by presenting a ranking of potential replacements. Other approaches mine statistical API recommendations include ARES [69], PART [70], Sysedminer [71], MAM [57], Staminer [72], and Java2CS [73]. This research primarily targets cross-language migrations, with most tools mining pairs of bilingual projects to learn API mappings between two languages. However, bilingual projects are increasingly rare. DeepAM [28] takes a different approach, using deep learning to find a common representation. It learns semantic representations using embeddings, based on surrounding API context and comments, to understand APIs in multiple languages.

2.3.2.C Example-based approaches

Instead of taking existing examples (e.g., from GitHub) to learn transformation rules, researchers have proposed *refactoring-by-example* techniques. These approaches infer the transformations as a program in a low-level DSL from input-output examples. The idea is for the developer to provide an example which is then generalized, rather than mining from client projects. For example, LASE [74] and Bluepencil [75] infer an edit script from two example edits. More recently, Overwatch [76] integrates refactoring by example ideas into core IDE infrastructure to learn edit sequences, not just from input-output examples but also from intermediate steps (i.e., using temporal context). Catchup! [77] follows a different approach. Instead of asking users for examples, it asks library developers to record refactoring operations. The idea is to provide a patch alongside the new library version so that clients can automatically update their dependencies.

3

Synthesis-Driven Refactoring with Documentation and Error Messages

Contents

| | |
|--|----|
| 3.1 Motivating Example | 16 |
| 3.2 SOAR's Algorithm | 17 |
| 3.3 Evaluation | 25 |
| 3.4 Discussion and threats to validity | 29 |
| 3.5 Key Takeaways and Contributions. | 31 |

In this chapter, I discuss my completed and published work on a “Synthesis Approach for Open Source API Refactoring” (SOAR) [42]. In this work, we leverage: 1. API documentation (including metadata), and 2. interpreter error messages to automatically refactor APIs. Unlike previous approaches, SOAR does not require any pairwise training data to find API mappings. Moreover, SOAR migrates code using a program synthesis based approach rather than using refactoring rules. This makes SOAR more powerful and expressive, as some transformations cannot be expressed as match-replace rules. SOAR is one the first end-to-end, synthesis-based approaches to API refactoring requiring minimal to no training data. In the evaluation, we show that it can

| | |
|---|---|
| <pre> 1 import tensorflow.keras.layers as tf 2 ... 41: 42 self.conv1 = tf.Conv2D(43 filters=32, 44 kernel_size=3, 45 strides=(2, 2)) 46 ... 47 48 self.dense1 = tf.Dense(10) 49 self.dense2 = tf.Dense(1568) 50 ... 51 52 self.deconv1 = tf.Conv2DTranspose(53 filters=64, 54 kernel_size=3, 55 strides=2) 56 ... 57 58 self.relu6 = tf.ReLU() </pre> | <pre> 1 import torch 2 ... 39: 40 self.var0 = torch.nn.Conv2d(41 in_channels=32, 42 out_channels=32, 43 kernel_size=(3, 3), 44 stride=(2, 2), 45 padding=(0, 0)) 46 ... 47 48 self.var5 = torch.nn.Linear(64, 10) 49 self.var6 = torch.nn.Linear(10, 1568) 50 ... 51 52 self.var8 = lambda t: t.permute(0, 3, 1, 2) 53 self.var9 = torch.nn.ConvTranspose2d(54 in_channels=32, 55 out_channels=64, 56 kernel_size=(3, 3), 57 stride=(2, 2), 58 padding=(0, 0)) 59 ... 60 61 self.var13 = torch.nn.ReLU() </pre> |
|---|---|

Figure 3.1: An example of how SOAR refactors a program written with TensorFlow (left) to using PyTorch (right). Note that the whole program consists of 15 APIs calls to TensorFlow, though we only show four blocks of them (i.e., A, B, C and D) for brevity. SOAR can migrate the full program in 161 seconds.

successfully migrate functions / programs of up to 45 lines within reasonable time frames.

3.1 Motivating Example

We illustrate some of the difficulties of manual API refactoring via example. Consider the code snippet depicted on the left-hand side of Figure 3.1. This code snippets features an autoencoder, a specific type of neural network, developed using the TensorFlow API. Our objective is to transition this code to the PyTorch API as shown in the right-hand side. For context, an autoencoder is an encoder-decoder style neural network designed for data compression. It is trained to transform data into a more compact form (i.e., a latent representation), and then reconstruct the original data as accurately as possible.

The example in Figure 3.1 shows only a portion of the program, for didactic purposes. To build the first layer of the encoder, it calls the Conv2D function, creating a convolution layer for 2D images. After further (elided) activation and convolution layers, it calls Dense to output a latent representation of the input image. Decoding this output follows roughly the same procedure as the encoding, but using Conv2DTranspose instead of Conv2D. The function ReLU appears in both the encoder (not shown) and decoder, is used to ensure non-linearity of the neural network.

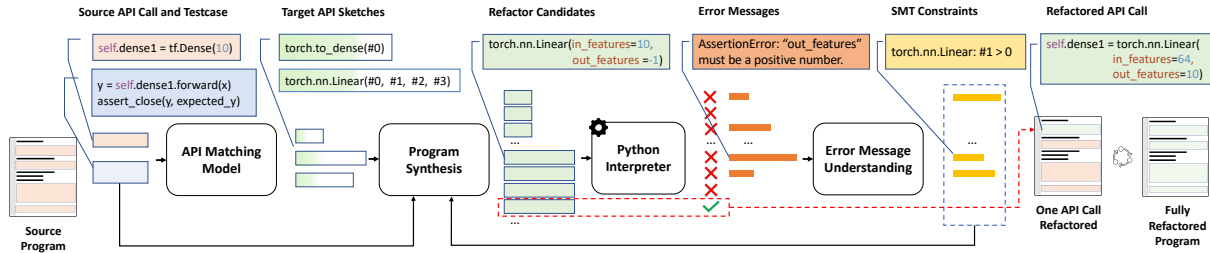


Figure 3.2: Overview of SOAR’s architecture.

The example in Figure 3.1 illustrates several of the core challenges in refactoring open-source APIs, as well as opportunities to inform an automated approach. First, the names of function calls implementing similar functionality may be very similar or even identical (such as those in blocks A, C, and D), or completely different (e.g., `Dense` versus `Linear` in block B). If a developer were performing this migration manually, they might reference the API documentation. For example, TensorFlow documentation describes the `Conv2D` class as a 2D convolution layer (e.g., spatial convolution over images)” [16]; the corresponding PyTorch documentation for the `Conv2d` call describes it similarly, as a 2D convolution over an input signal composed of several input planes” [78]. Here, the function names map well, but when this does not happen, the documentation should at least provide analogous descriptions for functions offering equivalent functionality.

Identifying the appropriate function is only part of the challenge in API refactoring. Even when the correct function is known, APIs mapped for the same functionality may have parameters with different names, types, conventions, and default values. This is evident in the majority of calls in our example (as seen in blocks A, B, and C). For instance, the `Conv2D` functions in both libraries take different parameters. There is some overlap — both include `kernel_size`, and `stride` corresponds with `strides` — yet they may expect different types (for example, `kernel_size` takes an integer in TensorFlow but a tuple in PyTorch). In some cases, new arguments must be inferred, varying based on context. For example, the first parameter of the `Linear` API calls in block B, mapped from `Dense`, require an extra argument to be dynamically computed. This makes it infeasible to write a simple match-replace rule for `Dense` to `Linear` mapping, as arguments must be dynamically generated. Finally, there are instances where no single function in the target API matches the semantics of a call from the source API, necessitating a one-to-many mapping, as illustrated in the conversion of the `Conv2DTranspose` call in block C.

3.2 SOAR’s Algorithm

This section describes SOAR, our approach for automatic API migration. We begin with an overview of the method (Section 3.2.1) before providing more detail on individual components (Section 3.2.2; 3.2.3; 3.2.5).

Algorithm 1 SOAR($\mathcal{I}, \mathcal{S}, \mathcal{T}, C$)

Input: \mathcal{I} : existing program, \mathcal{S} : source library, \mathcal{T} : target library, C : test cases

Output: O : refactored program

```
1:  $\vec{r}$  : API mapping = MAPAPI( $\mathcal{T}, \mathcal{S}$ )
2:  $O = \{\}$ 
3: for each  $l \in \mathcal{I}$  do
4:    $O = O \cup \text{REFACTORLINE}(l, \mathcal{T}, C, \vec{r})$ 
5: end for
```

Parameters

- **in_channels** (*int*) – Number of channels in the input image
- **out_channels** (*int*) – Number of channels produced by the convolution
- **kernel_size** (*int or tuple*) – Size of the convolving kernel
- **stride** (*int or tuple, optional*) – Stride of the convolution. Default: 1
- **padding** (*int or tuple, optional*) – Zero-padding added to both sides of the input. Default: 0
- **padding_mode** (*string, optional*) – 'zeros', 'reflect', 'replicate' or 'circular'. Default: 'zeros'
- **dilation** (*int or tuple, optional*) – Spacing between kernel elements. Default: 1
- **groups** (*int, optional*) – Number of blocked connections from input channels to output channels. Default: 1
- **bias** (*bool, optional*) – If `True`, adds a learnable bias to the output. Default: `True`

Figure 3.3: Description of the program parameters in `torch.nn.Conv2d` documentation [78].

3.2.1 Overview

Figure 3.2 shows an overview of the SOAR architecture, while Algorithm 1 provides an algorithmic view. SOAR takes as input a program \mathcal{I} consisting of a sequence of API calls from a source library \mathcal{S} , the source (\mathcal{S}) and target (\mathcal{T}) libraries and their corresponding documentation, and a set of existing test cases (C). Since the user wants to refactor code from \mathcal{S} to \mathcal{T} , we assume that the user already has test cases for \mathcal{I} that can be reused to check if the refactored code (O) has the same functional behavior as the original code (\mathcal{I}). Refactoring proceeds one line at a time in \mathcal{I} , finding/constructing an equivalent snippet of code (composed by one or more lines) that uses APIs of the target library \mathcal{T} ; the composition of all these translated lines comprises the output O .

For each API call in the input program, the first problem either a developer or a tool must face is to identify methods in the target API that implement the same functionality (i.e., for a given set of input parameters, the target API call must generate the same output). SOAR uses an *API matching model* to identify target API calls. This model is built using NLP techniques that analyze the provided API documentation for each call, and provides a mapping (\vec{r} in Algorithm 1) that computes the similarity between each target API function and

each potential source API function. SOAR uses this to find the most likely replacement methods in the target API for each source API call in the input program. We provide additional detail in Section 3.2.2.

Given a potential match call in the target API, the next step is to determine how to call it, in terms of providing the correct parameters, in the correct order, of the correct type. SOAR uses program synthesis to automatically write the refactored API call, using the provided test cases to define the expected behavior of the synthesized code and its constituent parts. The synthesis process can be assisted with additional automated analysis of API documentation, which often provides key information about each parameter, namely (1) whether it is required or optional, (2) its type, (3) its default value (if applicable), and (4) constraints between arguments, input and output (e.g., input and output tensor shapes). Figure 3.3 shows a snippet of the descriptions of all parameters for `torch.nn.Conv2d`. For example, the parameter `stride` is optional; it takes type `int` or `tuple`, and its default value is 1. Analysis of this documentation can produce a specification constraint for the `stride` parameter, assisting the program synthesis task. Section 3.2.3 describes the synthesis step.

Given a potential rewrite in the target API, a natural step for a developer would be to run the refactored code on test inputs. Unsuccessful runs can be quite informative, because many APIs (especially in the deep learning and data science domains) provide error messages that can be very helpful for debugging. SOAR simulates the manual debugging process by first adapting the input whole-program test cases to test partially refactored code, and then extracting both syntactic and semantic information from any error messages observed when running them. SOAR uses this information to add new constraints to the iterative synthesis process (Section 3.2.5).

After migrating all calls in the source API to the target API such that all input tests pass, SOAR outputs a fully refactored program. Subsequent sections provide additional detail on the previously described steps.

3.2.2 API Matching using Documentation

The first step in migrating a call in a source API is to identify candidate replacement calls in the target API with similar semantics. SOAR's API matching model achieves by analyzing the prose documentation associated with the APIs, rather than mining client projects.

At a high level, the model embeds each API method call in a source and target library into the same continuous high-dimensional space, and then computes similarity between two calls in terms of the distance between them in that space. We explored two information retrieval approach to obtain latent API representation: TF-IDF (term frequency – inverse document frequency) [79] and pretrained word embeddings [80].

TF-IDF. TF-IDF finds the most representative words in a sentence (which are usually different from the most frequent ones). The core idea is to discard irrelevant words when computing the API representation. For example, words like “the” or “this” convey very little information about an API.

For our TF-IDF model, we first derive a bag-of-words representation \mathbf{x}^i from a description of an API call after some stemming of the words with the Snowball Stemmer [81]. $\mathbf{x}^i = [x_1^i, x_2^i, \dots, x_n^i]$ where x_j^i denotes

the frequency with which word x_j appeared in the sentence \mathbf{x}^i , and n is the size of the vocabulary from the descriptions of all APIs considered. A TF-IDF representation of the call is computed as Equation (3.1):

$$\text{TF-IDF}(\mathbf{x}^i) = \left[\frac{x_1^i}{\sum_{t=0}^m x_1^t}, \frac{x_2^i y}{\sum_{t=0}^m x_2^t}, \dots, \frac{x_n^i}{\sum_{t=0}^m x_n^t} \right] \quad (3.1)$$

However, the major downside of TF-IDF is that it does not encode the similarities between words themselves. For example, consider two hypothetical call descriptions: (1) *Remove the last item of the collection*, and (2) *Delete one element from the end of the list*. They are semantically similar but since they have minimal overlapping words, a TF-IDF representation method would not recognize these two API calls as similar.

Tfidf-GloVe. We can fix this problem by adding the use of pretrained word embeddings. Specifically, we use the GloVe embedding [80], which is trained on a very large natural language corpus and learns to embed similar words closer in the embedding space. Since the paper was published many other embeddings models have emerged but they are fundamentally built using the same ideas.

To obtain sentence embeddings from individual words, we perform a weighted average of the word embeddings and use the TF-IDF scores of individual words as weight factors. It is a simple yet effective method to obtain sentence embedding for downstream tasks, as noted by previous work [82, 83]. This is shown in detail as Equation 3.2, where \mathbf{w}_j is the vector encoding the GloVe embedding of word x_j :

$$\text{Embedding}(\mathbf{x}^i) = \sum_{j=1}^n \frac{x_j^i \cdot \mathbf{w}_j}{\sum_{t=0}^m x_j^t} \quad (3.2)$$

By including the GloVe embedding, word similarity is preserved; by including the TF-IDF terms, the influence of embeddings of common words is greatly reduced. However, GloVe is trained with Common Crawl [84] which contains raw webpages, which is a mismatch from our domain of textual data (i.e., data science and programming). This causes a lot of OOV (out-of-vocabulary) problems.

API matching. Given the representation of two APIs $\text{Rep}(\mathbf{x}^i)$, $\text{Rep}(\mathbf{x}^j)$ in the same space $\text{Rep}(\cdot)$, we compute their similarity with cosine distance:

$$\text{sim}(\text{Rep}(\mathbf{x}^i), \text{Rep}(\mathbf{x}^j)) = \frac{\text{Rep}(\mathbf{x}^i) \cdot \text{Rep}(\mathbf{x}^j)}{|\text{Rep}(\mathbf{x}^i)| |\text{Rep}(\mathbf{x}^j)|} \quad (3.3)$$

For computational efficiency, we pre-compute the similarity matrix between the APIs across the source and target library. So we will be able to query the most similar API for the synthesizer to synthesize its parameters on the fly.

3.2.3 Program Synthesis

Instead of crafting match-replace rules using the API mappings, we refactor client code directly using program synthesis. This allows our refactoring engine to support expressive refactorings. Formally, given input test cases and an API matching model providing a ranked list \vec{r} of APIs in the target library, the synthesis model

Algorithm 2 REFACTORLINE($l, \mathcal{T}, C, \vec{r}$)

Input: l : line of code from \mathcal{I} , \mathcal{T} : target library, C : test cases, \vec{r} : ranked list of API matchings

Output: \mathcal{R} : refactored snippet

```
1:  $\mathcal{O} = \{\}$ 
2: for each  $l' \in \vec{r}$  do
3:    $\vec{s} = \text{GENERATE\_SKETCHES}(l', \mathcal{T})$ 
4:   for each  $s \in \vec{s}$  do
5:      $\mathcal{R} = \text{FILL\_SKETCH}(s)$ 
6:     if  $\text{PASSTESTS}(\mathcal{R}, C)$  then
7:       return  $\mathcal{R}$ 
8:     end if
9:   end for
10: end for
```

automatically constructs new, equivalent code, of one or more lines, that uses APIs of the target library \mathcal{T} . The refactored program \mathcal{O} has the same functionality as input program \mathcal{I} , and passes the same set of tests C .

To refactor each line of the existing program \mathcal{I} , we use techniques of programming-by-example (PBE) synthesis [40]. PBE is a common approach for program synthesis, where the synthesizer takes as specification a set of input-output examples and automatically finds a program that satisfies those examples. In the context of program refactoring, our examples correspond to the test cases for the existing code. For our experiments we restrict ourselves to straight-line code where each line returns an object that can be tested. With these assumptions, we can automatically generate new test cases for each line k of program \mathcal{I} . This can be done by using the input of the existing tests, running them, and using the output of line k as a new test case for the program composed by lines 1 to k .

Our program synthesizer for refactoring of APIs is presented in Algorithm 2 and it is based on two main ideas: (i) program sketching, and (ii) program enumeration. For each line l in program \mathcal{I} , we start by enumerating a program sketch (i.e., program with holes) using APIs from the target library \mathcal{T} (line 3). For each program sketch, we perform program enumeration on the possible completion of the API parameters (line 5). For each complete program, we run the test cases for the program up to line l . If all test cases succeed, then we found a correct mapping for line l between libraries \mathcal{S} and \mathcal{T} (line 6). Otherwise, we continue until we find a complete program that passes all test cases.

Program Sketching. Program sketching is a well-known technique for program synthesis [85] where the programmer provides a sketch of a program and the program synthesizer automatically fills the holes in this sketch such that it satisfies a given specification. We refactor one line of program \mathcal{I} at each time. Our first step is to use the ranked list of APIs to create a program sketch where the parameters are unknown. For

instance, consider the first layer from the motivating example that shows the network for an autoencoder using TensorFlow:

```
tf.keras.layers.Conv2D(filters=32, kernel_size=3, strides=(2, 2))
```

A possible sketch for this call using PyTorch is:

```
torch.nn.Conv2d(#1, #2, (#3,#4), stride=(#5, #6), padding=(#7, #8))
```

Where holes #*i* have to be filled with a specific value for the APIs to be equivalent. This approach works for one-to-one mappings but would not support common one-to-many mappings where the parameters often need to be transformed before being used in the new API. This is the case of the previous API where a reshaping operation must be performed before calling the PyTorch API. To support this common behavior, we include in our program sketch one API from the target library \mathcal{T} and common reshaping APIs (e.g., permute).

The sketch that corresponds to the refactoring solution of the Conv2D API from TensorFlow uses a reshaping API before calling the Conv2d API from PyTorch:

```
lambda x: x.permute(#9, #10, #11, #12)
torch.nn.Conv2d(#1, #2, (#3, #4), stride=(#5, #6), padding=(#7, #8))
```

Using Occam's razor principle, our program synthesizer enumerates program sketches of size 1 and iteratively increases the size of the synthesized program up to a specified limit.

Program Enumeration. For each program sketch \mathcal{P} , our program synthesizer enumerates all possible completions for each hole. Since each hole has a given type, we only want to enumerate well-typed programs. We encode the enumeration of well-typed programs into a Satisfiability Modulo Theories (SMT) problem using a combination of Boolean logic and Linear Integer Arithmetic (LIA). This encoding is similar to other approaches that use SMT-based enumeration for program synthesis [86] and encodes the following properties:

- Each hole contains exactly one parameter;
- Each hole only contains parameters of the correct type.

A satisfying assignment to the SMT formula can be translated into a complete program. The types for each hole can be determined by extracting this information from documentation, by performing static analysis, or by having this information manually annotated in the APIs. The available parameters and their respective types can be extracted automatically from the parameters used in the *k*-th line of program \mathcal{I} and by any default parameters that can be used in the API from \mathcal{T} that appears in the program sketch \mathcal{P} . For instance, for the Conv2d example presented in this section, we consider as possible values for the holes, the values that appear in the existing code (32, 3, 2) and default values for integer parameters (-1, 0, 1, 2, 3) that are automatically extracted from documentation.

Encoding the enumeration of well-typed programs in SMT has the advantage of making it easier to add additional logical constraints that can prune the search space.

Shape:

- Input: $(N, C_{in}, H_{in}, W_{in})$
- Output: $(N, C_{out}, H_{out}, W_{out})$ where

$$H_{out} = \left\lfloor \frac{H_{in} + 2 \times padding[0] - dilation[0] \times (kernel_size[0] - 1) - 1}{stride[0]} + 1 \right\rfloor$$
$$W_{out} = \left\lfloor \frac{W_{in} + 2 \times padding[1] - dilation[1] \times (kernel_size[1] - 1) - 1}{stride[1]} + 1 \right\rfloor$$

Figure 3.4: Relationship between the parameters of Conv2d API described in PyTorch documentation [78].

3.2.4 Documentating Metadata to guide Synthesis

As we described in Section 3.2.1, API documentation often provides additional useful information about parameters to function calls, including type and default values. For each considered API call, we scrape/process the associated documentation to extract these properties and encode them as SMT constraints to further limit the synthesizer search space.

Additionally, some APIs have complex relationships between parameters which if encoded into SMT may reduce the search space considerably. For instance, Figure 3.4 shows the relationship between the different parameters for the Conv2d API described in PyTorch documentation. For APIs with these kinds of shape constraints, we can encode these relationships into SMT to further prune the number of feasible completions. When we use these relationships in our experiments, we encode them manually (a one-time cost for an actual SOAR user or API maintainer), but we observe that in many cases they could be automatically extracted from documentation.

3.2.5 Error Message Understanding

Besides meta-data, we can also guide the refactoring engine using the error messages provided by the Python interpreter. The idea is to give feedback to the program synthesis engine during the refactoring process.

We use a simple natural language processing approach to extract data from compiler error messages. Specifically, we extract hyponymy relations and use Word2vec [87] to understand run-time error messages. We transform error messages into constraints for our synthesizer. Figure 3.5 illustrates the process.

