

Bringing Extract Method Refactoring to the Rust Programming Language: Experiments with REM

Matthew Britton¹, Supervisor: Alex Potanin²

Abstract

This report details my current progress on advancing extract method refactoring in the Rust programming language. This project builds on existing work completed by Ilya Sergey, Sewen Thy and others in Adventure of a Lifetime: Extract Method Refactoring for Rust. The first half of the project is focussed on bringing the Rusty Extraction Maestro (REM) tool to a usable state, with the second half focussed on researching evaluating and implementing new extraction methods and techniques.

Keywords

Rust, Extract Method Refactoring

¹ School of Engineering, The Australian National University. — matt.britton@anu.edu.au

² School of Computing, The Australian National University. — alex.potanin@anu.edu.au

Contents

1	Introduction
2	Work Completed to Date
3	Coding Progress
4	Future Work
5	Issues Encountered So Far
	Acknowledgments
	Source Code
	References
	Appendices
A1	Non Local Controlflow Repair [1]
A2	Ownership and Borrowing Repair [1]
A3	Lifetime Repair [1]

1. Introduction

Context

This project starts an extension of work done in a previous project by Sewen Thy and Ilya Sergey, in which they theorised and then implemented the beginnings of an automated “Extract Method” refactoring tool for the Rust programming language. Extract Method refactorings are a well known and widely used refactoring approach, however, due to Rust’s unique approach to memory safety of ownership and borrowing, and the lifetime system that enforces these rules, it is far more difficult to implement than in other languages. Ilya Sergey et al’s contribution focused on the theoretical aspects of the problem [1], whilst Sewen Thy’s contribution focused on the implementation, [2], with a pure Rust algorithm and a heavy reliance on IntelliJ IDEA’s type inference and static

code analysis tools. Their approach was dubbed “Rusty Extraction Maestro” (REM). Figure 1 gives a high-level overview of the REM algorithm.

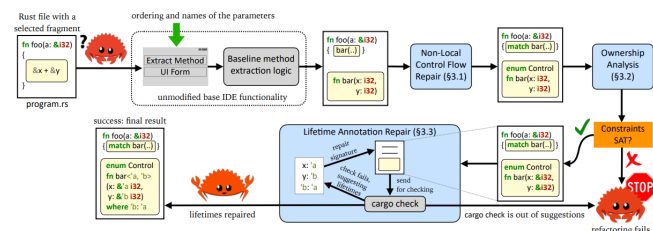


Figure 1. An overview of the REM Algorithm, from [1]

The REM algorithm is a three stage process, beginning after IntelliJ IDEA has performed an initial code extraction and type inference. If REM fails (i.e. an extraction isn’t possible), then the code is returned to its original condition. The stages are:

1. **Controller:** handles any non local control flow. This includes `return`, `break` and `continue`. As rust is an expressions based language [3] these statements can be used in arbitrary combinations, adding substantial complexity. The algorithm is included appendix A1.
2. **Borrower:** infers any ownership annotations. In Rust, variables can be either immutable (`imm`) or mutable (`mut`). Additionally, they are either owned by the current scope, or borrowed from another scope.

```
1 let x = 5; //x is an imm variable
2 let x_ref = &x; //imm reference to x
3 let mut y = 5; //z is a mut variable
4 let ref_y = &mut z; //mut reference to z
```

```

5 //Functions can take references
6 fn foo(s: &String) {println!("{}", s);}

```

The borrower is at its core, a dataflow analysis algorithm that infers the ownership of each variable, and adds the correct information to the function body and scope. The algorithms are included in appendix A2.

3. **Repairer:** focuses on correcting lifetime annotations in Rust functions, particularly those involving references. When a function has borrowed values, lifetimes need to be explicitly managed to ensure that the references remain valid. The repairer first explicitly annotates all lifetimes, using the rust compiler (rustc) to check for errors: 'a represents the lifetime of x ¹.

```

fn bar <'a 'b> (x: &'a i32, y: &'b i32)
    -> &'b i32 where 'a: 'b {...}

```

It then applies the lifetime elision rules [4] to remove any lifetimes the compiler can infer to bring the code back to a human readable state. The algorithm is included in appendix A3

One of the most difficult aspects of any refactoring / optimisation tool is ensuring the correctness of the output. A refactor is “correct” if the original behaviour of the program is preserved [5]. Unfortunately, for the REM algorithm, formal verification of the correctness is impossible as it would a model of the formal semantics of (safe) Rust, which doesn’t exist ². Instead, the toolchain relies on a combination of extensive testing, the rust compiler, and a fail-early approach to ensure that