Step 1: Extract hyponymy relation candidates from error messages. We perform an automatic extraction of customized hyponyms on each error message. Hyponyms are specific lexical relations that are expressed in well-known ways [88]. In encoding a set of lexico-syntactic patterns that are easily recognizable (i.e., hyponyms), we avoid the necessity for semantic extraction of a wide-range of error message text. We then use

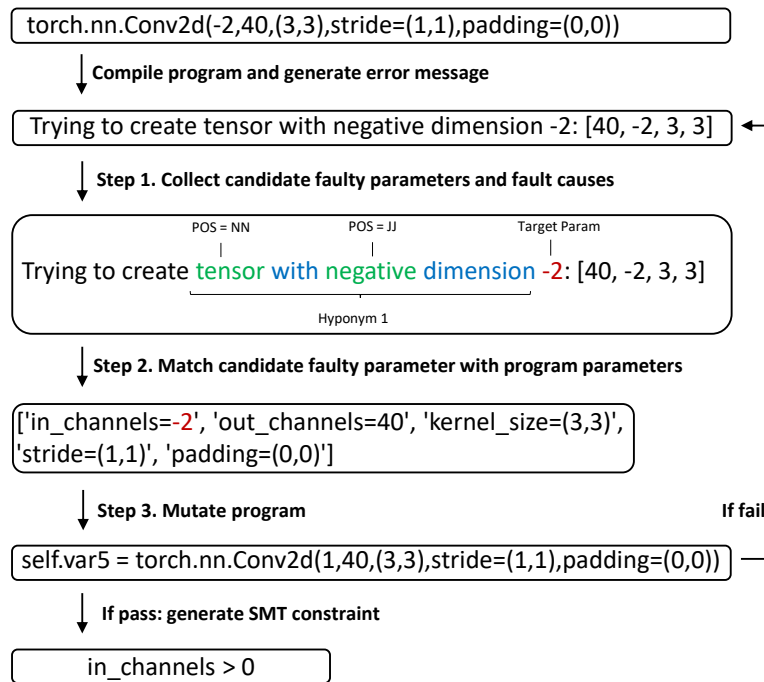


Figure 3.5: Example error message to SMT constraint pipeline using hyponym 1.

the collected hyponyms to map the error message to a single faulty parameter, and output a SMT constraint based on the faulty parameter.

We use four manually crafted lexico-syntactic patterns to identify hyponyms using *noun-phrases* (NP) and regular expressions frequently appearing in machine learning API error messages.

Step 2: Identify candidate faulty parameters and constraints. Step 2 uses different keywords based on the result of step 1 to identify the faulty parameter. As shown in Figure 3.5, an error message with hyponym 1 is likely to have the POS=JJ word as a parameter constraint (i.e., word “negative”). Based on the fault cause candidate, we then store all negative numbers as candidate faulty parameters (e.g., [40, -2, 3, 3] has -2 as the only faulty parameter). We then vectorize the candidate faulty parameter name (i.e., -2) and find the program parameter name with the closest vectorized distance. As shown in Figure 3.5, the parameter “in_channels = -2” has the nearest vectorized distance to the candidate faulty parameter -2. Based on the fault cause, we generate a candidate constraint. The example error message in Figure 3.5 has only one candidate constraint: “in_channels >= 0”.

Step 3: Mutate program. To validate the candidate faulty parameters and constraints, we mutate each faulty parameter according to each faulty parameter and constraints pair. We then re-compile the program for each mutation. If the error message remains the same, we discard the faulty parameter and constraint pair as a candidate. If the program passes, or if the error message changes, we store the faulty parameter and constraint pair as an SMT constraint. As shown in Figure 3.5, the API call mutator mutates the second parameter

("in_channels = -2") to a non-negative number. The mutator first attempts "in_channels = 0" and it encounters a different error message. From the new error message, we mutate this parameter to "in_channels = 1" and observe no further errors. Therefore, we refine our previous constraint to be "in_channels > 0", and store it as the final SMT constraint for the program in Figure 3.5.

3.3 Evaluation

We evaluate our approach by answering the following research questions:

RQ1. How effective is SOAR at migrating neural network programs between different libraries?

RQ2. How effective is API documentation to establish mappings? .

RQ3. How effective is API meta-data in guiding the refactoring process?

RQ4. How useful are error messages in guiding the refactoring process?

RQ5. Is SOAR generalizable to domains besides deep learning library migration?

3.3.1 Benchmarks and experimental setup

We collected 20 benchmarks for each of the two migration tasks. In particular, for the TensorFlow to PyTorch task, we gathered 20 neural network programs from TensorFlow tutorials [89], existing models implemented with TensorFlow [90] or its model zoo [91]. This set of benchmarks includes: Autoencoders for image and textual data, classic feed-forward image classification networks (i.e., the VGG family, AlexNet, LeNet, etc), convolutional network for text, among others. The average number of layers in our benchmark set is 11.80 ± 11.52 , whereas the median is 8. Our largest benchmark is the VGG19 network which contains 44 layers.

For the domain of table transformations, we collected 20 benchmarks from Kaggle [92], a popular website for data science. The programs in the benchmark set have an average of 3.05 ± 1.07 lines of code, and a median of 3 lines. Although the programs considered for this task are relatively small compared to the deep learning benchmarks, they are still relevant for data wrangling tasks as shown by previous program synthesis approaches [93].

All results presented in this section were obtained using an Intel(R) Xeon(R) CPU E5-2630 v2 @ 2.60GHz, with 64GB of RAM, running Debian GNU/Linux 10, and a time limit of 3600 seconds. To evaluate the impact of each component in SOAR, we run four versions of the tool. SOAR with TF-IDF (SOAR w/ TF-IDF) and SOAR with tfidf-GloVe (SOAR w/ Tfidf-GloVe) to evaluate the impact of API representation learning methods. SOAR without specification constraints (SOAR w/o Specs.) and SOAR without error message understanding (SOAR w/o Err. Msg.) to evaluate the impact of these components on the performance of SOAR.

Table 3.1: Execution time for the deep learning library migration task in each of the 20 benchmarks.

| | SOAR | SOAR w/o Specs. | SOAR w/o Err. Msg. |
|--------------------------|---------------|-----------------|--------------------|
| conv_pool_softmax(4L) | 1.60 | 23.02 | 14.35 |
| img_classifier(8L) | 12.82 | 336.00 | 65.66 |
| three_linear(3L) | 3.18 | 2.34 | 21.07 |
| embed_conv1d_linear(5L) | 5.27 | 123.85 | 16.90 |
| word_autoencoder(3L) | 1.81 | 1.46 | 2.64 |
| gan_discriminator(8L) | 12.80 | timeout | 252.20 |
| two_conv(4L) | 16.69 | timeout | 15.09 |
| img_autoencoder(11L) | 160.97 | 391.09 | 487.54 |
| alexnet(20L) | 425.22 | timeout | 66.13 |
| gan_generator(9L) | 412.47 | timeout | timeout |
| lenet(13L) | 280.91 | timeout | timeout |
| tutorial(10L) | 6.04 | timeout | 58.29 |
| conv_for_text(11L) | 9.04 | timeout | 32.29 |
| vgg11(28L) | 40.83 | timeout | 132.67 |
| vgg16(38L) | 82.05 | timeout | 139.27 |
| vgg19(44L) | 83.99 | timeout | 189.90 |
| densenet_main1(5L) | timeout | timeout | timeout |
| densenet_main2(3L) | timeout | timeout | timeout |
| densenet_conv_block(6L) | timeout | timeout | timeout |
| densenet_trans_block(3L) | timeout | timeout | timeout |

3.3.2 Implementation

The SOAR implementation integrates several technologies. Scrapy [94], a Python web-scraping framework, is used to collect documentation for the four libraries in our experiments. To enumerate programs in the synthesis step, we use the Z3 SMT solver [95]. For each target program call parameter, we extract an answer for the four parameter questions in Section 3.2.1 and generate corresponding SMT constraints. In both API matching model and the error message understanding model, the GloVe word embeddings [80] are used as an off-the-shelf representation of words. For the four libraries appearing in our two evaluation migration tasks, we use TensorFlow 2.0.0, PyTorch 1.4.0, dpLyr 1.0.1 (with R 4.0.0) and pandas 1.0.1, though our proposed method and associated implementation do not rely on specific versions.

3.3.3 Results

3.3.3.A RQ1: Overall SOAR effectiveness

Table 3.1 shows how long it takes to migrate each of the deep learning models from TensorFlow to PyTorch, using the various approaches. Our best approach (shown as SOAR) successfully migrates 16 of the 20 DL models with a mean run-time of 97.23 ± 141.58 seconds, and a median of 14.76 seconds. The average number of lines in the 16 benchmarks that we successfully migrate is 13.6 ± 12.14 , whereas the average number of lines in the output programs is 18.56 ± 16.40 . The reason the number of synthesized lines is higher than those in the original benchmarks is that we frequently do one-to-many mappings. In fact, 15 out of the 16 require at least one mapping that is one-to-many. In the 16 benchmarks, SOAR tests on average 4414.18 ± 5676 refactor candidates (i.e. program fragments tested for each mapping), and it needs to test a median 2111 candidates before migrating each benchmark. The reason 4 benchmarks timeout is that in each of these benchmarks there is at least one API in the benchmark that has a poor ranking (i.e., not in the top 200).

3.3.3.B RQ2: How effective is API documentation to establish mappings?

In Table 3.2, we show results of SOAR using different API representation learning methods, namely TF-IDF and TFIDF-GloVe, as described in Section 3.2. We can see that for these tasks of TensorFlow to PyTorch migration, using TF-IDF-based API matching model works better than adding pretrained GloVe embeddings. We believe this is because similar APIs are often named with same words (e.g., Conv2DTranspose vs. ConvTranspose2d) or even identical name (e.g., the APIs of creating a Rectified Linear Unit are both named as `ReLU(. . .)`), for TensorFlow and PyTorch. Thus simple word matching method like TF-IDF is suffice for API matching purposes.

Another interesting result worth noticing is that although the synthesis time differs for the two approaches, the average rankings are quite similar for most of the benchmarks. The reason is that despite the average rankings of correct target APIs being similar, the incorrect APIs ranked by the model before the correct one is different, and the time it takes to rule out those incorrect APIs varies greatly, determined largely by the number of parameters required for that API.

3.3.3.C RQ3: How effective is API meta-data in guiding the refactoring process?

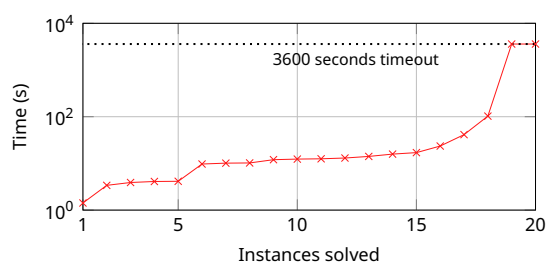
In Table 3.1, we also show the impact of specification constraints that describe the relationship between different parameters of a given API (see Section 3.2.3 for details). Even though, we only have these complex specifications for the 7 most common APIs, the impact on performance is significant. Without these specification we can only solve 6 out of 20 benchmarks. Relating the arguments of the APIs helps SOAR to significantly reduce the number of argument combinations that it needs to enumerate.

Table 3.2: Execution time and average API ranking for each of the 20 benchmarks using TF-IDF and GloVe models.

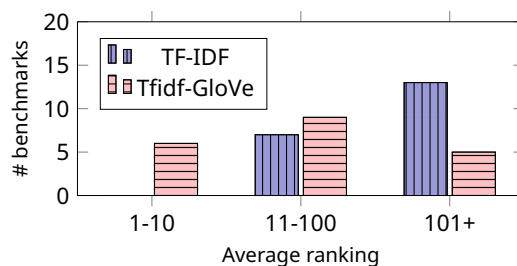
| | SOAR w/ TF-IDF | | SOAR w/ Tfidf-GloVe | |
|--------------------------|----------------|--------------|---------------------|--------------|
| | Time(s) | Avg. Ranking | Time(s) | Avg. Ranking |
| conv_pool_softmax(4L) | 1.60 | 1.0 | 1.56 | 1.0 |
| img_classifier(8L) | 12.82 | 2.8 | 31.04 | 2.8 |
| three_linear(3L) | 3.18 | 8.0 | 7.70 | 8.0 |
| embed_conv1d_linear(5L) | 5.27 | 2.4 | 7.75 | 2.4 |
| word_autoencoder(3L) | 1.81 | 1.0 | 1.52 | 1.0 |
| gan_discriminator(8L) | 12.80 | 2.8 | 37.01 | 2.8 |
| two_conv(4L) | 16.69 | 1.0 | 13.75 | 1.0 |
| img_autoencoder(11L) | 160.97 | 1.9 | 166.34 | 2.0 |
| alexnet(20L) | 425.22 | 2.1 | 428.42 | 2.1 |
| gan_generator(9L) | 412.47 | 2.0 | 1892.86 | 2.0 |
| lenet(13L) | 280.91 | 4.3 | timeout | 89.1 |
| tutorial(10L) | 6.04 | 2.3 | 21.31 | 2.4 |
| conv_for_text(11L) | 9.04 | 2.3 | 14.08 | 2.3 |
| vgg11(28L) | 40.83 | 1.8 | 73.92 | 1.8 |
| vgg16(38L) | 82.05 | 1.6 | 114.41 | 1.6 |
| vgg19(44L) | 83.99 | 1.5 | 114.98 | 1.5 |
| densenet_main1(5L) | timeout | 172.8 | timeout | 285.4 |
| densenet_main2(3L) | timeout | 16.0 | timeout | 387.5 |
| densenet_conv_block(6L) | timeout | 293.3 | timeout | 612.7 |
| densenet_trans_block(3L) | timeout | 291.0 | timeout | 480.0 |

3.3.3.D RQ4: How useful are error messages in guiding the refactoring process?

As shown in Table 3.1, SOAR performs significantly better when using the error message understanding model. We can observe that without this component, two of the benchmarks that SOAR could solve would timeout at the 1 hour mark. For the 14 benchmarks it still manages to solve, the synthesis time increases on average 4.66 \times . The number of performed evaluations also increase substantially for each benchmark. For the 16 benchmarks that SOAR successfully migrates, we evaluate an average of 43319.63 ± 61259.62 refactor candidates without the error message understanding model. This corresponds to a 9.81 \times increase in the number of necessary evaluations when compared to the full SOAR method. In summary, we can significantly reduce the search space by interpreting error messages.



(a) Execution time for each benchmark of the dplyr-to-pandas task with a timeout of 3600 seconds.



(b) Average ranking of the APIs for each of the 20 dplyr-to-pandas benchmarks.

Figure 3.6: Comparative results of dplyr-to-pandas task.

3.3.3.E RQ5: Generalizability of results

Our experiments so far concern deep learning library migration in Python. To study the generality of our proposed SOAR, we applied SOAR to another task of migrating from `dplyr`, a data manipulation package for R, to `pandas`, a Python library with similar functionality. Fig. 3.6b shows how the two API matching methods perform in this domain. While with `Tfidf-GloVe`, 30% of the correct APIs are ranked among the top 5, saving lots of evaluations for the synthesizer, none of the correct APIs are ranked by the TF-IDF-based matcher as its first 5 choices. Worse, nearly half of those are ranked above 100, making the synthesis time almost prohibitively long. We believe this is because the lexical overlap between the names of similar APIs in those two libraries is much smaller compared to the deep learning migration task. For example, `dplyr`'s `arrange` and `pandas`'s `sort_values` provide the same functionality (they both sort the rows by a given column), but the function names are different. In this way, `Tfidf-GloVe` can take advantage of the pretrained embeddings to explore the similarities between APIs beyond simple TF-IDF matching.

In Figure 3.6a, we show the time it takes to migrate each of the 20 benchmarks with a timeout of 3600 seconds when using word embeddings. We solve 18 out of 20 collected benchmarks in under 102.5 seconds. The average run time for 18 benchmarks is 17.31 ± 22.59 seconds and a median of 12.19 seconds. Note that for this task we did not consider error messages, nor specifications since we wanted to test how a basic version of SOAR would behave in a new domain. Moreover, for this domain, all the refactored benchmarks only used one-to-one mappings since no additional reshaping was needed before invoking `pandas` APIs. Even with these conditions, we show that we are able to successfully refactor code for a new domain across different languages.

3.4 Discussion and threats to validity

Overall, we focus our design and evaluation on deep learning and data science libraries. These libraries have properties that render them well-suited to our task in terms of common programming paradigms, and norms,

such as in the API documentation. However, we believe this is also a particularly useful domain to support, given the field's popularity and how quickly it moves, how often new libraries are released or updated, as well as the wide variety of skill sets and backgrounds present in the developers who write data science or deep learning code. Automation of migration and refactoring in this domain is very minimal, and we design SOAR as a step towards better tool support for this diverse and highly active developer population.

Next, we discuss the main limitations of our method and possible challenges for extending SOAR's ability to refactor new APIs, even potentially beyond the domain of data science.

Benchmarks. Our evaluation of SOAR uses benchmarks from well-known deep learning tutorials and architectures. However, they are all feed-forward networks, effectively sequences of API calls where the output of the current layer is the input of the next layer. There may be more applications that share this feature, but support for more complex structure is likely necessary to adapt to other domains.

Additionally, and naturally, the APIs in the benchmarks we collected may be biased and not reflect the set of APIs developers actually use. To assess this risk, we checked the degree to which the APIs used in our benchmarks appear to be widely used on other open-source repositories on GitHub. To do this, we collected the top 1015 starred repositories that have TensorFlow as a topic tag, which contains over 8 million lines of code and over 500K TensorFlow API calls. We found that 76% of the 1000+ repositories use API calls included in our benchmarks at least once, which validates some representativeness of our collected benchmarks.

Automatic testability. One benefit of the data science/scientific computing domain is that much of the input, output, and underlying methods are typically well-defined. As a result, it's particularly easy to test and verify the correctness of individually migrated calls, which can be processed in sequence. There may be other types of libraries that share these types of characteristics, like string manipulation or image processing libraries, whose intermediate outputs are strings/images. We also assume user-provided tests. Given the migration task, it is reasonable to assume the user has tests (the code must be sufficiently mature to justify migrating, after all), but a more general solution might benefit from automatically generating tests, which would both alleviate the input burden on the user and, potentially, reduce the risks of overfitting. In our current implementation, we moreover use the provided tests to construct smaller test cases for each mapping. This is particularly easy in this domain, because data science and deep learning API calls are often functional in their paradigm. Adapting the technique to other paradigms would require more complex test slicing or generation to support synthesis.

Correctness. Since we evaluate our migration tasks using test cases, it is always possible for our approach to overfit to said test cases. However, this threat can be mitigated if the user provides a sufficiently robust test suite that provides enough coverage.

Additionally, code written to different APIs may be functionally equivalent, but demonstrate different performance characteristics, which we do not evaluate. However, this fact is one reason users might find SOAR useful in the first place: a desire to migrate code from one library to another that is more performant for the

given use case.

Error message understanding. The error message understanding model is built on four domain specific lexico-syntactic patterns, which we identify as hyponyms when they appear in an error message. We propose the hyponyms based on the specific syntax of DL API error messages, thus take non-trivial human effort to make it generalize to error messages that appear when calling APIs from libraries of other domains. However, we believe the idea of program mutation (Step 3 of Fig. 3.5) is still widely applicable for the purpose of generating SMT constraints when dealing with error messages.

Synthesis. Our approach supports one-to-many mappings but it restricts the mapping to one API of the target library and one or more reshaping APIs. However, this could be extended to include many APIs of the target library at the cost of slower synthesis times. An additional challenge is to support many-to-one or many-to-many mappings since this would require extending our synthesis algorithm. However, even with the current limitations, our experimental results show that the current approach can solve a diverse number of benchmarks.

3.5 Key Takeaways and Contributions.

In this chapter, we demonstrate that API documentation can serve as a proxy to establish API mappings for migration. We used these mappings to guide a Synthesis approach for API Refactoring (SOAR). SOAR uses a generate-and-test strategy, as computing mappings alone is insufficient for completing a migration; API arguments and glue code also need to be mapped. Moreover, we cannot blindly trust the mappings derived from the documentation; indeed, the correct mappings are sometimes not correct (i.e., corresponding APIs may be apart in the embedding space, meaning they might be closer to other, non-semantically equivalent APIs). For these reasons, to complete migrations, we test the API mappings with multiple different combinations of arguments using synthesis. We determine that a particular API migration is complete if the synthesized code produces the same output on a set of auto-generated inputs. During the synthesis process, the interpreter also outputs warnings or errors due to API usage. We leverage this information with a simple error message understanding mode, which we use to prune the search space.

Overall, our approach successfully migrates API calls within reasonable time frames, particularly for small programs (under 100 lines) where it is feasible to compare objects across implementations. Since we synthesize code directly rather than generating migration scripts, this method is not generally applicable for large-scale refactorings of code bases comprising millions of lines of code. We aim to tackle this limitation in our proposed work (Chapter 6) on synthesizing of migration scripts for library migrations.

4

Mining API Refactoring Rules from the API Development Process

Contents

| | |
|---|----|
| 4.1 Motivating Example | 34 |
| 4.2 MELT's Approach | 36 |
| 4.3 Evaluation | 40 |
| 4.4 Discussion | 46 |
| 4.5 Key Takeaways and Contributions | 48 |

In this section, I discuss my completed and published work titled "Mining Effective Lightweight Transformations from Pull Requests" (MELT) [43]. In this work, we leverage the API development process from open-source libraries (i.e., pull requests) to mine rules for fixing breaking changes between library versions. One core idea of MELT is to identify pull requests (PRs) between library releases that break existing APIs. Upon identifying these PRs, we extract data from a variety of sources. First, we identify code changes to test cases that developers made after breaking the API. Our key observation is that if a library is well-tested, developers need to update test cases when they break an API; otherwise, the tests would fail. This data allows us to mine

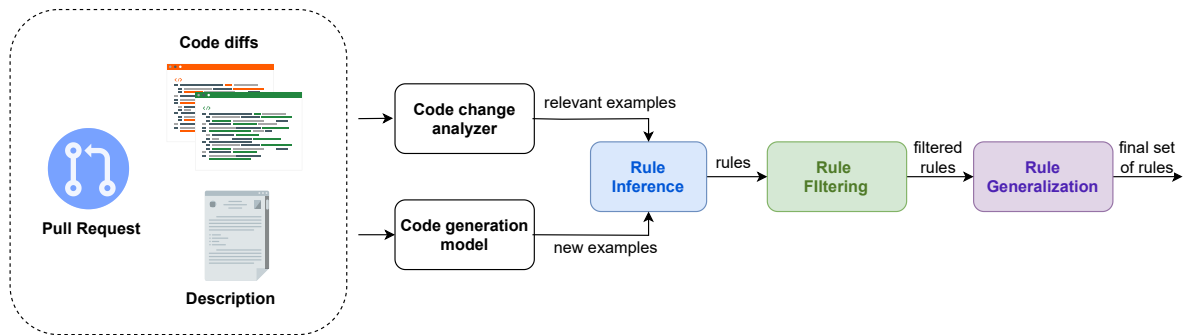


Figure 4.1: MELT takes as input a pull request (PR) and outputs a set of rules. The PR is processed in two ways: (1) the Code change analyzer identifies relevant code changes; (2) the Code generation model generates additional code examples. Rules are inferred from the code changes and examples, then filtered and generalized.

```

@@ -2427,7 +2427,7 @@ def test_datetimeindex(self):
2428 2428     # nat in index
2429 2429     s2 = Series(2, index=[Timestamp("20130111"), NaT])
2430 -     s = s2.append(s)
2430 +     s = pd.concat([s2, s])
2431 2431     result = s.to_string()

```

Figure 4.2: Code change in pull request #44539 [96] from the pandas-dev/pandas repository.

rules directly from the library instead of relying on downstream clients like the majority of state-of-the-art methods. The migrations themselves are expressed in the comby language for broader applicability. This means the synthesis only needs to happen once.

4.1 Motivating Example

Figure 4.1 provides a high-level overview of MELT and its main components. We delve into the specifics of each component in Section 4.2.

Pull requests are the input of MELT, as they are the key source that informs our approach. Pull requests generally contain all the code changes related to a given new feature. For example, Figure 4.2 shows an example code change from a pull request [96] submitted to pandas [97] that deprecates two popular APIs: `DataFrame.append` and `Series.append`.¹ MELT identifies code changes, such as the one shown in Figure 4.2, within the pull request using its *Code Change Analyzer* (Section 4.2.1) and inputs them into the *Rule Inference* algorithm (Section 4.2.3.A) to generate rules. The top portion of Table 4.1 shows two of the rules

¹Both APIs were later removed from pandas in version 2.0.0.

Table 4.1: Top: comby rules extracted from pandas pull request #44539, deprecating DataFrame.append and Series.append. Bottom: Rules extracted from sci-py pull request #14419, including original specific (“Spec”) and generalized (“Gen”) versions. Template variable constraints are omitted for brevity.

| | Match Template | Rewrite Template |
|------|---|---|
| | <pre> : [[s2]].append(: [[s1]]) where : [[s1]].type == Series and : [[s2]].type == Series </pre> | <pre> pd.concat(: [[s2]], : [[s1]]) </pre> |
| | <pre> : [[df]].append(: [[s]]) where : [[df]].type == DataFrame and : [[s]].type == Series </pre> | <pre> pd.concat(: [[df]], DataFrame(: [[s]].T.infer_objects()) </pre> |
| Type | Match Template | Rewrite Template |
| Spec | <pre> : [[s]].spline.cspline2d(: [[x]], : [y]) </pre> | <pre> : [[s]].cspline2d(: [[x]], : [y]) </pre> |
| Gen | <pre> : [[s]].spline.cspline2d(: [args]) </pre> | <pre> : [[s]].cspline2d(: [args]) </pre> |

MELT infers from the code changes for this specific pull request.

The rules in Table 4.1 are expressed in comby’s domain specific language [98]. The match template (left column) is the code structure for which comby searches. The rewrite template (right column) shows how to transform the matched code based on the variables in the match template [23]. comby uses template variables, i.e., placeholders that can be matched with certain language constructs. For example, a template variable to match alphanumeric characters is represented by `:[x]`, where `x` is the name of the template variable. The template variables in the match template can be constrained in multiple ways using a `where` clause. In particular, to prevent spurious matches, template variables can be constrained to be a certain type (like `: [s2]].type == DataFrame`). Although type information is not strictly required, it is useful when working with common API names such as `append` and `concat` (both are part of Python’s `stdlib`).

Code diffs in pull requests provide valuable information, however, they do not always contain the necessary code examples for rule inference. Fortunately, pull requests offer alternative sources of information that can be used to extract further details about the changing APIs. Figure 4.3 shows an informative comment left by a developer in a code file when deprecating namespace `scipy`’s [99] namespace `scipy.signal.spline` in favor of `scipy.signal`. To leverage all available information in the pull request, MELT uses a *Code Generation Model* to generate additional code examples and test cases for this change (Section 4.2.2). Figure 4.4 shows a simplified version of code GPT-4 [39] (a state-of-the-art model) generates from the pull request in Figure 4.3. The generated examples enable us to both infer and test the rules.

Since the test case executes successfully, MELT uses the code example to generate a rule by abstracting


```

26  scipy/signal/spline.py
@@ -0,0 +1,26 @@
1  + # This file is not meant for public use and will be
   + removed in the future
2  + # versions of SciPy. Use the `scipy.signal` namespace for
   + importing the
3  + # functions included below.
4  +
5  + import warnings

```

Figure 4.3: Pull Request #14419 [100] from `scipy/scipy`. This pull request was part of SciPy 1.8, released in Feb 2022.

| | |
|--|---|
| <pre> def old_usage1(image): return signal.spline.cspline2d(image, 8.0) def new_usage1(image): return signal.cspline2d(image, 8.0) </pre> <p style="text-align: center;">(a) Old and new usage of <code>cspline2d</code>.</p> | <pre> class TestEquiv(unittest.TestCase): def test_assert1(self): np.random.seed(181819142) image = np.random.rand(71, 73) assert np.allclose(old_usage1(image), new_usage1(image)) </pre> <p style="text-align: center;">(b) Test case for the transition.</p> |
|--|---|

Figure 4.4: Code generated by GPT-4 showcasing how to transition from the old `cspline2d` usage and a test case.

concrete identifiers and literals. For this case, MELT generates the rule in the third row of Table 4.1. This rule accurately reflects the deprecation made in the pull request (i.e., replaces the deprecated namespace with the new one). Nevertheless, a closer inspection reveals that the rule is too specific: it will only match usages where: (1) the first argument of `cspline2d` is an identifier (`[[s]]` only matches with identifiers), and (2) the function is called with two or more arguments. The `cspline2d` function can accept multiple combinations of arguments, including keyword arguments with default values.

To guard against overly-specific rules, MELT applies *Rule Generalization* (Section 4.2.3.C). For example, the template holes `[[x]]` and `[[y]]` in the rule in the first row of the bottom of Table 4.1 remain unchanged in the match and rewrite templates, indicating that they are not relevant to the change at hand. To enhance the rule’s applicability, MELT generalizes the specific argument combination, resulting in an updated version of the rule (shown in the last row of Table 4.1). The revised rule uses a more permissive match template using `[[args]]`, which can match any number of function arguments.

4.2 MELT’s Approach

In this section, we provide a brief overview of MELT’s approach.

4.2.1 Extracting Code Examples from Diffs in Pull Requests

MELT's input is a pull request \mathcal{P} , which contains both natural language descriptions and a set of code diffs, each of which corresponds to changed code snippets. However, not all diffs in a pull request are relevant to an API change, as they may encompass unrelated refactoring actions. Therefore, MELT first identifies which changes in the pull request are relevant to the API of interest.

MELT determines which code changes are relevant using its *Code Change Analyzer*. MELT starts by pinpointing which public APIs are affected by the pull request by examining the scope of each code diff to identify the affected function and its corresponding class. For example, for the code change in Figure 4.2, MELT identifies the function name `test_datatimeindex` and the class where the function comes from `TestSeriesFormatting`. MELT filters out test functions and private namespaces, to exclude API names that are not the main focus of the change.² On this example, the test class and method will be filtered, but other changes in the same PR (not shown) affect the `append` and `concat` methods, so MELT considers those methods relevant.

MELT then filters the code diffs to retain only those diffs and surrounding code that contain at least one of the relevant keywords. This produces a set of code examples to serve as inputs to rule inference. For the pandas example, although the test method itself is not a relevant API name, the code change in that test method does concern relevant API calls, and so these diffs will be retained for use in inference. A strength of this approach is its generalizability across multiple libraries and languages, since it works at token level.

4.2.2 Generating Examples for Mining using Natural Language

As illustrated in Section 4.1, pull requests sometimes lack sufficient code examples to infer migration rules. In a preliminary study, we analyzed 174 pull requests related to breaking changes and deprecations from pandas' release notes. We discovered that only 41 (23.6%) of these pull requests contained at least one meaningful code example showcasing the transition from old to new usage. However, pull requests offer other information sources about API changes, including natural language descriptions in comments, developer discussions, and documentation. Our key insight is that this additional data can also be leveraged to generate and test more code examples. MELT uses a *Code Generation Model* to produce extra code examples from this data. Generating code examples rather than the rules directly is advantageous, because we can test and validate the generated code, enhancing confidence in the rules inferred from it. Additionally, the code examples may enhance interpretability by demonstrating the provenance of inferred rules to MELT users.

We developed prompts and conducted experiments with GPT-4 8K [39], which is well-versed in our target libraries' code, to process PR information (code diffs, title, description, discussion). For each PR, we asked the model to generate: 1. transition examples, and 2. test cases asserting that the old behavior was the same as the new one. Full prompts and algorithmic description can be found in the paper. However, the key idea

²Although our experiments do not exercise this setting, developers can also provide the names of affected APIs when submitting the pull request, which MELT can use directly to eliminate irrelevant code changes.

| | |
|---|---|
| <p>(a) Code before migration</p> <pre>r = pd.read_csv(f, compression=comp, encoding=enc, index_col=0, - squeeze=True)</pre> | <p>(b) Code after migration</p> <pre>r = pd.read_csv(f, compression=comp, encoding=enc, index_col=0). + squeeze()</pre> |
|---|---|

Figure 4.5: Example code change from PR #43242 [102] in pandas

is to check if automatically generated tests are correct according to the automatically generated test suite. To make sure the test suite is not spurious, we sample multiple tests. If any test fails, the example is not considered for mining.

4.2.3 Rule Mining from Examples

MELT uses the comby language [23] and toolset [98] to express refactoring match-replace rules. We introduced some elements of the language in Section 4.1, with examples of comby’s syntax-driven match and rewrite templates. Formally, a rewrite rule in comby is of the form $M \rightarrow R$ where c_1, c_2, \dots, c_n , where M is the match template, R is the rewrite template, and c_1, c_2, \dots, c_n are constraints in the rule language. The key structure of comby rules are template variables, which are holes in the match and rewrite templates that can be filled with code. Template variable types include, e.g., `:[x]` matching alphanumeric characters (similar to `\w+` in regex), and `:[x]` matching anything between delimiters (e.g., `[]`, `()`, `{}`). comby also supports a small rule language to add additional constraints, like types or regular expression matches, on the template variables. comby’s website [98] provides the full syntax reference. Although language agnostic, comby is still language aware, and can deal with comments and other language-specific constructs. Its rules are also close to the underlying source, and thus typically easier to read than, e.g., transformations over ASTs.

4.2.3.A Rule Inference

Given a set of code examples, MELT infers a set of comby rules that can be used to automatically migrate APIs in client code. First, MELT parses the code files corresponding to each code diff into an abstract syntax tree (AST), identifying the nodes corresponding to the change before and after. MELT then uses a variation of InferRules’s algorithm [101] (adapted to Python) that always returns a single rule, and never abstracts away class names, method names, and keyword arguments.

To illustrate, consider the code change in Figure 4.5, where a library maintainer transforms a keyword argument into a function call. The smallest unit MELT considers for a comby rule is a source code line. Given the two assignment nodes corresponding to the change, rule inference then abstracts away child nodes with template variables. When a construct has the same character representation, MELT uses the same template variable. For the example, MELT abstracts the left-hand side and right-hand side of both assignments, yielding: `:[a]= :[b]`, and `:[a]= :[c]`. Notice that the template variable for the target of both assignments

is the same, `:[a]`, because their source representation is the same. However, MELT cannot match the right-hand side of the assignments (`:[b]`, and `:[c]`). It, therefore further decomposes the AST nodes' children:

```
:[a] = :[i].read_csv( [d],          :[a] = :[i].read_csv(:[d],
    compression=[e], encoding=[f], →    compression=[e], encoding=[f],
    index_col=[g], squeeze=[h])        index_col=[g]).squeeze()
```

MELT never abstracts away class names, function names, and keyword arguments, as preserving these details is crucial for API migration. Additionally, MELT consistently yields a single, all-encompassing rule. In this case, MELT can match every template variable in the match template with a corresponding node in the rewrite template except `:[h]`. Consequently, it attempts to further decompose the nodes, but still fails to match `:[h]`, ultimately reverting it and generating the final rule:

```
:[a] = :[i].read_csv([d],          :[a] = :[i].read_csv(:[d],
    compression=[e], encoding=[f], →    compression=[e], encoding=[f],
    index_col=[g], squeeze=True)      index_col=[g]).squeeze()
```

```
where :[h].type == int,
      :[i].type == pandas
```

After inferring a rule, MELT incorporates type guards. The goal is to constrain each template hole to its respective observed type. This step is crucial in preventing the misapplication of rules for common API names (e.g., matching `List.append` when the rule targets `DataFrame.append`). In contrast to previous rule synthesis approaches [101, 103], MELT directly incorporates type constraints into comby's rule language. This integration is possible because we extend comby to support Language Server Protocol (LSP) type inference. MELT uses the Jedi [104] type inference language server, making it available for client usage.

4.2.3.B Rule Filtering

Occasionally, MELT infers spurious rules (e.g., rules that contain variables in the rewrite template that might not be in scope). First, MELT discards duplicate rules within the same pull request (post generalization, as well). A rule is considered a duplicate if all of the match, rewrite template and template variable constraints are the same. MELT then further filters by:

API Keywords MELT discards transformation rules that do not contain the name of any affected APIs. This can occur when a developer modifies the surrounding context of a code block, for example, by wrapping a statement in a try-catch block (e.g., `:[x]→ try:\n:[x]`). These rules are considered spurious because they can match arbitrary code and are not specific to API migration.

Unsafe Variable and Private Namespaces MELT discards rules where a rewrite template uses either vari-

ables from private namespaces (indicated by calls with underscores, Python’s convention for private attributes/-functions/namespaces), or variables not present in the match template. This ensures that the rules do not rely on private or internal functionality that is not accessible to client code.

4.2.3.C Generalizing Rules

Rules inferred from single code examples may be too specific, as demonstrated in our rule for the `squeeze` example so far. This change is specific to a particular argument combination. However, the `read_csv` function has numerous optional arguments, and the rule should therefore be versatile. Moreover, it can only be applied to assignments, even though the migration applies to other contexts.

Therefore, our approach generalizes rules for broader applicability by abstracting irrelevant context and generalizing arguments. MELT removes common context in the source and match templates unrelated to the API. For our example, it unwraps the assignment statement and simply keeps the API call. MELT also uses InferRules [101] algorithm to find mappings between call nodes in the match and rewrite template, and generalizes common arguments. If there are multiple consecutive arguments between the match and rewrite template of the call node, we replace the arguments with a generic template variable `: [args]`. For our running example, the final rule is:

```
: [[i]].read_csv(: [args], squeeze=True) → : [[i]].read_csv(: [args]).squeeze()  
where : [[i]].type == pandas
```

Generalization is crucial to ensuring broader rule applicability. The paper describes the generalization algorithm in more detail.

4.3 Evaluation

We answer the following research questions:

RQ1. How effectively can MELT generate transformation rules from code examples in pull requests?

RQ2. How do code examples generated automatically complement code examples in pull requests?

RQ3. What is the impact of rule generalizability?

RQ4. Are the rules effective for updating client code?

4.3.1 Experimental Setup

4.3.1.A Implementation

Although our approach is largely language-agnostic, we implement it for Python libraries because: (1) Python is one of the most popular programming languages [105], and (2) there exists a gap in migration tools for

Python [27]. We implemented rule inference using the Python abstract syntax tree (AST) module. InferRules [101] was originally implemented for Java AST; we brought native implementation to Python. We also perform rule generalization at the Python AST level. For code generation, we used the state-of-the-art GPT-4 [39]. We extended comby to support Language Server Protocol (LSP)-based type inference over match templates [106] with Jedi [104], a state-of-the-art static analysis tool. MELT’s source code, data, and logs used for the evaluation are available at Zenodo [107].

4.3.1.B Methodology

We evaluated MELT using four of the most popular Python data science libraries: `numpy`, `scipy`, `sklearn`, and `pandas`. We collected a total of 722 pull requests for `pandas`, 141 for `sklearn`, 186 for `numpy`, and 130 for `scipy` using the GitHub QL API and web crawlers over release notes. We took a convenience sampling approach to find PRs concerning API or breaking changes, or deprecation-related PRs, moving backwards from the version of each library (as of April 2023); this includes merged PRs intended for future library releases, as well as those that have been released. We collected more PRs for `pandas` than other libraries because it had a higher number of pull requests, and breaking changes in `pandas` are particularly well documented. We then executed MELT on each pull request.

For our manual assessment of rule correctness and relevancy, two authors of this paper manually labeled a set of rules independently. We defined a rule to be correct if (1) it correctly reflects the change in the pull request, and (2) it is generally applicable to client code and does not overgeneralize (i.e., it will not produce incorrect migrations even if it matches the correct APIs in some cases). This procedure requires analyzing the pull request discussion, changes, source code, and documentation when necessary. The annotators discussed five representative examples together and then individually labeled 151 unique rules, achieving an inter-rater reliability (IRR) with a Cohen’s kappa of 0.84 (almost perfect agreement) [108]. Due to the high agreement, the first author labeled the remaining rules to cover all research questions.

4.3.2 Results

4.3.2.A RQ1: Mining Rules from Code Examples in PRs

Table 4.2 summarizes MELT’s rule inference algorithm on 1179 PRs (722 `pandas`, 130 `scipy`, 186 `numpy`, 141 `sklearn`). MELT’s ability to extract code examples from pull requests largely depends on the libraries’ testing practices. Nonetheless, a significant number of pull requests contain valuable examples for rule extraction. Previous studies [109] found that only 27.1% of migrations in a different set of libraries were potentially fully automatable. MELT generates correct migration rules for 12.2% of analyzed pull requests, indicating room for improvement (further explored in RQ2).

Running MELT’s rule inference algorithm to the 1179 PRs results in 5504 rules. After filtering and gen-

Table 4.2: RQ1. Left: Pull requests per library, with mined rules and correct rules. Right: Filtered and generalized rules mined per library, with total and correct counts.

| Library | # PRs | PRs with | | Mined Rules | |
|--------------|-------|-------------|---------------|-------------|-------------|
| | | Mined Rules | Correct Rules | Total | Correct (%) |
| pandas | 722 | 169 | 102 | 521 | 359 (68.9%) |
| scipy | 130 | 21 | 11 | 33 | 19 (57.6%) |
| numpy | 186 | 20 | 10 | 47 | 27 (57.4%) |
| sklearn | 141 | 38 | 21 | 82 | 56 (68.3%) |
| Total | 1179 | 248 | 144 | 683 | 461 (67.5%) |

eralization, we ended up with 683 rules. The right-most columns of Table 4.2 show the number of mined rules after generalization and filtering for each library, and their correctness based on manual validation. On 67.5% of the cases, our mined rules are correct and do not overgeneralize. However, on 32.5% of the cases, MELT derived incorrect, non-generally applicable, or over general rules. We observed three primary reasons for incorrect rules: (1) Code change not generally applicable, such that the rule cannot capture the context in which it is applicable. For example, in numpy PR #9475 [110], the `np.rollaxis` is deprecated in favor of `np.moveaxis`. Migrating from one API to another depends on the actual content of the variables used in the API, as it behaves differently depending on the variables' content. Our rule cannot capture this, as it only considers types, not content. (2) Overgeneralization of rule arguments. For instance, pandas PR #21954 [111] says "*read_table is deprecated. Instead, use pandas.read_csv passing sep='t' if needed.*". However, one of the inferred rules is `read_table(:[args])` \mapsto `read_csv(:[args])`, because the algorithm abstracts all arguments based on the code example. and, (3) Unrelated changes not caught by filtering.

4.3.2.B RQ2: Mining Rules from Autogenerated Code Examples

To evaluate the role example generation played in rule inference, we sampled 50 pull requests for each library (limited by budget). We used a template to create a prompt to ask the model to generate both code examples and test cases/inputs for the examples, per pull request. The prompt includes the title, description, discussion, and code changes. We used OpenAI's API to prompt GPT-4, with a (default) temperature of 0.2, and sampled the model 5 times to generate transition examples. We then followed up with the model to ask for test cases for each sample (in total 10 requests per PR).

The left side of Table 4.3 shows the number of unique examples generated for each library and the number of examples that passed the test suite. MELT produced 248 unfiltered and ungeneralized rules on these examples; filtering and generalization produced 156 unique rules. We also assessed whether these rules could

Table 4.3: Left: Code examples generated and passing tests per library. Middle: Pull requests with mined and correct (“Corr.”) rules from generated examples. Right: Filtered and generalized rules per library. Note: Limited to 50 PRs per library for budgetary reasons.

| Library | Code Examples | | PRs with Rules | | Mined Rules Correct | | |
|--------------|---------------|--------|----------------|-------|---------------------|------|-----|
| | Total | # pass | Total | Corr. | Total | Prev | New |
| pandas | 285 | 134 | 25 | 19 | 45 | 7 | 30 |
| scipy | 194 | 68 | 15 | 13 | 30 | 4 | 18 |
| numpy | 222 | 114 | 21 | 14 | 46 | 2 | 31 |
| sklearn | 187 | 63 | 21 | 13 | 35 | 5 | 17 |
| Total | 888 | 379 | 82 | 59 | 156 | 18 | 96 |

have been generated from the pull request code directly, by checking (1) whether they were mined in RQ1 (Section 4.3.2.A), or (2) whether they could be directly applied to their corresponding pull request (meaning that they could have been mined in RQ1, but may have been heuristically filtered away).

Table 4.3 summarize rule mining success using generated examples by pull request (middle columns); the right-hand side shows the number of rule mined. We categorized correct rules into those that could have been mined without new examples (prev), and those that are new with the generated examples. Like in the previous RQ, MELT can generate incorrect rules in some scenarios. Consider the following example rule: `::[[aah]].shift(:[aae], fill_value=:[aaf]) → ::[[aah]].shift(:[aae], fill_value=pd.Timestamp(:[aaf]))`.³ This rule is derived from pandas pull request number #49362 [112]. The release notes for the PR state: “Enforced disallowing passing an integer fill_value to DataFrame.shift and Series.shift with datetime64, timedelta64, or period dtypes”. This transformation is only valid if the series has a `datetime64` dtype object, a condition not captured by the rule. While the transformation correctly preserves behavior in this instance, it is incorrect for general application. More diverse tests for the code example could likely increase coverage and filter more incorrect rules.

4.3.2.C RQ3: Generalizability

Of the 156 rules we manually validated in RQ2, 41 had generalized arguments, and only 9 (22%) were incorrect. To further evaluate the impact of generalizability with an ablation study, by disabling the generalization procedure. We selected 15 rules that had been generalized, along with their non-generalized counterparts. Using Sourcegraph’s code search [113],⁴ we searched for repositories containing a given keyword in the rule (e.g., for `readcsv(..., squeeze=True)`, we searched for `squeeze=True`). We then cloned 50

³Template variables are omitted for brevity.

⁴Note SourceGraph only indexes repositories with at least two stars.

Table 4.4: Comparison of Non-General and Generalized Rules

| Library | Original Rule | | Generalized Rule | |
|---------|---|---------|---|---------|
| | Match Template | Matches | Match Template | Matches |
| pandas | <code>::[[x]].set_index(:[a], drop=[[b]], inplace=True)</code> | 2 | <code>::[[x]].set_index(:[args], inplace=True)</code> | 370 |
| | <code>::[[x]].read_csv(:[[a]], compression=[[b]], encoding=[[c]], index_col=[d], squeeze=True)</code> | 0 | <code>::[[x]].read_csv(:[args], squeeze=True)</code> | 21 |
| | <code>::[[aai]].apply(:[a], axis=[[b]], reduce=True)</code> | 3 | <code>::[[aai]].apply(:[args], reduce=True)</code> | 4 |
| scipy | <code>jaccard_similarity_score(:[[a]], :[[b]])</code> | 94 | <code>jaccard_similarity_score(:[args])</code> | 226 |
| | <code>::[[x]].filters.gaussian_filter(:[a], :[b], mode=[[c]])</code> | 0 | <code>::[[x]].filters.gaussian_filter(:[args])</code> | 86 |
| | <code>::[[x]].query(:[[a]], :[[b]], n_jobs=[c])</code> | 0 | <code>::[[x]].query(:[args], n_jobs=[y])</code> | 0 |
| | <code>::[[x]].hanning(:[[a]], :[[b]])</code> | 0 | <code>::[[x]].hanning(:[args])</code> | 0 |
| numpy | <code>::[[x]].alltrue(:[a], axis=[b])</code> | 7 | <code>::[[x]].alltrue(:[args])</code> | 208 |
| | <code>::[[x]].histogram(:[[a]], bins=[b], range=[c], normed=[y])</code> | 2 | <code>::[[x]].histogram(:[args], normed=[y])</code> | 66 |
| | <code>::[[x]].complex(:[[a]], :[[b]])</code> | 17 | <code>::[[x]].complex(:[args])</code> | 20 |
| sklearn | <code>BaggingClassifier(base_estimator=[[a]], n_estimators=[[b]], random_state=[[c]])</code> | 26 | <code>BaggingClassifier(base_estimator=[x], :[args])</code> | 220 |
| | <code>BaggingRegressor(base_estimator=[[a]], n_estimators=[[b]], random_state=[[c]])</code> | 7 | <code>BaggingRegressor(base_estimator=[x], :[args])</code> | 116 |
| | <code>KMeans(n_clusters=[a], init=[[b]], n_init=[[c]], algorithm='full')</code> | 0 | <code>KMeans(:[args], algorithm='full')</code> | 38 |
| | <code>AgglomerativeClustering(n_clusters=[a], linkage=[b], affinity=[c])</code> | 4 | <code>AgglomerativeClustering(:[args], affinity=[c])</code> | 28 |
| | <code>OneHotEncoder(sparse=[[aac]], categories=[[aan]], drop=[[aaz]])</code> | 0 | <code>OneHotEncoder(sparse=[x], :[args])</code> | 66 |

random repositories for each rule, and ran the generalized and non-generalized rules on these repositories, counting matches.

Table 4.4 shows matches for original and generalized rules, showing that generalization significantly improves rules applicability. For instance, the number of matches for the `set_index` case increased from 2 to 370 (185x) with generalization. Generalization is important because it abstracts context unrelated to API changes. As we focus on API migration in Python, where there can be many argument combinations (e.g., APIs with as many as 10 keyword arguments), generalization helps capture the essence of the change by abstracting arguments. Some rules had 0 matches because comby was unable to infer types (comby does not apply rules when it cannot infer types of a template match), or the query was poorly constructed.

4.3.2.D RQ4: Updating Client Code

To evaluate the effectiveness of our approach to updating developer code, we migrated outdated library API usage in developer projects found on GitHub for the `sklearn`, `pandas`, and `scipy` libraries. Collecting and running client projects requires significant manual effort: many projects do not specify dependencies or

Table 4.5: RQ4. Effects of rule application on developer projects.

| Library | Total Projects | Affected Projects | Unique Rules | Rule Applications | Additional Warnings | Resolved Warnings | Additional Passing Tests | Additional Failures | Resolved Failures |
|----------------|-----------------------|--------------------------|---------------------|--------------------------|----------------------------|--------------------------|---------------------------------|----------------------------|--------------------------|
| sklearn | 20 | 10 | 6 | 27 | 9 | 598 | 2 | 1 | 1 |
| pandas | 20 | 10 | 4 | 23 | 0 | 44 | 7 | 81 | 7 |
| scipy | 20 | 6 | 5 | 23 | 0 | 266 | 0 | 1 | 0 |
| Total | 60 | 26 | 15 | 73 | 9 | 908 | 9 | 83 | 8 |

provide tests. We therefore did not evaluate numpy API usage, but we can expect similar results.

We found client projects by searching GitHub for public repositories that used outdated versions of each library, and included code that matched to at least one of the match templates of an inferred rule from RQs 1 and 2. We applied a total of 15 unique rules across the three libraries. We provide detail on specific rules and projects in Zenodo [107]. For each library, we identified 20 client projects that used outdated versions, and between one and three rules applied. We cloned each project, updated its library dependencies to a version with the breaking change, installed necessary dependencies, and ran all tests to note passing tests, failures, errors, and warnings. We then used comby to automatically update the outdated API usage, and reran the tests to compare results post-migration. We did this separately for each applicable rule.

Table 4.5 summarizes results. Total Projects refers to the total number of projects to which we applied rules and tested. Affected Projects refers to the number of evaluated projects that had a change in the tests after rule application from new or resolved warnings, passed tests, or failures. Not all of the projects had tests affected by rule application, either because test coverage was incomplete or because persistent failing tests in developer projects obscured the effect of rule application.

For sklearn, slightly less than half the developer project tests were affected by rule application. Only two of the projects showed a negative impact of rule application, where one project had an additional failing test and another project had nine new warnings. The sklearn rules were applied without type information, which is one potential cause for the negative impact. The other affected projects had warnings resolved, ranging from 1 to 563 warnings resolved for a single project. One project had additional passing tests.

For pandas, rule application affected half of client projects. While there were 81 additional failures from pandas rules, they were isolated to four projects and a single rule. These new failures occurred because of a lack of type information, meaning one rule was erroneously applied to API calls unrelated to the pandas library. In other projects, the same rule was applied correctly, even without type information, and successfully resolved warnings. The other three unique pandas rules were applied with type information. No pandas rules introduced new warnings.

For scipy, rules were also applied absent type information, but only one application introduced an error. All six affected scipy projects had warnings resolved by rule application, and none of the scipy rule

applications caused additional warnings.

Of the 60 evaluation repositories, 34 had no change in the tests or warnings. However, this does not indicate that rule transformation was incorrect or unnecessary: most projects had failing tests and errors unrelated to API usage, which can obscure the effect of rule application. Overall, the resolved warnings and failures demonstrate MELT's potential to help developers more easily maintain large projects.

4.4 Discussion

In this section, we address the main limitations of our approach.

4.4.1 Limitations and threats

4.4.1.A Rule correctness.

We used manual validation to assess rule correctness, with a process that entailed high IRR kappa indicating agreement. One approach for further validation could involve upgrading client projects to newer library versions and applying the rules on projects using these libraries. In RQ4, we use this method to demonstrate that some rules are indeed correct. However, this process is challenging. MELT does not mine rules for all breaking changes in a given release, so upgrading client projects may break multiple aspects in ways automatic find-and-replace rules cannot address [109]. However, automating a large part of migration in ways that entail minimal additional technology or effort on the part of the client developer holds promise for reducing the challenge of upgrading library dependencies. Our rules could also potentially be validated using differential testing techniques or by requesting more tests from the code generation model. However, it is also important to note that we are limited by the expressiveness of the language in which we represent the changes.

4.4.1.B Code generation model.

Our approach relies on a code generation model to generate examples when none are available. We selected GPT-4, a state-of-the-art model trained on data before September 2021. We successfully evaluated on pull requests opened after September 2021, demonstrating the risk of data leakage in these experiments is low. The model, however, is paid and not open-source. As AI research advances, we anticipate better models being made public. We opt for a model-based code generation approach over generating comby rules directly because rules can be validated with code examples (if the code does not pass, we discard the example). Additionally, the model is not fine-tuned and has limited exposure to comby, and is likely to work better on commonly-used languages like Python. For less popular APIs, however, fine-tuned versions of the model on library code might be necessary.

4.4.1.C Generalization.

Our generalization procedure removes context and arguments that appear unrelated to the change, only considering diffs. Removing too much context and type information may result in spurious rules. Conversely, insufficient generalization can make the rule too specific. This limitation stems from the expressiveness that comby language provides, rather than MELT’s approach. Regardless, MELT can return both kinds of rules to the user (i.e., specific or generalized), allowing them to decide what to keep. Currently, developers must manually validate rules to ensure they make sense. To facilitate this, we developed a CI solution on GitHub for integrating our tool. Rules can be validated and modified, if necessary, by whoever merges the PR, or automatically validated, as previously discussed.

4.4.2 Comparison against prior work

Few API migration tools target Python, challenging direct comparison to prior work. MELT adapts its inference algorithm from InferRules [101], designed for type migration in Java. Consequently, MELT without generalization and filtering serves as a baseline equivalent to InferRules. The most closely related approach, PyEvoLve [103], builds on InferRules using comby as an intermediate representation. PyEvoLve focuses on general refactoring, and adapts rules to different control variants, requiring more complex client code analysis. This is in contrast to MELT’s lightweight approach, which aims to minimize overhead on client developers. Since most of our rules are 1:1 and 1:n transformations, adapting rules for control flow variants is less relevant. Overall, while PyEvoLve is more powerful in the types of rules it can infer, fundamentally it serves a different goal as compared to MELT.

Our evaluation differs from closely-related prior work [18, 21] in two ways. First, our manual validation process is able to consider more information in the form of the PR and library documentation. That is, rather than looking at rules in isolation or limiting attention to syntactic validity, we can consider whether the change actually reflects PR intent. Second, we provide an end-to-end evaluation of automatically inferred rules on a number of client code repositories, complementing manual rule validation.

As we discuss in Section 2.3, most prior approaches for automatic API migration (or code evolution generally) mine migration examples from client projects or their source control histories. MELT relies solely on the changed library, looking at internal code changes to inform rule mining. This allows MELT to apply earlier in the library update process. However, libraries do not always include sufficient changed code examples to inform migration, which is why MELT also prompts an LLM to generate extra examples, along with tests to validate those examples. Other approaches may also benefit from using LLMs this way, particularly those whose use cases entail fewer available examples, like A3 [19] (focusing on Android API migration), or APIFix [18] (evaluated on changes to library code, similar to MELT). APIFix in particular could likely benefit from the LLM-generated examples and tests, because it uses edit examples in its program synthesis algorithm. Other tools are evaluated across many more example changes to client code, like Meditor [21]. These approaches

may not require new examples, but leveraging LLMs may allow them to apply earlier in the update process, or in scenarios where migration examples are scarce. Indeed, as models with larger context windows become available (e.g., CLAUDE 100K token context [114]), it becomes possible to include more comprehensive data in prompts, such as full API documentation. This suggests a promising avenue for generating higher-quality, context-rich examples for rule mining, particularly when extant migration examples are scarce.

4.5 Key Takeaways and Contributions

In this chapter, we demonstrated that the library development process itself can be used to mine migration rules to address simple API breaking changes. It is not necessary to rely on commit data from client projects that have already undergone migration. This assumption had previously limited the migration approaches, especially immediately after a library version is released. By integrating the mining process directly with the library workflow in CI, migration rules can be provided to clients immediately after the source code is updated.

MELT represents migration rules in `comby`, a lightweight match-replace language for large-scale code transformation. Since mining can result in very specific rules, we developed a generalization procedure to increase rule applicability. However, this approach comes with trade-offs: some rules are overly generalized and thus may produce incorrect code, while others are too specific. Although the work was successful, some migrations could not be expressed, primarily due to limitations of the language we used for the synthesis (`comby`). In the next, we further explain and exemplify these problems, and introduce a more suitable language to express complex API migrations.

5

A Language for Automating Logically Related Changes

Contents

| | |
|---|----|
| 5.1 Motivation | 50 |
| 5.2 Preliminaries | 52 |
| 5.3 Language Syntax and Overview | 54 |
| 5.4 Language Runtime | 58 |
| 5.5 Evaluation | 61 |
| 5.6 Discussion | 69 |
| 5.7 Key Takeaways and Contributions | 71 |

In this chapter, I present my published and completed work on a new code transformation language. This language is designed to allow expressing flow, dependencies, and the composition of match-replace rules. This effort is driven by the observation that most code changes (including library migrations) tend to be cascading and interconnected; yet, modern languages for code transformations do not inherently offer support for expressing sequences of changes. Indeed, we faced this problem in Chapter 4, where finding

a balance between overgeneralization and overly specific rules was a persistent challenge, largely due to expressiveness limitations stemming from the comby language.

5.1 Motivation

Library migrations usually result in a web of cascading and interdependent code changes that span and propagate across multiple files or repositories [115]. For example, consider the library migration from Figure 5.1, where the goal is to replace `log4j` (a logging library for Java) with an alternative `slf4j`. We divide the migration in four steps:

1. **Migrating Imports:** Replace the import statement for the `Logger` type, `org.apache.log4j.Logger`, with its `slf4j` equivalent, `org.slf4j.Logger`.
2. **Migrating Instantiation:** In contrast to `log4j`, in `slf4j` logger objects are instantiated using a factory method pattern (implemented in the `slf4j.LoggerFactory` class). Therefore, it is necessary to replace the `Logger.getLogger` call with `LoggerFactory.getLogger`.
3. **Adding Missing Import:** Since `getLogger` is a method of a different class (`slf4j.LoggerFactory`), it is also necessary to include the appropriate import statement for this class.
4. **Migrating Associated Method Calls:** Migrate all method calls associated with the logger to their `slf4j` equivalents. In this example, it is only necessary to update the `logger.info` usage.

While it might be feasible to express these migration steps as individual transformation rules (for example using comby [98]), such a strategy is suboptimal for multiple reasons.

Runtime Performance. Applying transformation rules indiscriminately to an entire code base, regardless of the file's relevance to the type or object in question, can lead to significant overhead and result in slow and resource-intensive tool executions. For example, suppose our goal is to migrate a code base from `log4j` to `slf4j` as described above. Here, we notice that despite there being multiple steps in the migration, only files with logger objects need to be touched. One proxy for detecting such files is by looking up an import statement, as imports typically indicate usage. This means that our migration scripts can be designed to only attempt to migrate a file (and execute the rules corresponding to the migration), if an import statement to `log4j` is present. This is particularly relevant when dealing with large-scale migrations across codebases with millions of lines of code.

Accuracy and Precision. Match-replace rules typically provide limited control over which code should be transformed. For example, the import statement in (Line 2, Step 3) is necessary only if the `getLogger` API is migrated within the file and the import is not already present. Additionally, the migration of the `info` API (Line 8, Step 4) should only affect the logger object that was migrated in (Line 5, Step 2). Calls to other objects with an `info` API from another object should not be affected. Writing such constraints in

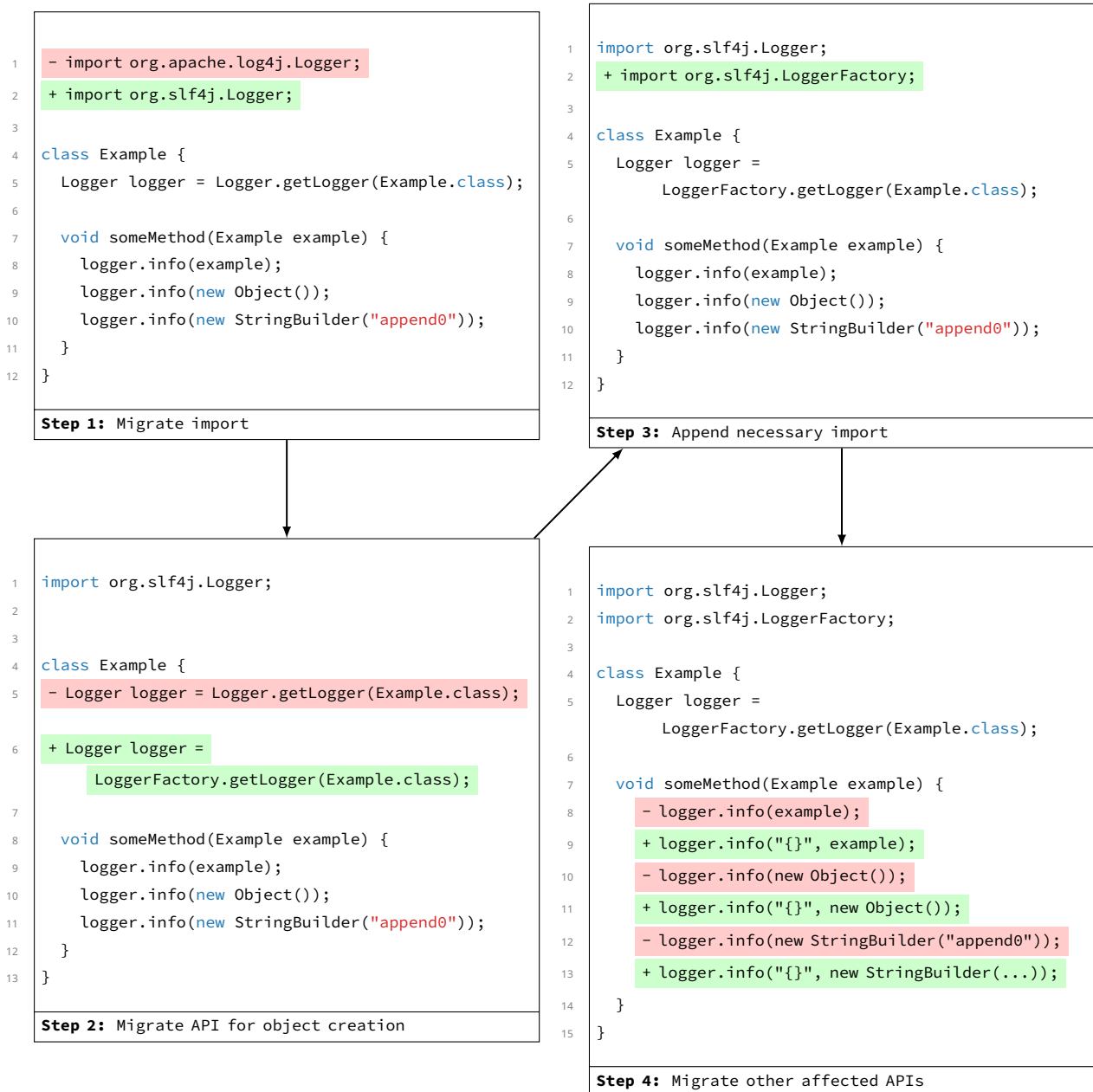


Figure 5.1: Migration steps to move from two popular logging libraries in java: log4j and slf4j.

the comby language would require additional scripting. An attempt at writing this rule in comby could be `: [x].info(: [args]) -> : [x].info("{} ", : [args])`. However, such a rule would apply across the entire codebase and not just to the previously migrated logger object. On the other hand, substituting `logger` for `: [x]` would overly specialize the rule to this example: `logger.info(: [args]) -> logger.info("{} ", : [args])`. This is particularly challenging to address when migrating common API names like `append` and `concat` as discussed in Chapter 4.

Expressiveness. Current techniques cannot express complex automations as match-replace rules. Generally, to perform complex automations, developers have to rely on imperative code transformation frameworks instead. These frameworks (e.g., [116, 117]) provide APIs for manipulating code at the AST level, allowing for arbitrary code transformations. The APIs let users control where, when, and how code should be rewritten based on context, symbol information, and more complex analyses. However, imperative frameworks present a twofold problem. First, imperative frameworks are typically monuments of engineering, demanding significant time and effort to learn [118]. Second, the frameworks are typically language-specific, as they rely on compiler and build infrastructure to be able to get meaningful information from the code. This results in significant burdens as multiple developer experts are necessary for automating the same task in different languages [119].

5.2 Preliminaries

5.2.1 Existing Code Transformation Languages

Frameworks for automating code transformation vary widely. At one end of the spectrum, *lightweight techniques* [23, 63, 120] offer declarative languages to rewrite code with simple match-replace rules. The key advantage of lightweight techniques is language agnosticism, which stems from these techniques being independent of the underlying compiler infrastructure. Moreover, match-replace rules are often syntactically close to the target language, making them easy to write and use [23]. However, lightweight techniques are often limited to atomic context-free code changes, lacking support for tasks requiring cascading and interdependent code changes. On the other hand, *imperative frameworks* [121] for AST-level manipulation allow for arbitrary code transformations. As discussed earlier, imperative frameworks provide infrastructure and finer control over where, when, and how code should be rewritten based in analyses. These frameworks may not be the best fit for learning migration scripts due to their large APIs, and their imperative nature. Synthesizing arbitrary imperative programs is undecidable.

5.2.2 A novel Lightweight Language for Cascading Transformations

To address the limitations of existing lightweight match-replace tools (as described in 5.1), we have developed a new language called POLYGLOTPIRANHA. At a high level, POLYGLOTPIRANHA allows for sequencing of rule applications as well as propagating information across rules using a *directed graph of match-replace rules*.

Programs in POLYGLOTPIRANHA are graphs, where the nodes represent individual transformation rules (expressed in a code transformation language of choice), and the edges determine the order for applying these rules. Each edge is also associated with a label that defines the scope within which the target rule is applied with respect to the source rule. For example, an edge $\mathcal{R}_1 \xrightarrow{\text{class}} \mathcal{R}_2$ reads as, "Apply rule \mathcal{R}_1 and then

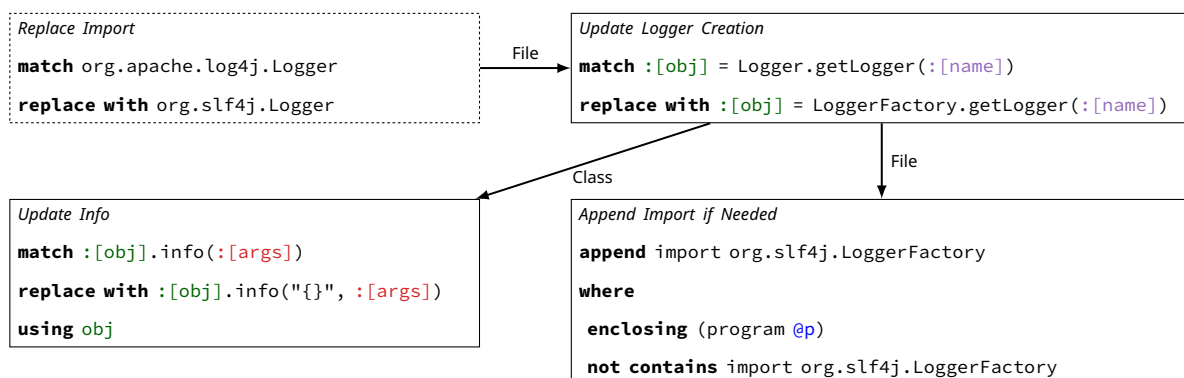


Figure 5.2: Program in the DSL to described the migration for Figure 5.1. The dashed line represent a seed rule, i.e., the rule that triggers the migration. Each edge is annotated with a scope. The scope determines where the target rule will be applied with respect to the source.

apply rule \mathcal{R}_2 within the enclosing class where \mathcal{R}_1 was applied." The ability to cascade transformations using match-replace rules makes it an ideal candidate for writing API refactoring rules.

Note that the goal of this research is not to design a new matching syntax language per se, but rather a meta-language that allows for the combination and interleaving of existing match-replace languages, as well as provide infrastructure to combine and sequence them. We build our approach on top of this idea of graph of match replace-rules. In POLYGLOTPIRANHA:

1. We aim to allow users to express match-replace rules in *any* source code matching language. We provide alternative syntaxes, catering to different users and different needs. The current implementation supports the *tree-sitter* query language, *comby*-like syntax, as well as regular expressions.
2. We enable match-replace languages to be interleaved in the graph (e.g., the first rule uses regular expressions, while the second uses *comby*). Some tasks might be accomplished using simple regular expressions, but for other *comby* might be a better option.
3. We support composition of match-replace rules using a set of language-agnostic filter primitives. The goal is to enhance rule precision by leveraging the surrounding code context. Instead of writing a match-replace rule, the idea is to write multiple rules and only transform the code if all conditions are satisfied. In this sense, filters act as a set of pre-conditions.

5.2.2.A DSL Motivating Example

Figure 5.2 illustrates a program to automate the migration process outlined in Figure 5.1 in our proposed language. In this example, rules are expressed in the *comby* syntax as explained in Chapter 4. A program in our DSL is a graph of match-replace rules. The graph has a source/seed rule `Replace Import`, i.e., this rule initiates and triggers the code transformations by replacing the `log4j` import. `Replace Import` is connected

| | |
|---|--|
| <pre> <rule_graph> ::= <rule>+ <edge>* <edge> ::= from string to (string, scope) <scope> ::= Global File n-Ancestors Method Class <rule> ::= name string match match [replace [template_variable] with <replace>] [where <filters>] [using <holes>] [is_seed bool] </pre> | <pre> <replace> ::= string <replace> template_variable <replace> <> <filters> ::= enclosing match [<contains>] <filters> not_enclosing match <filters> <> <contains> ::= contains match [at_least int] [at_most int] not_contains match <> <holes> ::= template_variable <holes> <> </pre> |
|---|--|

Figure 5.3: Syntax of our DSL for cascading code transformations. The elements inside square brackets are optional. The symbol `match` is an expression for pattern matching (e.g., `comby`); the symbol `template_variable` represents named capture groups from the match pattern.

to another rule `Update Logger Creation` with an edge labelled `File`. This means that the rule `Update Logger Creation`'s application is restricted to the files affected by the `Replace Import` rule.

Furthermore, `Update Logger Creation` has two outgoing edges: (1) `Update Info` - for updating the info API usage between the libraries, and (2) `Append Import if Needed` - for appending the import of `LoggerFactory` if necessary. First, we seek to update the info API usage only for the `Logger` field that was previously updated. Thus, `Update Info` takes as input the `:[obj]` name from the previous rule, as indicated by the `using` keyword. The value of `:[obj]` will be instantiated at runtime based on the previous rule. Moreover, the rule will only be applied within the same `Class` as indicated by the edge label. Second, the import is appended to the `File` where the `getLogger` API was applied as indicated by the edge between `Update Logger Creation` and `Append Import if Needed`.

Notice that the rule `Append Import if Needed` uses an extra clause called *enclosing*. Enclosing is a filter that ensures the import is not added to the file twice. In this case, the filter is written in a combination of *tree-sitter* and regular expressions. The *enclosing* pattern (`program @p`) is a *tree-sitter* query that matches the entire source file. The *not contains* clause specifies a regular expression to check against enclosing source file. The rule will only match if the source file does not already contain the import statement.

5.3 Language Syntax and Overview

Figure 5.3 describes the grammar of our DSL. A program in the DSL is a graph of match-replace rules. The rule graph is captured as a list of *directed* and *labelled* edges. Each node represents an individual transformation rule that structurally matches and rewrites code. Rules can also just match code without transforming it. The edges between rules specify *which* rule to apply next and the scope where it should be applied.

5.3.1 Edges

As shown in Figure 5.3, the edges are *directed and labelled*. Each *edge* connects either two rules or a rule to a rule group, defining the order in which they should be applied, akin to the `andThen` operator¹. The edge label specifies the *scope* of application, selecting the portion of the code base upon which the target rewrite rule is applied, with respect to the code that the source rule matched. For example, given an edge from $\mathcal{R}_1 \xrightarrow{\text{method}} \mathcal{R}_2$, reads as “apply \mathcal{R}_1 and then apply \mathcal{R}_2 within the enclosing *method* where \mathcal{R}_1 was applied”.

The DSL supports two *language-agnostic* predefined scopes as shown in Figure 5.3. 1. *Global* the target rule is applied across the codebase, 2. *File* the target rule is applied in the enclosing file. It is also possible to support other *language-specific* scopes that depend on the granularity of the internal representation of code within the implementation. The current implementation represents code internally with `tree-sitter` [122], and supports three other scopes: 3. *n-Ancestors* the target rule is applied to the parent parse tree nodes of the code fragment that the origin rule matched, 4. *Method* the rule is applied to the enclosing method where the preceding rule was applied, and 5. *Class* scope refers to the enclosing class. Note that the language specific scopes are set up only once per language.

5.3.2 Rules

Besides the name, a match-replace rule has 4 major components 1. *match* - a pattern to match source code, 2. *replace* - a pattern to rewrite the matched code, 3. *filter* - to filter out certain matches based on the surrounding code, and 4. *holes* - variables referenced in the rule, these are filled at run time and serve as the *dynamic* component of the rule. Furthermore, a rule could be a *seed_rule*. These seed rules are entry points to the graph. This graph is traversed in a depth-first manner at each location where the rule was applied. A valid rule graph contains at least one *seed rule*.

5.3.2.A Match.

The *match* expression is a declarative pattern that captures a code snippet with a specific *structure* or *shape* (based on its parse tree). The match also labels portions of the matched parse tree like the *named captured groups* [123] in regular expressions. Our DSL can support multiple structural matching languages, as long as they support named capture groups (used to label portions of the code). The current implementation supports concrete patterns, structural queries, and regular expressions.

Concrete Patterns: A concrete pattern is a string with template variables / holes, that is matched to concrete syntax nodes [124] from the program’s parse tree². Formally, let s be a concrete pattern containing holes of the form $:[var1]$, where each hole can represent syntactically valid sub-trees. A Concrete Syntax Tree

¹<https://docs.oracle.com/javase/8/docs/api/java/util/function/Function.html#andThen-java.util.function.Function>

²We use Concrete Syntax Trees (CST) over Abstract Syntax Trees (AST) because we must preserve all syntactic structures within the source code, which are necessary for source code matching

| Rule | Matched code snippet | Capture Groups |
|--|---|---------------------------------------|
| <code>match :[x].info:[args]</code> | <code>logger.info(example)</code> | <code>x: logger, args: example</code> |
| <code>match import org.apache.:[name]</code> | <code>import org.org.apache.log4j.Logger</code> | <code>name: log4j.Logger</code> |

Figure 5.4: Example rules using concrete patterns applied to motivation example in Figure 5.1.

| Rule | Matched code snippet | Capture Groups |
|---|--|--|
| <pre>match (method_invocation object: (identifier) @obj name: (_) @m arguments: (object_creation))</pre> | <code>logger.info(new StringBuilder("append0"))</code> | <code>obj: logger, m: info</code> |
| <pre>match ((field_declaration (name: (identifier) @obj (value: (method_invocation name: (_) @m))) (#eq? @m "getLogger"))</pre> | <code>Logger logger = Logger.getLogger(A.class)</code> | <code>obj: logger, m: getLogger</code> |

Figure 5.5: Examples of rules using simplified structural queries applied to motivation example in Figure 5.1.

(CST) node t matches s if, traversing t in depth-first order yields leaf nodes with a string representation that aligns with s from left to right. Each hole can represent entire sub-tree structures (i.e., multiple sequential leaf nodes under an internal node). This paradigm of matching is supported by multiple other tools (e.g., [23, 63]). POLYGLOTPIRANHA adopts the syntax proposed by [23] in their tool Comby. However, our concrete patterns have stricter semantics compared to Comby. In our concrete pattern, a template hole, `:[x]`, matches whole syntactic structures / CST nodes, whereas Comby templates can represent arbitrary strings. Figure 5.4 shows two examples.

Structured Query Language: A query consists of one or more patterns, where each pattern is an *s-expression* that matches a certain set of nodes in a parse tree. These queries capture the structure of the target pattern in terms of AST node types and string based predicates. This paradigm is programming language agnostic, and is supported by systems like `tree-sitter`. POLYGLOTPIRANHA supports the *s-expression* based `tree-sitter` queries [125]. Figure 5.5 shows two examples.

Each matching paradigm has distinct advantages and disadvantages. By construction structural queries are more precise than concrete syntax because they can leverage node-types or absence of particular nodes, and therefore leave less room for ambiguity (e.g., it is possible to differentiate between a field and a local variable declaration). For example, matching method declarations is easier with structural query, because we

| Rule | Source Code Update |
|---|---|
| match <code>:[object].info(:[args])</code> | <code>logger.info(new Object())</code> |
| replace <code>:[object] with other_: [name]</code> | <code>other_logger.info(new Object())</code> |
| match <code>import org.slf4j.: [rest]</code> | <code>- import org.slf4j.Logger</code> |
| replace with <code>-</code> | <code>- import org.slf4j.LoggerFactory</code> |

Figure 5.6: Examples of replacement rules using concrete syntax. Notice that in the first example, the code is only partially rewritten. Whereas in the second example, the entire code snippet is deleted.

would not need to account for all its syntactic variations (e.g., modifiers like `public`, `static`, `final`) like in concrete syntax. In contrast, matching API invocation pattern like `logger.info(example)` (from Figure 5.1) the *concrete pattern* is convenient and more succinct. The structural query for this pattern is verbose, and requires knowledge of the target language's grammar. Regex matching is more suitable for semi-structured documents like markdown files. Note that POLYGLOTPIRANHA is not tied to these three languages, more can be supported.

5.3.2.B Replacement.

The replacement pattern decides on how a matched code snippet should be transformed. It is possible to either replace the entire matched code or just segments identified by a named capture group. The replacement expression / pattern can be seen as partial function that is instantiated at run time by substituting a referenced named groups or template variables with their values from either the initial match in the rule, or inputs to the rules declared with the `using` keyword (i.e., code snippets captured in previous rule applications). In Figure 5.6, we show two examples of replacement rules. In the first one, the code is only partially rewritten. In the second example, the matched code is completely deleted because no target node is specified.

5.3.2.C Filters.

To make the *rules* more precise and context-aware, our DSL provides `filters` to control the application of a rule based on the surrounding code. First, the candidate code to transform is checked against the `matcher` of the rule. Then, at each matched location, the filters will check if the surrounding code of this location satisfies certain criteria.

There are two primitive filters: 1. `enclosing` – checks if the primary match is enclosed by a parse tree node that satisfies the given matcher, and 2. `not_enclosing` – checks if the primary match is *not* enclosed by parse tree node that satisfies the given matcher. The `enclosing` filters can be further refined by specifying `contains` and `not_contains` expressions. The `contains` (`not_contains`) expressions specify `matchers`

| Rule | Source Code Update |
|--|--|
| Delete unused local variable match <code>:[var_name] = :[rhs];</code> replace with <code>-</code> where enclosing (method_declaration) contains <code>:[var_name] atleast 1</code> | <pre>def some_python_function() { - int unused_variable = 42 execute() }</pre> |
| Add import statement if absent match (import_declaration) <code>@p</code> replace with <code>@p \n import org.slf4j.LoggerFactory;</code> where enclosing_node (compilation_unit) not_contains <code>import org.slf4j.LoggerFactory</code> | <pre>import org.slf4j.Logger; import org.slf4j.LoggerFactory; class A { Logger logger = LoggerFactory.getLogger(A.class); }</pre> |

Figure 5.7: Example rules using *filters*. Note how these rules leverage both *concrete pattern* and *structural query*. In the first example, we use a contains filter inside the enclosing method declaration. This allows us to check if a variable is used only *once*. If this is true, the usage corresponds to its declaration, and thus, can be safely deleted. In the second example, the import is added only if it is not already in the code, as indicated in the not contains predicate. Since the import is already present, the code is not rewritten.

that should (not) match at least once inside the `enclosing_node`. The user can also specify the frequency of these matches with `at_least` and `at_most` attributes.

Holes. These serve as dynamic components within a rule. They describe input variables to the rule. At run time, their corresponding values are populated from a symbol table (which maintains the bindings from named captured groups to code snippets from current and previous applications). The example in Figure 5.2 showcases the usage of these holes. The rule `Update_Info` declares the hole `:[obj]`, which will be instantiated during the transformation based on the code matched in the previous rule `Update_Logger_Creation`.

5.4 Language Runtime

5.4.1 Algorithmic Overview

Algorithm 3 provides a high level overview for the language implementation and runtime. The core idea is to maintain a queue of *seed rules*, and traverse the graph and the files in the codebase starting from each seed rule. First, we validate the rule graph to prevent unexpected behavior using a data-flow analysis and syntactic checks on the rules (Line 1). After the validation, we push the seed rules into a global queue and initialize an environment / symbol table with the input substitutions (Line 4 - 5). The environment is used to store both

Algorithm 3 Core procedure for transforming codegiven a graph of rules

Input:

```
( $\mathcal{R}$  : RuleGraph,  $S$  : substitutions)
 $C$  : path to codebase
1: if  $\neg$ VALIDATE( $\mathcal{R}, S$ ) then
2:   return
3: end if
4:  $Q \leftarrow$  SEEDRULES( $\mathcal{R}$ )
5:  $env \leftarrow S$ 
6: while NOTEMPTY( $Q$ ) do
7:    $rule, _ \leftarrow$  POP( $Q$ )
8:   loop
9:      $isApplied \leftarrow$  false
10:    for file in RELEVANT( $C, rule, env$ ) do
11:       $isApplied \vee =$ 
12:      EXECUTERULEGRAPH( $rule, file, \mathcal{R}, Q, env$ )
13:    end for
14:    if  $\neg isApplied$  then
15:      break
16:    end if
17:  end loop
18: end while
```

Algorithm 4 EXECUTERULEGRAPH function

```
1: function EXECUTERULEGRAPH( $rule, file, \mathcal{R}, mut\ Q, mut\ env$ )
2:    $rulesStack \leftarrow [(rule, file)]$ 
3:    $isApplied \leftarrow$  false
4:   while NOTEMPTY( $rulesStack$ ) do
5:      $rule, scope \leftarrow$  POP( $rulesStack$ )
6:      $rule \leftarrow$  INSTANTIATE( $rule, env$ )
7:     while HASMATCH( $rule, scope$ ) do
8:        $match \leftarrow$  GETMATCH( $rule, scope$ )
9:        $isApplied \leftarrow$  true
10:      APPLYEDITS( $match, rule, mut\ env$ )
11:      for  $rule, scScope$  in SUCCESSORS( $rule, \mathcal{R}$ ) do
12:        if  $scScope \equiv$  GLOBAL then
13:          PUSH( $Q, (rule, GLOBAL)$ )
14:        else
15:           $scope \leftarrow$  RESOLVE( $match, file, scScope$ )
16:          PUSH( $rulesStack, (rule, scope)$ )
17:        end if
18:      end for
19:    end while
20:  end while
21:  return  $isApplied$ 
22: end function
```

the initial set of substitutions as well as the captured groups of from rule executions, which can be used as dynamic elements in subsequent rules. Each seed rule is applied across the entire codebase recursively in a depth-first fashion (Line 6), until no rules match (Line 8 - 15). For each *relevant* file (e.g., a file that is likely to contain the match template of the rule, see Section 5.4.1.C), we invoke `ExecuteRuleGraph` (Algorithm 4). In this step, the tool traverses over the CSTs and transforms the source code. For each match, it explores the rule graph and stacks the rules in a DFS-manner (Line 12), applying them exhaustively within the scope. The function `ExecuteRuleGraph` is *not pure*, it updates the environment, transforms the source code in-place, and pushes new rules into the queue (Q). We detail each function of the algorithm more thoroughly in subsequent sections.

5.4.1.A Graph Validation

The first step in the core algorithm is to verify the graph (Line 1). In our implementation, POLYGLOTPIRANHA statically validates the constructed graph to prevent unexpected behavior when the graph is applied to the codebase. First, POLYGLOTPIRANHA checks if the individual rules' matchers and filters are well-formed. For example, POLYGLOTPIRANHA ensures that each regex compiles and that each s-expression parses correctly according to the language's grammar. It also conducts a data-flow analysis to ensure that no path in the graph traversal leads to a rule where an input variable is not initialized correctly. This is implemented as a definite assignment analysis [126]. If the graph is incorrect, POLYGLOTPIRANHA alerts the user to prevent panics that could result from accessing undefined variables.

5.4.1.B Environment

The environment is a simple symbol table, which is initialized with the substitutions from the program (Figure 5.3). Rules can access symbol table variables if they have been declared. If a rule is triggered and a match is found, the symbol table is updated by binding the matched source code to the corresponding named captured group in the symbol table. If a variable already exists in the symbol table, its entry gets over written. Therefore a rule always gets instantiated with the most recent binding of the referenced symbol from the environment. This kind of dynamic variable scoping can also be observed in languages like LaTeX or Bash.

5.4.1.C Relevancy check for performance

In rewriting large code bases, repeatedly parsing the entire codebase is inefficient, especially in monorepos with millions of lines. The goal of the function `relevant` is to optimize code rewriting by only parsing files whose content matches the concrete values assigned to the holes of the *global rules* (Line 10). In practice, the concrete values to the input substitutions are used to filter out files that are not relevant to the transformation using string matching. This simple insight improves POLYGLOTPIRANHA's overall performance. The implementation of POLYGLOTPIRANHA further boosts this by parallelizing the lookup using fork-join frameworks (like Comby). Note that, POLYGLOTPIRANHA circumvents this optimization for holes that are referenced inside the `not_contains` or `not_enclosing` clause.

5.4.2 Rule Graph Execution

Algorithm 4 describes the procedure `EXECUTERULEGRAPH` that applies a given *rule* across a *file*. Each time a *seed rule* is triggered, we initialize a *stack* (`ruleStack`) for depth-first traversal of the rule graph (Line 2). Then, we pop rules from stack and apply each rule exhaustively within the specified scope (until `hasMatch` is `false` as shown in Line 7). For each match, we update the environment with the new capture groups and transform

the source code by applying the rule (Line 10). Finally, we add the successors of the current rule in the graph to the local stack or the global queue, and continue this until fix point (Lines 11 - 18).

5.4.2.A Optimizations.

POLYGLOTPIRANHA uses the `tree-sitter` [122] framework for parsing the source code. POLYGLOTPIRANHA maintains only one parse tree *object* in its memory, and updates this object sequentially leveraging the `tree-sitter`'s incremental parsing feature. This eliminates the need to parse the file again from scratch after the rewrite, thus optimizing POLYGLOTPIRANHA's overall performance. Additionally, to minimize the impact on the parse tree, by default our approach 1. orders the rules from inner to outer scope: starting from the parent, to method, class, file, and finally to global scope, and 2. rewrites code bottom up. In the future, we plan to support alternative transformation strategies.

5.5 Evaluation

Our language implementation is merged into the POLYGLOTPIRANHA repository, which is maintained on GitHub. The tool is used internally at Uber with multiple use cases. In our evaluation, we aim to show key desirable properties of the language for our use case. To do this, we evaluate POLYGLOTPIRANHA on three case studies related to migration and code cleanup. In particular, we answer the following research questions:

- RQ1. [Expressiveness]** How expressive is the DSL for real-world code transformation tasks? *We assess this through three case studies. We highlight the complexity of each, and how to encode it in the DSL.*
- RQ2. [Effectiveness]** How effective is our language at automating code changes? To what extent is it useful in practice? *We run the above tools across Uber's proprietary codebase, and measure the percentage of Pull Requests (PRs) that pass Continuous Integration (CI) and are merged without manual intervention. For PRs with intervention, we measure the LoC changed by tools versus developer.*
- RQ3. [Comparison with state-of-the-art]** How do tools built upon the DSL compare to similar tools built upon state-of-the-art frameworks? *We compare the POLYGLOTPIRANHA-based implementation against the imperative variants developed upon `ErrorProne` [116] and `OpenRewrite` [127], and against its declarative variants developed upon `Comby` [23] (a lightweight tool). We compare implementations in terms of size, complexity and performance.*

5.5.1 RQ1. Expressiveness

5.5.1.A Experimental Setup

To showcase the *expressiveness* of the DSL, we present *three* real-world case studies where we automate complex code transformation tasks using POLYGLOTPIRANHA. In each case study, we highlight the complexity of

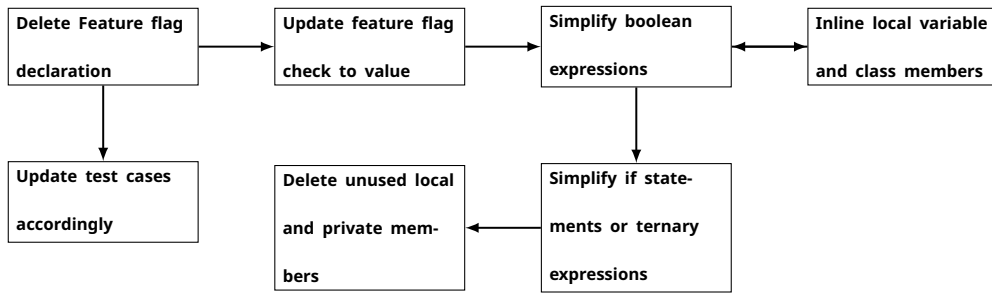


Figure 5.8: Strongly connected components of the rule graphs for feature flag cleanup. The graph structure is language-agnostic. Implementations across languages require some adaptations.

the task, and how the DSL can be used to encode it. We chose these three case studies because they are high-impact tasks crucial to Uber’s operational needs and they are representative of the tasks that Uber or other software companies would want to automate. Moreover, these tasks are not trivial to automate using existing frameworks.

5.5.1.B Case study: Stale Feature Flag Cleanup

Feature flagging is a widely adopted and highly encouraged practice at Uber³, and other major software companies [128, 129]. It allows developers to modify configurations without redeploying, supporting A/B testing in production. However, feature flags often become stale, and retaining them beyond their original purpose can lead to technical debt. Therefore, it is important to automate their removal. Indeed, researchers [119] have developed the Piranha tool for this purpose. Piranha is built on top of the ErrorProne [116] frameworks for java and SwiftSyntax [130] for Swift. However, Uber’s codebase uses Kotlin and Go too. Instead of developing two new language-specific tools, we used POLYGLOTPIRANHA to implement this transformation as **one** tool supporting java, Kotlin, Swift and Go.

Figure 5.8 shows the strategy that we implemented for automating the cleanup of stale feature flags at Uber. Each node in this figure is a strongly connected component or sub-graph of the original large graph implementing the transformation. Here, each subgraph is a cleanup category. For instance, `Simplify boolean expressions` contains rules that simplify nested boolean expressions with conjunctions, disjunctions and negations. These rules are recursively applied until the expression cannot be further simplified. It should be noted how the `Simplify boolean expressions` and `Inline local variables and members` call each other, until no more simplification is possible. The `Cleanup tests` sub-graph is particularly interesting. In this sub-graph we identify all the tests that explicitly set the feature flag to a specific Boolean value. If the set value is the same as the status of the feature flag we elide the setter, else we delete the test case.

³In fact, our motivating example is a simplified version of feature flag cleanup we performed internally.

```

1 - public enum IUIModesEnum {
2 + public interface IUIModes {
3   - DARK_MODE,
4   + @Param(key="DARK_MODE")
5   + BoolParam isDarkMode();
6   - LIGHT_MODE,
7   + @Param(key="LIGHT_MODE")
8   + BoolParam isLightMode(); }

```

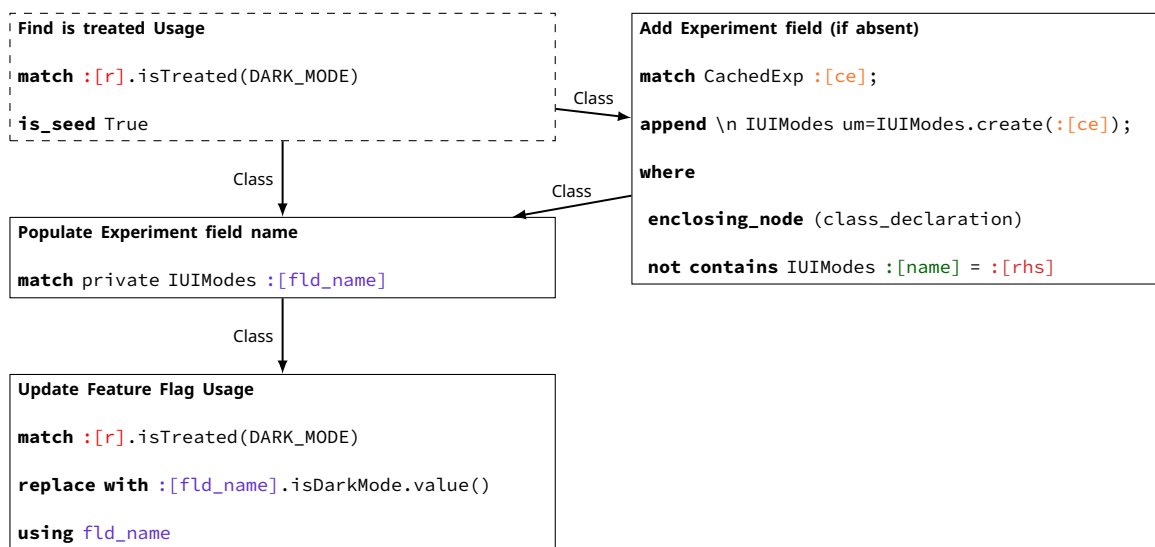
(a) Example migration from enum-based feature flag declaration to annotations.

```

1 class Consumer {
2   CachedExp ce = new Experiment();
3 + IUIModes um = IUIModes.create(ce);
4   public String color() {
5 - return ce.isTreated(DARK_MODE)
6 + return um.isDarkMode().value()
7     ? "Black" : "White";
8   }
9 }

```

(b) Source code update after the migration of enums to interfaces as shown in Figure 5.9a.



(c) Part of the original rule graph that migrates usages of the `isTreated` API. The input substitutions in the bottom right instantiates this graph to migrate the `DARK_MODE` feature flag described in this figure.

Figure 5.9: Experimentation API usage update after the migration from enum-based feature flag declarations.

5.5.1.C Case Study: Experimentation API Migration

The *Experimentation* team at Uber developed a new *feature flagging* API to support its growing needs. It was imperative for Uber to transform thousands of lines of their *Android* code to use this new API.

Figure 5.9 showcases the code changes required for the migration. The previous *feature-flag* API declared feature flags using enum data types. To adapt the code to the new API, these enums need to be rewritten as *annotated abstract methods* (as shown in Figure 5.9a). These annotations were added to specify metadata information for a feature flag such as *key* and *namespace*. After migrating the enum to an interface, this change has to be propagated. For example, consider the feature flag usage in Figure 5.9b. Previously, the `isTreated` method (Line 5) was invoked to check the status of the feature flag by passing the enum `DARK_MODE`, declared in Figure 5.9a. However, with the new design clients are expected invoke the feature flag method `isDarkMode()` as shown in Line 6.

```

company_kotlin_android_module(
    name = "src_release",
    plugins = [
        "- //libraries/compiler:processor" ,
        "//libraries/utilities",
    ],
    + kotlin_plugins = [
    + "//libraries/processor-kt:processor"],
    tests = [":test_release"],
    visibility = ["PUBLIC"],
)

```

```

- import com.co.ParameterUtils
interface UIParams{
    @JvmStatic
    fun create(cp:CachedParams): UIParams =
- ParameterUtils.create(UIParameters::class.java,cp)
+ UIParamsProvider.create(cp)
}
}

```

(a) Changes to the BUCK file. Here the java dependency is replaced with the Kotlin counterpart

(b) Changes in the source code illustrating the usage of the new Kotlin-based processor

Figure 5.10: Examples of modifications in the BUCK and Kotlin files for the annotation processor migration.

In practice, this migration has to accommodate many other *caveats*. To complete this migration, it is necessary to also add new fields (e.g., `IUIModes` (Line 3, Figure 5.9b). This is handled by writing two rules as shown in Figure 5.9c: 1. Add Experiment field - adds a field of type `IUIMode` (if absent), and 2. Populate Experiment field name captures the name of the field of type `IUIMode`. The field name (i.e.: `[fld_name]`) is used in the following rule Update Feature Flag Usage, which is the actual rule used to replace the `isTreated` API. Other nuances include deleting consequently unused members and imports and adapting test cases accordingly.

5.5.1.D Case Study: Annotation Processor Migration

The goal of this migration is to transition the Android codebase from a java-based annotation processor to a Kotlin-based system to improve overall performance. The changes required for this migration are described in Figure 5.10. This migration requires changing all the build configurations (written in BUCK [131]) to be adapted by replacing the old processor dependency with the new one as depicted in Figure 5.10a. Besides the build files, it is necessary to migrate all Kotlin files that initially used the java processor (shown in Figure 5.10b). For example, Figure 5.10b shows how `ParameterUtils.create` is replaced with a Kotlin equivalent method, `create`, and the unnecessary import statement is deleted.

5.5.2 RQ2. Effectiveness and Usefulness

5.5.2.A Experimental Setup

To evaluate the effectiveness of POLYGLOTPIRANHA's framework, we show its effectiveness on the three case studies above by applying them to Uber's proprietary code-bases. Specifically, the Android codebase is composed of 7.5M LoC of java and 2.5M LoC of Kotlin, while the iOS codebase is composed of 7.5M LoC of

| Application | Language | Effectiveness | | Usefulness | | |
|--------------------------------|-------------------|---------------|-------------------|------------------|-----------------|-------------------------------|
| | | # PRs | # PRs (CI passes) | # PRs (Accepted) | # files updated | #+/- Lines |
| Stale Feature Flag Cleanup | ☪ Java & Kotlin | 2515 | 1413 | 817 | 2952 | +15032 /-107635 [†] |
| | 🍏 Swift | 2186 | 1309 | 614 | 1733 | + 21230 /-104721 [‡] |
| Experimentation API migration | ☪ Java | 155 | 89 | 155 | 2146 | +19157 /-19041 [§] |
| Annotation processor migration | ☪ Kotlin & Python | 25 | 25 | 25 | 2042 | +2809 /-3897 |

[†] 85.7% was automated [‡] 95.3% was automated [§] 73.4% was automated ^{||} 100% was automated

Table 5.1: PRs created and merged by the tool, as well as the % of LOC automatically deleted for each.

Swift. The PRs produced by our tools are reviewed by the appropriate teams, and merged if they pass the *Continuous Integration* checks and tests. The PRs that fail CI are expected to be manually fixed by the respective team before merging. Note that this data represents months of company wide-effort.

5.5.2.B Results

Table 5.1 summarizes the overall results we obtained by running POLYGLOTPIRANHA based tools over our proprietary corpora. For each application, it reports the number of PRs created, PRs accepted (and merged), and PRs that pass the Continuous Integration checks and tests. At large, the three tools produced 4881 PRs in the last six months of which 1611 have been accepted and merged into the main codebase at the time of writing this paper. Particularly, for *stale feature flag cleanup* our acceptance rate is 52.5% (of PRs that pass CI) while for the migrations it is unsurprisingly 100% (because the migrations were orchestrated centrally). These PRs have deleted over 200k LoC of dead code and migrated over 20k LoC of old code to use the new APIs.

A - Stale feature flag cleanup The data for this experiment was collected between April and November of 2023. POLYGLOTPIRANHA created a total of 4701 PRs, and reviewers did some kind of activity on 1727 (36.7%) of the total number of PRs. These activities include, accepting the PR and merging it, commenting the PR, or patching the PR before accepting it. There are still 1410 PRs that pass all CI checks and are still in queue for review. Further, the reviewers have marked 114 PRs as *Needs Changes* status indicating that they expect extra cleanup from the tooling. For most of these PRs, the reviewers have reported issues with new features and bugs. The reviewers abandoned 182 PRs, to assert that the cleaned up feature flags are not stale.

We observed that 56.2% of all the Android PRs and 59.9% of the iOS PRs passed all CI checks. Uber's CI not only builds and tests the change, but it also employs over a hundred linters and bug-checkers to ensure the quality of the change meets the Uber's high standards. These checkers ensure there are no unreachable and unreferenced elements (e.g. `UnusedMethod` check [132]), no sub-optimal code (e.g. `ComplexBooleanCon-`

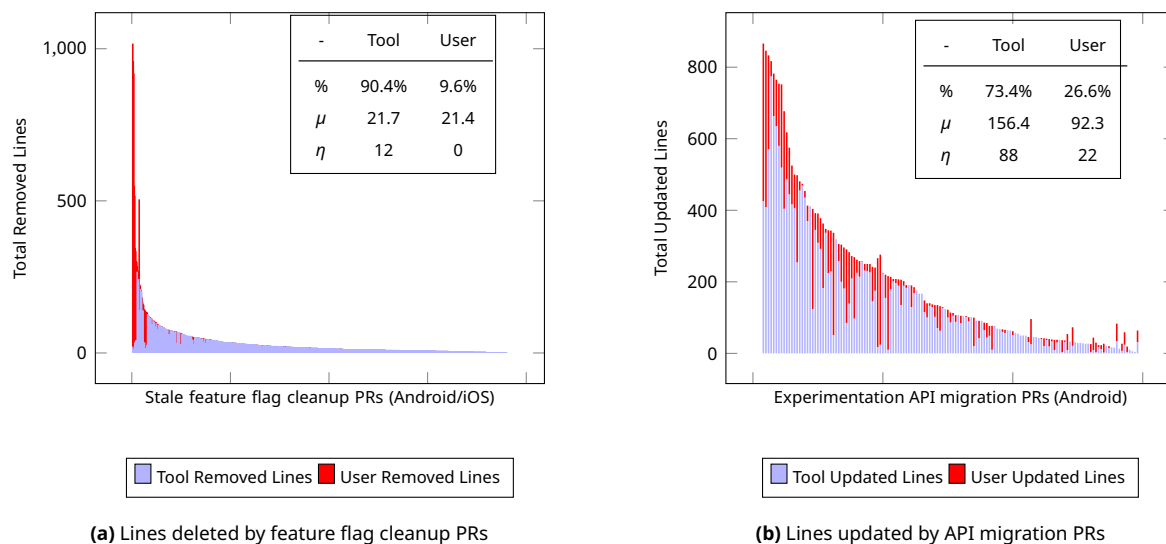


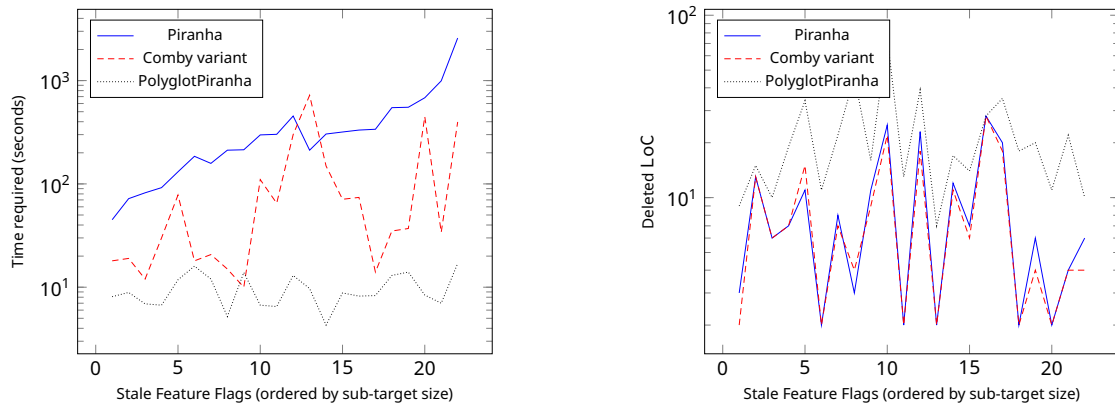
Figure 5.11: Lines deleted/updated by tool (blue) vs users (red)

stant check [133]) and no nullability errors [134].

The tool deleted 85.7% and 95.3% of all the total deleted lines across the Android and iOS codebases (90.4% gross) respectively across all merged PRs, as shown in Figure 5.11a. We observed that 75.9% of the PRs that were merged required no user intervention. However when the developer did intervene, they deleted a lot of code before merging the PR, hence the mean number of lines deleted by user is skewed ($\mu = 21.4$, $\eta = 0$). In few outlier cases developer deleted more than 900 lines of code. Probing further into these outlier PRs, we discovered that developers had removed a collection of top-level classes that were guarded by the flag. Some of these scenarios will be incorporated into the next version of our tool. However, very precise and general support for such cleanups is impractical in our lightweight approach.

B - Experimentation API For the *Experimentation API migration*, we observed that 89 (59.9%) PRs passed all CI checks. The main reason migration PRs to fail was non-standardized usage of the API and usage of some specific API patterns that were not automated. Nonetheless, the tool still automated 73.4% of all lines deleted. The migration was driven centrally by the team, therefore the all the PRs were immediately acted upon after creation. The team reviewed these PRs, patched them if necessary and merged them.

The tool migrated 73.4% of the total lines deleted, however we observed that more than 74.8% PRs needed some manual intervention. In these cases developers on average updated another 92 lines upon the changes proposed by the tool. We observed that Uber developers also made manual changes to the PRs that pass CI. These changes include class deletions, removing unused data files, updating comments and method names.



(a) Time required to perform the cleanup. Flags are ordered by target size. (b) Lines of code deleted for each tool. Flags are ordered by target size.

Figure 5.12: Comparative analysis of Comby, Piranha, and PolyglotPiranha for stale feature flag cleanup.

While refining rules can resolve certain scenarios, some require symbol or type information, and others, such as method renaming and updating documentation, are beyond the scope of traditional tools. We also observed that the team knowingly used the tool to perform partial migrations even for cases where all APIs were not supported. The small spikes towards the tail end of the chart show these scenarios.

C - Annotation processor migration All the 25 PRs for this migration passed CI and were merged automatically, without any user intervention.

5.5.3 RQ3. Comparison with state-of-the-art code rewrite frameworks

5.5.3.A Performance

A - Experimental Setup We compare POLYGLOTPIRANHA-based stale feature flag cleanup against Piranha [119] and an equivalent we develop based upon Comby [23]. The Comby implementation has 29 rewrite rules for java. It was particularly easy to develop the Comby variant because POLYGLOTPIRANHA's concrete syntax DSL is inspired by Comby. For this evaluation, we chose 24 stale feature flags randomly from the PRs that 1. passed CI but were not accepted (at the time of writing this paper) 2. were used in java files (because Piranha only supports java). Note that we only chose 24 feature flags because it takes significant manual effort to integrate Piranha within our infrastructure due to Piranha depending on compilation⁴. For each feature flag, we noted the affected sub-targets and their sizes. We then applied the three tools across the sub-targets and the execution time was recorded. These experiments were performed on an enterprise-class VM in Google Cloud Platform. Note that we neither compare the quality of the cleanups nor precision because by construc-

⁴While Piranha was developed and was previously integrated at Uber, however, both Uber's feature flag API and developer infrastructure have changed since then

Table 5.2: Comparison of POLYGLOTPIRANHA against existing tools

| Tool | Metric | Bug Fixes | | | Feature Flag Cleanup | |
|-----------------|---------|-----------|-----------|---------------|----------------------|-------------|
| | | CWE-338 | slf4j | java.security | Android | iOS |
| POLYGLOTPIRANHA | LoC | 68 | 31 | 23 | 654 | 1156 |
| | # rules | 4 | 3 | 3 | 31 | 42 |
| Error-Prone | LoC | - | - | - | 3467 | - |
| SwiftSyntax | LoC | - | - | - | - | 1316 |
| OpenRewrite | LoC | 145 | 87 | 92 | - | - |
| Comby | # rules | - | - | - | 29 [†] | - |

[†] This feature flag cleanup variant was developed for the experiments.

tion Comby uses a more loose representation of code, based on Dyck-extended grammars [23] (i.e., balanced parenthesis grammars), whereas POLYGLOTPIRANHA uses language-specific grammars from the tree-sitter repertoire, hence POLYGLOTPIRANHA transformations are more powerful and precise. Conversely, ErrorProne and OpenRewrite can leverage semantic information like symbol/name resolution, for higher precision and applicability, but are not polyglot.

B – Results The line chart in Figure 5.12a shows the performance of each of the tools for the set of flags we identified above (ordered by the size of the corresponding *sub-targets*, ranging from 1.2K to 1.8M LoC). POLYGLOTPIRANHA took an average of 9.74 ± 3.46 seconds, Comby 121.67 ± 179.03 seconds (12.32 \times), and Piranha 413.91 ± 521.94 seconds (42.5 \times). We can see Piranha’s execution increases almost linearly with the target size (due to the fact that Piranha relies on building the target). POLYGLOTPIRANHA and Comby depended on the number of passes and files affected for the refactoring. The fact that POLYGLOTPIRANHA is faster than the Comby-based variant is surprising because Comby has a string based matching approach with minimal overhead. These results can be attributed to the fact that Comby has no sense of ordering between rules (nor scope), therefore, the match-replace rules are applied across the entire subtarget. POLYGLOTPIRANHA’s performance is also attributed to optimizations discussed in Section 5.4.1.C and 5.4.2.A.

Figure 5.12b shows the number of lines deleted by each of the tools for the same set of flags (in the same order). POLYGLOTPIRANHA deletes more lines of code because it’s able to delete trailing commas and comments. Note that we manually vetted that there are no over deletions in these PRs. In summary, POLYGLOTPIRANHA is consistently faster while deleting more lines than its imperative alternative Piranha and a lightweight Comby-based alternative.

5.5.3.B Expressiveness and Ease-of-use

A – Experimental Setup We compare the implementation of POLYGLOTPIRANHA-based tools against their imperative variants. Specifically, we compare the POLYGLOTPIRANHA-based Stale Feature flag cleanup pro-

gram against the implementation of Piranha [119]. Further we also encode three *pre-existing* code transformation recipes developed by *professional tool builders*, specifically OpenRewrite - 1. (JHipster Upgrade) Fix CWE-338 with SecureRandom [135] 2. (Slf4j) Loggers should be named for their enclosing classes [136] 3. (java.security) Use secure temporary file creation [137]. The selected patterns are 1. related to a popular java library 2. involve multiple interdependent changes 3. have associated test cases for validation 4. clearly fix a bug or security vulnerability.

B – Stale feature flag cleanup. As discussed in Section 5.5.3.A, Piranha is a stale feature flag cleanup tool with multiple implementations, one for each language supported. This is because Piranha is built upon language-specific imperative frameworks for code analysis and rewriting. To compare the expressiveness and conciseness of both approaches, we qualitatively and quantitatively compare the PiranhaJava and PiranhaSwift variants against their POLYGLOTPIRANHA-based counterparts.

PiranhaJava is built upon the ErrorProne [116] framework, whereas PiranhaSwift uses SwiftSyntax [130]. Table 5.2 (right) shows that POLYGLOTPIRANHA based approach is significantly more concise in terms of LoC. Moreover, rules can be re-used across languages. POLYGLOTPIRANHA-based Swift variant is more powerful than PiranhaSwift (e.g., supports variable inlining and cleanup of the unused members).

In contrast, our Comby implementation for feature flag cleanup in Java comprises 29 rules. Due to Comby's limitations, we were unable to express 10 transformations from our POLYGLOTPIRANHA implementation, including inlining singly-used boolean variables, deleting unused fields and variables, removing unnecessarily nested blocks, and deleting files under certain conditions and enum blocks. Despite this, the rule count difference is minor: 31 for POLYGLOTPIRANHA versus 29 for Comby. This is because POLYGLOTPIRANHA allows for the use of different, more powerful transformation languages. For instance, *tree-sitter* queries provide a syntax for complex alternations. Therefore, the Comby variant ends up being more verbose, requiring additional rules for the same task.

C – OpenRewrite OpenRewrite project is a semantic code search and transformation ecosystem. Its platform allows writing code transformation recipes for common framework migration and stylistic consistency tasks. We picked three relevant recipes written by professional developers, corresponding to high-impact transformations. We implemented the same refactoring actions using POLYGLOTPIRANHA. Table 5.2 shows the LoC count and number of rules for both POLYGLOTPIRANHA and OpenRewrite recipes. Our implementations pass the tests of the OpenRewrite recipes.

5.6 Discussion

Transformation Correctness. POLYGLOTPIRANHA does not guarantee that the transformed code will compile, be semantically correct, or precisely reflect the developer's intent. This limitation is common to other syntax-

driven code transformation tools such as [23, 63, 120]. While our dataflow analysis verifies the rule graph’s consistency and grammatical accuracy (Section 5.4.1.A), the effectiveness and accuracy of transformations ultimately rely on the quality of the rule graph itself.

Syntactic Limitations. POLYGLOTPIRANHA’s purely syntactic approach limits its ability to perform transformations that require semantic information of the code. In practice, this means that code rewrites that require type resolution, class hierarchy analysis, and/or control-flow analysis can not be expressed in the DSL today. Specifically, POLYGLOTPIRANHA: 1. *lacks precise def-use information.* We designed rules conservatively to identify def-use relationships within the syntactic scope of the variable declaration. However, due to the lack of SSA representation and dominator information POLYGLOTPIRANHA cannot reason about variable shadowing or re-initialization. 2. *lacks precise type information.* We approximate type information by analyzing declarations within a scope. This falls short when dealing with language features that obscure type information, such as Java’s `var` keyword or dynamically typed languages like Python. 3. *lacks call-graph analysis.* We approximate caller-callee relationships using method names and their number of arguments, resulting in imprecision in the presence of interfaces, class hierarchies, and method overloading. 4. *cannot handle advanced language features* that require semantic analysis, such as reflection.

Despite these limitations, our evaluation showcases that POLYGLOTPIRANHA is effective at automating three real-world code transformation tasks. Though imperfect, even in cases where it was partial, this automation substantially alleviated developers’ load as seen in Figure 5.11a.

Supporting New Languages. POLYGLOTPIRANHA supports languages beyond the ones listed in the evaluation, including Go, Python, Scala, TypeScript, as well as protocol formats like Thrift. POLYGLOTPIRANHA uses `tree-sitter` for code parsing, thus supporting a new language requires: 1. incorporating the `tree-sitter` grammar within POLYGLOTPIRANHA, and 2. authoring scope-capturing rules in a configuration file (i.e., one rule per scope such as class, method, or file). POLYGLOTPIRANHA uses these scopes when applying rules from the rule graph. Note that `tree-sitter` officially supports 133 programming languages [122], including functional languages like Haskell and Scheme. In fact, we support Scheme as a language in POLYGLOTPIRANHA, and use it within POLYGLOTPIRANHA’s implementation for rewriting its structural queries (a subset of Scheme). The implementation burden for this support was comparable to other languages.

Adapting POLYGLOTPIRANHA-based tools, like those for feature flag cleanup, to new languages may require additional work. For example, a rule for simplifying a disjunction (`true || :[a]`) in Java needs to be customized for Python as `true or :[a]`. However, we observed that some rules are reusable within a broad family of languages (Java, Kotlin, etc).

POLYGLOTPIRANHA’s Usability. To assist users in debugging and root-causing failures due to errors in the rule graph, POLYGLOTPIRANHA outputs detailed reports of all executed rules (in order) including their corresponding matched LoC ranges, and runtime arguments in an easily queryable format. This allows for step-by-step

replay and analysis. Our repository contains examples that explain how to enable debugging mode. We have also developed a playground for rule experimentation that allows users to easily experiment with rules and rules graphs on code snippets. This playground is publicly available on our artifact.

5.7 Key Takeaways and Contributions

In this chapter, we introduced a novel code transformation language, POLYGLOTPIRANHA. The language and toolset were designed to support complex code transformations. We demonstrated desirable properties of the language – namely, its expressiveness, usefulness, and run-time efficiency – through three case studies. By construction, this language is more expressive than *comby*, while maintaining its lightweight and declarative nature. Our goal was to enable complex code transformation to be expressed in a lightweight declarative language, rather than resorting to language-specific and imperative toolsets for code manipulation. In the next chapter, I introduce my proposed work on synthesizing complex API migration scripts in this language for automating logically related changes (like the ones depicted in Figure 5.1 and Figure 5.9).

6

Automating Logically Related Changes

Contents

| | |
|---|----|
| 6.1 Proposed Research Overview | 74 |
| 6.2 Generating Synthetic Pairwise Examples | 74 |
| 6.3 Generating Rules Graphs in the POLYGLOTPIRANHA language | 75 |
| 6.4 Evaluation Strategy | 76 |

In this chapter, I present my proposed work on automating logical code changes using migration graphs. My objective is to generate migration scripts that encapsulate complex code transformations, going beyond the straightforward API migration rules discussed in Chapter 4. These transformations typically involve multiple changes; for instance, in the migration example from Figure 5.9b, it is necessary to add a new field to the class to use as an argument. Similarly, in Figure 5.1, the migration of the logger object triggers several additional changes at various locations throughout the code.

I aim to synthesize these migration scripts and express them in the POLYGLOTPIRANHA language introduced in Chapter 5. In order to collect data to synthesize these scripts, I propose to develop a novel method to generate diverse and equivalent API usage examples across different libraries.

6.1 Proposed Research Overview

My proposed approach is divided in two steps:

1. **Generating Pairwise Training Examples:** As previously noted, automated API refactoring tools typically work by mining data from clients projects that have undergone migration between libraries. However, such examples are scarce. To overcome this, I propose an approach to automatically generate diverse migration pairwise examples that can be used for mining. In this approach, we seek to go beyond the simple example generation approach proposed introduced in Chapter 4. To do this, I plan to collect a set of existing API examples from GitHub (and simplify them if necessary) in the source library. Next, I will use a large language model pretrained on code from both the source and target libraries to generate and test migration attempts from the model.
2. **Synthesizing Rules Graphs in POLYGLOTPIRANHA:** I aim to synthesize migration scripts in the POLYGLOTPIRANHA language using the equivalent pairs generated from the previous step. The migration scripts will be synthesized using a set of predefined recipes or sketches.

6.2 Generating Synthetic Pairwise Examples

To synthesize rule graphs, we first need to collect a set of equivalent¹ API usages in the source and target libraries. In Chapter 4, we showed that LLMs are capable of synthesizing equivalent API usage examples in similar libraries. Our goal is to extend this approach.

We will use synthesis to generate and simplify API usage examples (e.g., [138]) from existing repositories on GitHub. The plan is to collect realistic API usage examples in the source library and minimize or remove irrelevant parts (similar to program slicing), and then convert these mined usages into minimal and functional examples. Since the code will be simplified, we may need to adapt the code snippets to ensure they remain functional, using either large language models or other synthesis techniques. This step allows us to mine simple API usage examples in the source library.

Next, we plan to use an LLM pre-trained on the source and target libraries, to attempt to migrate the API usage examples in the source to the target library. We will follow a generate and test approach. In many cases, we expect the model to fail to migrate the usage examples. Thus, the programs will be tested for equivalence before being taken as equivalent pairs of API usage. To compare the programs' behaviour, we will first extract comparable properties of the output of the APIs. For example, if the output of the programs are list objects, a property could be its length, or its contents. Using properties allows us to easily compare objects with different implementations and representations. The comparison functions will either have to be synthesized or manually crafted depending on the objects the libraries manipulate. Then, we will test whether

¹The term "equivalent" is used here loosely, as establishing true equivalence between programs is generally undecidable. We consider two programs to be equivalent if they exhibit similar behavior across a comprehensive test suite.

```

1 from Crypto.Cipher import AES as AES_PyCrypto
2 from Crypto.Random import get_random_bytes
3
4 cipher_pycrypto = AES_PyCrypto.new(aes_key, AES_PyCrypto.MODE_CTR, nonce=nonce)
5
6 ciphertext = cipher_pycrypto.encrypt(data)

```

(a) Encryption using PyCryptodome

```

1 from cryptography.hazmat.primitives.ciphers import Cipher, algorithms, modes
2 from cryptography.hazmat.backends import default_backend
3
4 cipher_cryptography = Cipher(algorithms.AES(aes_key), modes.CTR(nonce + (b'\x00' * 8)),
    backend=default_backend())
5
6 encryptor_cryptography = cipher_cryptography.encryptor()
7 ciphertext = encryptor_cryptography.update(data) + encryptor_cryptography.finalize()

```

(b) Encryption using Cryptography

Figure 6.1: Equivalent implementation of AES encryption using PyCryptodome and Cryptography libraries.

the pairs of programs are equivalent using a set of automatically generated inputs using either fuzzing [139] or property-based testing approaches [140].

Figure 6.1 shows an example of two equivalent API usage examples that GPT-4 generates in two different cryptography libraries. In this case, to test for equivalence, we simply need to check whether the output string (`ciphertext`) is the same for the same input data.

6.3 Generating Rules Graphs in the POLYGLOTPIRANHA language

Using the auto-generated examples, we will identify dependencies between individual APIs and understand which portions are related (e.g., by finding data-flow and control-flow matchings between the APIs in the source and the target [141]). Finding which APIs are related in the source and target may not be straightforward and require analyzing multiple API usages pairs simultaneously. For example, the instantiation of the encryptor object in Line 6 of Figure 6.1b directly results from invoking the encrypt API in Line 6 of Figure 6.1a, rather than the cipher creation.

Our approach to learning migration rules will be twofold:

1. **Simple One-to-One Mappings:** For straightforward mappings, we will create individual transformation rules using the approach described in Chapter 4 for MELT. For example, in Figure 6.1, the migration of `cipher_pycrypto` (Line 4, Figure 6.1a) from PyCryptodome to its equivalent in Cryptography (Line 5, Figure 6.1b) can be represented in a single comby rule. Note, however, that the API arguments differ

between both libraries; the target API (Line 5, Figure 6.1b) requires an `algorithms.AES` object instead of a direct `aes_key`, and also requires an additional argument, `backend=default_backend()`. This extra argument is created in place. If it had been necessary to instantiate a class field instead, it would not have been possible to represent this change in the comby language.

2. **Complex Multi-Step Migrations:** API migrations can be inherently related to others and result from cascading changes. Therefore, they are better expressed within an ordered rule graph. For instance, the migration of `cipher_pycrypto.encrypt` (Line 6, Figure 6.1a) is directly related to the migration of the `cipher_pycrypto` object. Moreover, this API call is replaced with two API calls (Lines 7-8, Figure 6.1b) in target libraries. Our approach will identify these complex migrations using data and control-flow dependency analysis to generate concise migration subgraphs from sketches or migration recipes. In this example, we would have an edge between the rule for object migration, and all the migration rules corresponding to associated API calls. This allows us to have precise control over the transformed code. The sketches for synthesizing the rules graphs will be manually crafted, incorporating specific structures or patterns for migration that we have identified from experience in writing migration graphs, as well as from previous research [127].

6.4 Evaluation Strategy

I will focus my efforts on libraries that manipulate objects that are in principle comparable across libraries. These include libraries for: XML manipulation (e.g., `xml.etree.ElementTree` [142] and `lxml` [143]), cryptography (e.g., `Pycryptodome` [144] and `Cryptography` [145]), compression, statistics, tensor manipulation, logging, and JSON manipulation. To evaluate the approach, I will answer the following research questions:

- RQ1. How accurate are the examples generated using our approach?** This question aims to evaluate the accuracy of the example generation approach. I will generate comprehensive test suites to test if the programs are equivalent, as well as perform manual validation on a sample of examples. The goal is to measure the percentage of programs that are equivalent.
- RQ2. How effective is our approach in generating migration graphs in practice?** I plan to assess this by synthesizing migration graphs and applying them on a test set (i.e., use the rule graph to migrate examples other than the ones used for learning). The goal is to test if our generated rule graphs can accurately migrate programs. I also plan to manually validate rule graphs.
- RQ3. How does our approach compare to existing state-of-the-art automated refactoring tools?** I will compare our approach with alternative tools in terms of (1) capability (i.e., which migration tasks can each approach handle); (2) accuracy of the refactorings (i.e., correctness of the produced code), and (3) run-time efficiency of the tools (i.e., time to execute the refactoring script). In particular, I intend to compare the approach against my previous work on SOAR (3) and MELT (4).

7

Conclusions

Refactoring is an important but laborious task. It ensures that code is better adaptable to future changes and is robust. However, developers often struggle with refactoring and organizing code, as most companies prioritize and incentivize feature development over refactoring. To help developers be more efficient and refactor faster, we propose a modern language and toolset for large-scale code transformation. To further facilitate the refactoring process, we also propose three different approaches to automated API refactoring. Our approaches advance the state of the art in multiple dimensions: first, they do not rely on or require training examples, which are scarce; second, the migrations are expressed in a simple, easy-to-read, and modifiable language that can be interpreted and changed by developers; third, they are fast, scalable, and can express complex, multi-step migrations.

Bibliography

- [1] C. Jaspan and J. Aldrich, "Checking framework interactions with relationships," in *ECOOP*, ser. Lecture Notes in Computer Science, vol. 5653, Springer, 2009, pp. 27–51.
- [2] C. R. De Souza, D. Redmiles, L.-T. Cheng, D. Millen, and J. Patterson, "How a good software practice thwarts collaboration: The multiple roles of apis in software development," *ACM SIGSOFT Software Engineering Notes*, vol. 29, no. 6, pp. 221–230, 2004.
- [3] C. R. de Souza and D. F. Redmiles, "On the roles of apis in the coordination of collaborative software development," *Computer Supported Cooperative Work (CSCW)*, vol. 18, no. 5-6, p. 445, 2009.
- [4] W. Maalej and M. P. Robillard, "Patterns of knowledge in api reference documentation," *IEEE Transactions on Software Engineering (TSE)*, vol. 39, no. 9, pp. 1264–1282, 2013.
- [5] A. Ketkar, D. Ramos, L. Clapp, R. Barik, and M. K. Ramanathan, "A Lightweight Polyglot Code Transformation Language," in *Proceedings of the ACM on Programming Languages*, ser. PLDI '24, ACM, 2024, doi: [10.1145/3656429](https://doi.org/10.1145/3656429).
- [6] V. Rajlich, "Software evolution and maintenance," in *Future of Software Engineering Proceedings*, 2014, pp. 133–144.
- [7] N. Chapin, J. E. Hale, K. M. Khan, J. F. Ramil, and W.-G. Tan, "Types of software evolution and software maintenance," *Journal of software maintenance and evolution: Research and Practice*, vol. 13, no. 1, pp. 3–30, 2001.
- [8] J. H. Perkins, "Automatically generating refactorings to support API evolution," in *PASTE*, ACM, 2005, pp. 111–114.
- [9] T. Mens and T. Tourwé, "A survey of software refactoring," *IEEE Transactions on Software Engineering (TSE)*, vol. 30, no. 2, pp. 126–139, 2004.
- [10] D. Silva, N. Tsantalis, and M. T. Valente, "Why we refactor? Confessions of Github contributors," in *Proc. ACM SIGSOFT Symposium on the Foundations of Software Engineering (FSE)*, 2016, pp. 858–870.

- [11] N. A. Ernst, S. Bellomo, I. Ozkaya, R. L. Nord, and I. Gorton, "Measure it? manage it? ignore it? software practitioners and technical debt," in *Proc. ACM SIGSOFT Symposium on the Foundations of Software Engineering (FSE)*, E. D. Nitto, M. Harman, and P. Heymans, Eds., ACM, 2015, pp. 50–60.
- [12] S. G. Eick, T. L. Graves, A. F. Karr, J. S. Marron, and A. Mockus, "Does code decay? assessing the evidence from change management data," *IEEE Transactions on Software Engineering (TSE)*, vol. 27, no. 1, pp. 1–12, 2001.
- [13] H. Krasner, *The cost of poor code quality in the us: A 2022 report*, Accessed: 28-09-2023, CISQ, 2022, [Online]. Available: <https://www.it-cisq.org/wp-content/uploads/sites/6/2022/11/CPSQ-Report-Nov-22-2.pdf>.
- [14] M. Kim, T. Zimmermann, and N. Nagappan, "A field study of refactoring challenges and benefits," in *Proc. ACM SIGSOFT Symposium on the Foundations of Software Engineering (FSE)*, 2012, pp. 1–11.
- [15] *Pytorch*, <https://pytorch.org/>, Accessed: 2023-11-20, 2023.
- [16] *Tensorflow*, <https://www.tensorflow.org/>, Accessed: 2023-11-20, 2023.
- [17] E. Duala-Ekoko and M. P. Robillard, "Asking and answering questions about unfamiliar apis: An exploratory study," in *2012 34th International Conference on Software Engineering (ICSE)*, IEEE, 2012, pp. 266–276.
- [18] X. Gao, A. Radhakrishna, G. Soares, R. Shariffdeen, S. Gulwani, and A. Roychoudhury, "Apifix: Output-oriented program synthesis for combating breaking changes in libraries," in *Proc. ACM SIGPLAN Object-Oriented Programming, Systems, Languages & Applications (OOPSLA)*, vol. 5, 2021, pp. 1–27.
- [19] M. Lamothe, W. Shang, and T. P. Chen, "A3: assisting android API migrations using code examples," *IEEE Transactions on Software Engineering (TSE)*, pp. 417–431, 2022.
- [20] M. Fazzini, Q. Xin, and A. Orso, "Automated api-usage update for android apps," in *Proc. of the International Symposium on Software Testing and Analysis*, D. Zhang and A. Møller, Eds., ACM, 2019, pp. 204–215.
- [21] S. Xu, Z. Dong, and N. Meng, "Meditor: Inference and application of API migration edits," in *Proc. of the International Conference on Program Comprehension*, Y. Guéhéneuc, F. Khomh, and F. Sarro, Eds., IEEE / ACM, 2019, pp. 335–346.
- [22] M. Fazzini, Q. Xin, and A. Orso, "Apimigrator: An api-usage migration tool for android apps," in *Proc. IEEE/ACM 7th International Conference on Mobile Software Engineering and Systems*, D. Lo, L. Mariani, and A. Mesbah, Eds., 2020, pp. 77–80.
- [23] R. van Tonder and C. Le Goues, "Lightweight multi-language syntax transformation with parser parser combinators," in *Proc. ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, 2019, pp. 363–378.

- [24] GitHub, *GitHub: Where the world builds software*, <https://github.com>, Accessed: 2023-11-22, 2023.
- [25] R. G. Kula, D. M. German, A. Ouni, T. Ishio, and K. Inoue, "Do developers update their library dependencies? an empirical study on the impact of security advisories on library migration," *Empirical Software Engineering*, vol. 23, pp. 384–417, 2018.
- [26] T. McDonnell, B. Ray, and M. Kim, "An empirical study of api stability and adoption in the android ecosystem," in *Proc. IEEE International Conference on Software Maintenance*, 2013, pp. 70–79.
- [27] M. Lamothe, Y.-G. Guéhéneuc, and W. Shang, "A systematic review of api evolution literature," *ACM Computing Surveys (CSUR)*, vol. 54, no. 8, pp. 1–36, 2021.
- [28] X. Gu, H. Zhang, D. Zhang, and S. Kim, "Deepam: Migrate apis with multi-modal sequence to sequence learning," *arXiv preprint arXiv:1704.07734*, 2017.
- [29] H. Zhong, L. Zhang, T. Xie, and H. Mei, "Inferring resource specifications from natural language API documentation," in *Proc. International Conference on Automated Software Engineering (ASE)*, IEEE Computer Society, 2009, pp. 307–318.
- [30] G. Uddin and M. P. Robillard, "How API documentation fails," *IEEE Software*, vol. 32, no. 4, pp. 68–75, 2015.
- [31] R. B. Watson, M. Stamnes, J. Jeannot-Schroeder, and J. H. Spyridakis, "API documentation and software community values: A survey of open-source API documentation," in *Proceedings of the 31st ACM international conference on Design of communication, Greenville, NC, USA, September 30 - October 1, 2013*, M. J. Albers and K. Gossett, Eds., ACM, 2013, pp. 165–174, [Online]. Available: <https://doi.org/10.1145/2507065.2507076>.
- [32] T. Lethbridge, J. Singer, and A. Forward, "How software engineers use documentation: The state of the practice," *IEEE Softw.*, vol. 20, no. 6, pp. 35–39, 2003, [Online]. Available: <https://doi.org/10.1109/MS.2003.1241364>.
- [33] W. Jin, D. Zhong, Z. Ding, M. Fan, and T. Liu, "Where to start: Studying type annotation practices in python," in *36th IEEE/ACM International Conference on Automated Software Engineering, ASE 2021, Melbourne, Australia, November 15-19, 2021*, IEEE, 2021, pp. 529–541, [Online]. Available: <https://doi.org/10.1109/ASE51524.2021.9678947>.
- [34] S. M. Sohan, C. Anslow, and F. Maurer, "A case study of web API evolution," in *2015 IEEE World Congress on Services, SERVICES 2015, New York City, NY, USA, June 27 - July 2, 2015*, L. Zhang and R. Bahsoon, Eds., IEEE Computer Society, 2015, pp. 245–252, [Online]. Available: <https://doi.org/10.1109/SERVICES.2015.43>.
- [35] G. Gousios, M. Pinzger, and A. van Deursen, "An exploratory study of the pull-based software development model," in *Proc. ACM IEEE International Conference on Software Engineering (ICSE)*, P. Jalote, L. C. Briand, and A. van der Hoek, Eds., ACM, 2014, pp. 345–355.

- [36] Y. Yu, H. Wang, V. Filkov, P. T. Devanbu, and B. Vasilescu, "Wait for it: Determinants of pull request evaluation latency on github," in *Proc. ACM IEEE International Conference on Mining Software Repositories (MSR)*, M. D. Penta, M. Pinzger, and R. Robbes, Eds., IEEE Computer Society, 2015, pp. 367–371.
- [37] M. D. Ernst, "Natural language is a programming language: Applying natural language processing to software development," in *2nd Summit on Advances in Programming Languages, SNAPL 2017, May 7-10, 2017, Asilomar, CA, USA*, B. S. Lerner, R. Bodík, and S. Krishnamurthi, Eds., ser. LIPIcs, vol. 71, Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2017, 4:1–4:14, [Online]. Available: <https://doi.org/10.4230/LIPIcs.SNAPL.2017.4>.
- [38] T. F. Bissyandé, D. Lo, L. Jiang, L. Réveillère, J. Klein, and Y. L. Traon, "Got issues? who cares about it? A large scale investigation of issue trackers from github," in *IEEE 24th International Symposium on Software Reliability Engineering, ISSRE 2013, Pasadena, CA, USA, November 4-7, 2013*, IEEE Computer Society, 2013, pp. 188–197, [Online]. Available: <https://doi.org/10.1109/ISSRE.2013.6698918>.
- [39] OpenAI, *GPT-4*, <https://openai.com/research/gpt-4/>, Accessed: May 2, 2023.
- [40] S. Gulwani, O. Polozov, and R. Singh, "Program synthesis," *Found. Trends Program. Lang.*, vol. 4, no. 1-2, pp. 1–119, 2017, [Online]. Available: <https://doi.org/10.1561/25000000010>.
- [41] K. L. Gwet, *Handbook of inter-rater reliability: The definitive guide to measuring the extent of agreement among raters*. Advanced Analytics, LLC, 2014.
- [42] A. Ni *et al.*, "SOAR: A synthesis approach for data science API refactoring," in *Proc. ACM IEEE International Conference on Software Engineering (ICSE)*, IEEE, 2021, pp. 112–124.
- [43] D. Ramos, H. Mitchell, I. Lynce, V. M. Manquinho, R. Martins, and C. L. Goues, "MELT: mining effective lightweight transformations from pull requests," in *38th IEEE/ACM International Conference on Automated Software Engineering, ASE 2023, Luxembourg, September 11-15, 2023*, IEEE, 2023, pp. 1516–1528, [Online]. Available: <https://doi.org/10.1109/ASE56229.2023.00117>.
- [44] M. Reddy, *API Design for C++*. Elsevier, 2011.
- [45] J. Bloch, *A brief history of the api*, InfoQ, Presentation at QCon San Francisco, 2018, [Online]. Available: <https://www.infoq.com/presentations/history-api/>.
- [46] J. Ofoeda, R. Boateng, and J. Effah, "Application programming interface (API) research: A review of the past to inform the future," *Int. J. Enterp. Inf. Syst.*, vol. 15, no. 3, pp. 76–95, 2019, [Online]. Available: <https://doi.org/10.4018/IJEIS.2019070105>.
- [47] J. J. Bloch, "How to design a good API and why it matters," in *Companion to the 21th Annual ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2006, October 22-26, 2006, Portland, Oregon, USA*, P. L. Tarr and W. R. Cook, Eds., ACM, 2006, pp. 506–507, [Online]. Available: <https://doi.org/10.1145/1176617.1176622>.

- [48] K. Cwalina, J. Barton, and B. Abrams, *Framework design guidelines: conventions, idioms, and patterns for reusable .net libraries*. Addison-Wesley Professional, 2020.
- [49] L. Murphy, M. B. Kery, O. Alliyu, A. Macvean, and B. A. Myers, "API designers in the field: Design practices and challenges for creating usable apis," in *2018 IEEE Symposium on Visual Languages and Human-Centric Computing, VL/HCC 2018, Lisbon, Portugal, October 1-4, 2018*, J. Cunha, J. P. Fernandes, C. Kelleher, G. Engels, and J. Mendes, Eds., IEEE Computer Society, 2018, pp. 249–258, [Online]. Available: <https://doi.org/10.1109/VLHCC.2018.8506523>.
- [50] D. Dig and R. Johnson, "How do apis evolve? a story of refactoring," *Journal of software maintenance and evolution: Research and Practice*, vol. 18, no. 2, pp. 83–107, 2006.
- [51] B. Ellis, J. Stylos, and B. A. Myers, "The factory pattern in API design: A usability evaluation," in *29th International Conference on Software Engineering (ICSE 2007), Minneapolis, MN, USA, May 20-26, 2007*, IEEE Computer Society, 2007, pp. 302–312, [Online]. Available: <https://doi.org/10.1109/ICSE.2007.85>.
- [52] V. Rajlich, "Software evolution and maintenance," in *Future of Software Engineering Proceedings*, 2014, pp. 133–144.
- [53] H. C. Benestad, B. Anda, and E. Arisholm, "Understanding cost drivers of software evolution: A quantitative and qualitative investigation of change effort in two evolving software systems," *Empirical Software Engineering*, vol. 15, pp. 166–203, 2010.
- [54] R. Minelli, A. Mocci, and M. Lanza, "I know what you did last summer: An investigation of how developers spend their time," in *Proceedings of the 2015 IEEE 23rd International Conference on Program Comprehension, ICPC 2015, Florence/Firenze, Italy, May 16-24, 2015*, A. D. Lucia, C. Bird, and R. Oliveto, Eds., IEEE Computer Society, 2015, pp. 25–35, [Online]. Available: <https://doi.org/10.1109/ICPC.2015.12>.
- [55] G. Suryanarayana, G. Samarthyam, and T. Sharma, *Refactoring for software design smells: managing technical debt*. Morgan Kaufmann, 2014.
- [56] M. Kim, T. Zimmermann, and N. Nagappan, "An empirical study of refactoring challenges and benefits at Microsoft," *IEEE Transactions on Software Engineering (TSE)*, vol. 40, no. 7, pp. 633–649, 2014.
- [57] H. Zhong, S. Thummalapenta, T. Xie, L. Zhang, and Q. Wang, "Mining API mapping for language migration," in *Proc. ACM IEEE International Conference on Software Engineering (ICSE)*, J. Kramer, J. Bishop, P. T. Devanbu, and S. Uchitel, Eds., ACM, 2010, pp. 195–204.
- [58] D. Dig and R. E. Johnson, "How do apis evolve? A story of refactoring," *J. Softw. Maintenance Res. Pract.*, vol. 18, no. 2, pp. 83–107, 2006, [Online]. Available: <https://doi.org/10.1002/smr.328>.
- [59] J. Lawall and G. Muller, "Coccinelle: 10 years of automated evolution in the linux kernel," in *2018 USENIX Annual Technical Conference (USENIX ATC 18)*, 2018, pp. 601–614.

- [60] L. Wasserman, "Scalable, example-based refactorings with refaster," in *Workshop on Refactoring Tools*, 2013, [Online]. Available: <http://dx.doi.org/10.1145/2541348.2541355>.
- [61] H. J. Kang, T. Ferdian, J. Lawall, G. Muller, L. Jiang, and D. Lo, "Semantic patches for java program transformation (artifact)," 2019.
- [62] I. Uber Technologies, *Go-patch: Structured code diffs and refactors*, <https://github.com/uber-go/gopatch>, GitHub repository, 2023, [Online]. Available: <https://github.com/uber-go/gopatch>.
- [63] aftandilian2012buildingOpen-source community, *Ast-grep: Write code to match code*, Accessed: 2023-09-22, 2023, [Online]. Available: <https://ast-grep.github.io>.
- [64] R. Rolim *et al.*, "Learning syntactic program transformations from examples," in *Proc. ACM IEEE International Conference on Software Engineering (ICSE)*, S. Uchitel, A. Orso, and M. P. Robillard, Eds., IEEE / ACM, 2017, pp. 404–415.
- [65] W. Ni, J. Sunshine, V. Le, S. Gulwani, and T. Barik, "Recode : A lightweight find-and-replace interaction in the IDE for transforming code by example," in *Proc. ACM Symposium on User Interface Software and Technology, Virtual Event, USA, October 10-14, 2021*, J. Nichols, R. Kumar, and M. Nebeling, Eds., ACM, 2021, pp. 258–269.
- [66] T. Petsios, A. Tang, S. Stolfo, A. D. Keromytis, and S. Jana, "Nezha: Efficient domain-independent differential testing," in *2017 IEEE Symposium on security and privacy (SP)*, 2017, pp. 615–632.
- [67] R. G. Kula, D. M. Germán, A. Ouni, T. Ishio, and K. Inoue, "Do developers update their library dependencies? - an empirical study on the impact of security advisories on library migration," *Springer Empirical Software Engineering (ESE)*, vol. 23, no. 1, pp. 384–417, 2018.
- [68] B. Dagenais and M. P. Robillard, "SemDiff: Analysis and recommendation support for API evolution," in *Proc. ACM IEEE International Conference on Software Engineering (ICSE)*, 2009, pp. 599–602.
- [69] G. Dotzler, M. Kamp, P. Kreutzer, and M. Philippsen, "More accurate recommendations for method-level changes," in *Proc. ACM SIGSOFT Symposium on the Foundations of Software Engineering (FSE)*, E. Bodden, W. Schäfer, A. van Deursen, and A. Zisman, Eds., ACM, 2017, pp. 798–808.
- [70] M. Asaduzzaman, C. K. Roy, S. Monir, and K. A. Schneider, "Exploring API method parameter recommendations," in *Proc. IEEE International Conference on Software Maintenance and Evolution (ICSME)*, R. Koschke, J. Krinke, and M. P. Robillard, Eds., 2015, pp. 271–280.
- [71] T. Molderez, R. Stevens, and C. D. Roover, "Mining change histories for unknown systematic edits," in *Proc. ACM IEEE International Conference on Mining Software Repositories (MSR)*, J. M. González-Barahona, A. Hindle, and L. Tan, Eds., 2017, pp. 248–256.

- [72] A. T. Nguyen, H. A. Nguyen, T. T. Nguyen, and T. N. Nguyen, "Statistical learning approach for mining API usage mappings for code migration," in *Proc. International Conference on Automated Software Engineering (ASE)*, I. Crnkovic, M. Chechik, and P. Grünbacher, Eds., ACM, 2014, pp. 457–468.
- [73] H. D. Phan, A. T. Nguyen, T. D. Nguyen, and T. N. Nguyen, "Statistical migration of API usages," in *Proc. ACM IEEE International Conference on Software Engineering (ICSE)*, S. Uchitel, A. Orso, and M. P. Robillard, Eds., 2017, pp. 47–50.
- [74] N. Meng, M. Kim, and K. S. McKinley, "Lase: Locating and applying systematic edits by learning from examples," ser. ICSE '13, San Francisco, CA, USA: IEEE Press, 2013, pp. 502–511, ISBN: 9781467330763.
- [75] A. Miltner *et al.*, "On the fly synthesis of edit suggestions," *Proc. ACM Program. Lang.*, vol. 3, no. OOPSLA, Oct. 2019, [Online]. Available: <https://doi.org/10.1145/3360569>.
- [76] Y. Zhang *et al.*, "Overwatch: Learning patterns in code edit sequences," *Proc. ACM Program. Lang.*, vol. 6, no. OOPSLA2, pp. 395–423, 2022, [Online]. Available: <https://doi.org/10.1145/3563302>.
- [77] J. Henkel and A. Diwan, "CatchUp!: Capturing and replaying refactorings to support API evolution," in *Proc. ACM IEEE International Conference on Software Engineering (ICSE)*, G. Roman, W. G. Griswold, and B. Nuseibeh, Eds., 2005, pp. 274–283.
- [78] *Pytorch conv2d api documentation*, <https://pytorch.org/docs/stable/generated/torch.nn.Conv2d.html>, Dec. 2023.
- [79] G. Salton and C. Buckley, "Term-weighting approaches in automatic text retrieval," *Information processing & management*, vol. 24, no. 5, pp. 513–523, 1988.
- [80] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *EMNLP, ACL*, 2014, pp. 1532–1543.
- [81] M. F. Porter, *Snowball: A language for stemming algorithms*, 2001.
- [82] C. S. Perone, R. Silveira, and T. S. Paula, "Evaluation of sentence embeddings in downstream and linguistic probing tasks," *arXiv preprint arXiv:1806.06259*, 2018.
- [83] S. Arora, Y. Liang, and T. Ma, "A simple but tough-to-beat baseline for sentence embeddings," in *ICLR (Poster)*, OpenReview.net, 2017.
- [84] *Common crawl*, <https://commoncrawl.org/>, Dec. 2023.
- [85] A. Solar-Lezama, "The sketching approach to program synthesis," in *APLAS*, ser. Lecture Notes in Computer Science, vol. 5904, Springer, 2009, pp. 4–13.
- [86] P. Orvalho, M. Terra-Neves, M. Ventura, R. Martins, and V. M. Manquinho, "Encodings for enumeration-based program synthesis," in *CP*, ser. Lecture Notes in Computer Science, vol. 11802, Springer, 2019, pp. 583–599.

- [87] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *NIPS*, 2013, pp. 3111–3119.
- [88] M. A. Hearst, "Automatic acquisition of hyponyms from large text corpora," in *COLING*, 1992, pp. 539–545.
- [89] *Tensorflow tutorial*, <https://www.tensorflow.org/tutorials>, Dec. 2023.
- [90] *Tensorflow applications*, <https://www.tensorflow.org/apidocs/python/tf/keras/applications>, Dec. 2023.
- [91] *Tensorflow models*, <https://github.com/tensorflow/models>, Dec. 2023.
- [92] *Kaggle*, <https://www.kaggle.com>, Dec. 2023.
- [93] Y. Feng, R. Martins, J. V. Geffen, I. Dillig, and S. Chaudhuri, "Component-based synthesis of table consolidation and transformation tasks from examples," in *PLDI*, ACM, 2017, pp. 422–436.
- [94] *Scrapy: A fast and powerful scraping and web crawling framework*, <https://scrapy.org/>, Dec. 2023.
- [95] L. M. de Moura and N. Bjørner, "Z3: an efficient SMT solver," in *TACAS*, ser. Lecture Notes in Computer Science, vol. 4963, Springer, 2008, pp. 337–340.
- [96] J. Reback, *DEPR: Deprecate DataFrame.append and Series.append*, <https://github.com/pandas-dev/pandas/pull/44539>, [Online; accessed May 02, 2023], 2022.
- [97] pandas-dev, *pandas-dev/pandas: Powerful data structures for data analysis, time series, and statistics*, <https://github.com/pandas-dev/pandas>, [Online; accessed May 02, 2023], 2022.
- [98] R. van Tonder. "Comby." (Jul. 2022), [Online]. Available: <https://comby.dev/docs/overview>.
- [99] *Scipy: Open source scientific tools for python*, <https://www.scipy.org/>, Accessed: May 2, 2023, 2023.
- [100] *Deprecate signal.spline in favor of signal*, <https://github.com/scipy/scipy/pull/14419>, Accessed: May 2, 2023.
- [101] A. Ketkar, O. Smirnov, N. Tsantalis, D. Dig, and T. Bryksin, "Inferring and applying type changes," in *Proc. ACM IEEE International Conference on Software Engineering (ICSE)*, ACM, 2022, pp. 1206–1218.
- [102] *DEPR: Squeeze() argument in read_csv/read_table #43242*, GitHub Pull Request Pandas #43242, Accessed on May 5, 2023 <https://github.com/pandas-dev/pandas/issues/43242>.
- [103] M. Dilhara, D. Dig, and A. Ketkar, "Pyevolve: Automating frequent code changes in python ml systems," in *Proc. ACM IEEE International Conference on Software Engineering (ICSE)*, 2023.
- [104] D. Halter, *Jedi: An autocompletion tool for python*, <https://jedi.readthedocs.io/en/latest/>, Accessed: May 2, 2023.
- [105] T. S. BV, *Tiobe index*, <https://www.tiobe.com/tiobe-index/>, Accessed on May 3, 2023, 2023.

- [106] R. van Tonder, *Comby with types*, <https://comby.dev/blog/2022/08/31/comby-with-types>, Accessed on 3 May 2023, Aug. 2022.
- [107] D. Ramos, *Replication of MELT: Mining Effective Lightweight Transformations from Pull Requests*, <https://doi.org/10.5281/zenodo.8226234>, 2023, DOI: 10.5281/zenodo.8226234.
- [108] J. Cohen, "A Coefficient of Agreement for Nominal Scales," *Educational and Psychological Measurement*, vol. 20, no. 1, pp. 37–46, 1960, issn: 0013-1644.
- [109] B. Cossette and R. J. Walker, "Seeking the ground truth: A retroactive study on the evolution and migration of software libraries," in *Proc. ACM SIGSOFT Symposium on the Foundations of Software Engineering (FSE)*, W. Tracz, M. P. Robillard, and T. Bultan, Eds., ACM, 2012, p. 55.
- [110] *Dep: Deprecate rollaxis*, GitHub Pull Request Numpy #9475, Accessed on May 4, 2023 <https://github.com/numpy/numpy/pull/9475>.
- [111] *DEPR: Pd.read_table*, GitHub Pull Request Pandas #21954, Accessed on May 4, 2023 <https://github.com/pandas-dev/pandas/pull/21954>.
- [112] *DEPR: Disallow int fill_value in shift with dt64/td64 #49362*, GitHub Pull Request Pandas #49362, Accessed on May 5, 2023 <https://github.com/pandas-dev/pandas/issues/49362>.
- [113] *Sourcegraph*, <https://sourcegraph.com/>, Accessed on May 5, 2023.
- [114] Claude, *100k context windows*, Anthropic, Accessed on July 24, 2023, <https://www.anthropic.com/index/100k-context-windows>, 2023.
- [115] T. Winters, T. Manshreck, and H. Wright, "Large-scale changes," in *Software engineering at google: Lessons learned from programming over time*. O'Reilly Media, 2020, ch. 22.
- [116] E. Aftandilian, R. Sauciu, S. Priya, and S. Krishnan, "Building useful program analysis tools using an extensible java compiler," Riva del Garda, Italy: IEEE Computer Society, 2012, pp. 14–23, [Online]. Available: <https://doi.org/10.1109/SCAM.2012.28>.
- [117] Clang. "Clang libtooling." Accessed: 2024-03-06. (2024).
- [118] J. Kim, D. Batory, and D. Dig, "Scripting parametric refactorings in java to retrofit design patterns," in *Proc. IEEE International Conference on Software Maintenance and Evolution (ICSME)*, IEEE, 2015, pp. 211–220.
- [119] M. K. Ramanathan, L. Clapp, R. Barik, and M. Sridharan, "Piranha: Reducing feature flag debt at uber," in *Proc. International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*, 2020, pp. 221–230.
- [120] M. Inc., *Fastmod*, Accessed: 2023-09-22, 2023, [Online]. Available: <https://github.com/facebookincubator/fastmod>.

- [121] A. Bosch and Contributors, *Detekt: A static code analyzer for Kotlin*, Accessed: 2023-09-21, 2023, [Online]. Available: <https://detekt.dev>.
- [122] GitHub, *Tree-sitter*, Accessed: 2023-09-24, 2023, [Online]. Available: <https://tree-sitter.github.io/tree-sitter/>.
- [123] J. E. F. Friedl, *Mastering Regular Expressions*, 3rd ed. Sebastopol, CA, USA: O'Reilly Media, 2006, ISBN: 978-0-596-52812-6.
- [124] A. V. Aho, R. Sethi, and J. D. Ullman, *Compilers: Principles, Techniques, and Tools* (Addison-Wesley series in computer science / World student series edition). Boston, MA, USA: Addison-Wesley, 1986, ISBN: 0-201-10088-6, [Online]. Available: <https://www.worldcat.org/oclc/12285707>.
- [125] M. Brunsfeld, *Tree-sitter: A new parsing system for programming tools*, Strange Loop Conference, Available online: <https://thestrangeloop.com/2018/tree-sitter---a-new-parsing-system-for-programming-tools.html>, St. Louis, MO: Strange Loop, Sep. 2018.
- [126] N. G. Fruja, "The correctness of the definite assignment analysis in c#," *J. Object Technol.*, vol. 3, no. 9, pp. 29–52, 2004, [Online]. Available: <https://doi.org/10.5381/jot.2004.3.9.a2>.
- [127] Moderne. "OpenRewrite: Semantic code search and transformation tool." Accessed: 2024-03-06, Moderne Inc. (2024), [Online]. Available: <https://docs.openrewrite.org>.
- [128] M. T. Rahman, L. Querel, P. C. Rigby, and B. Adams, "Feature toggles: Practitioner practices and a case study," in *Proc. ACM IEEE International Conference on Mining Software Repositories (MSR)*, M. Kim, R. Robbes, and C. Bird, Eds., Austin, TX, USA: ACM, 2016, pp. 201–211, [Online]. Available: <https://doi.org/10.1145/2901739.2901745>.
- [129] M. T. Rahman, P. C. Rigby, and E. Shihab, "The modular and feature toggle architectures of google chrome," *Springer Empirical Software Engineering (ESE)*, vol. 24, no. 2, pp. 826–853, 2019, [Online]. Available: <https://doi.org/10.1007/s10664-018-9639-0>.
- [130] Apple. "Swiftsyntax documentation." Accessed: 2024-03-06, Apple Inc. (2024), [Online]. Available: <https://swiftpackageindex.com/apple/swift-syntax/509.0.0/documentation/swiftsyntax>.
- [131] Meta, "Build faster with buck2: Our open source build system," Apr. 2024, Accessed: 2023-11-14, [Online]. Available: <https://engineering.fb.com/2023/04/06/open-source/buck2-open-source-large-scale-build-system/>.
- [132] Online. "ErrorProne ComplexBooleanConstant." Available: <https://errorprone.info/bugpattern/UnusedMethod>. (2024), (visited on 03/06/2024).
- [133] Online. "ErrorProne ComplexBooleanConstant." Available: <https://errorprone.info/bugpattern/ComplexBooleanConstant>. (2024), (visited on 03/06/2024).

- [134] S. Banerjee, L. Clapp, and M. Sridharan, "Nullaway: Practical type-based null safety for java," in *Proc. ACM SIGSOFT Symposium on the Foundations of Software Engineering (FSE)*, M. Dumas, D. Pfahl, S. Apel, and A. Russo, Eds., Tallinn, Estonia: ACM, 2019, pp. 740–750, [Online]. Available: <https://doi.org/10.1145/3338906.3338919>.
- [135] Online. "Fix cwe-338 with securerandom." Available: <https://docs.openrewrite.org/recipes/java/security/fixcwe338>. (2024), (visited on 03/06/2024).
- [136] Online. "Loggers should be named for their enclosing classes." Available: <https://docs.openrewrite.org/recipes/java/logging/slf4j/loggersnamedforenclosingclass>. (2024), (visited on 03/06/2024).
- [137] Online. "Use secure temporary file creation." Available: <https://docs.openrewrite.org/recipes/java/security/securetempfilecreation>. (2024), (visited on 03/06/2024).
- [138] R. P. L. Buse and W. Weimer, "Synthesizing API usage examples," in *34th International Conference on Software Engineering, ICSE 2012, June 2-9, 2012, Zurich, Switzerland*, M. Glinz, G. C. Murphy, and M. Pezzè, Eds., IEEE Computer Society, 2012, pp. 782–792, [Online]. Available: <https://doi.org/10.1109/ICSE.2012.6227140>.
- [139] A. Zeller, R. Gopinath, M. Böhme, G. Fraser, and C. Holler, *The fuzzing book*, 2019.
- [140] G. Fink and M. Bishop, "Property-based testing: A new approach to testing for assurance," *ACM SIGSOFT Software Engineering Notes*, vol. 22, no. 4, pp. 74–80, 1997.
- [141] H. A. Nguyen, T. N. Nguyen, D. Dig, S. Nguyen, H. Tran, and M. Hilton, "Graph-based mining of in-the-wild, fine-grained, semantic code change patterns," in *Proceedings of the 41st International Conference on Software Engineering, ICSE 2019, Montreal, QC, Canada, May 25-31, 2019*, J. M. Atlee, T. Bultan, and J. Whittle, Eds., IEEE / ACM, 2019, pp. 819–830, [Online]. Available: <https://doi.org/10.1109/ICSE.2019.00089>.
- [142] Python Software Foundation. "xml.etree.ElementTree - the elementtree xml api." (Accessed 2024), [Online]. Available: <https://docs.python.org/3/library/xml.etree.elementtree.html>.
- [143] lxml.de Contributors. "lxml - processing xml and html with python." (Accessed 2024), [Online]. Available: <https://lxml.de/>.
- [144] The PyCryptodome Team. "PyCryptodome documentation." (Accessed 2024), [Online]. Available: <https://pycryptodome.readthedocs.io/en/latest/>.
- [145] The Cryptography Developers. "Cryptography documentation." (Accessed 2024), [Online]. Available: <https://cryptography.io/en/latest/>.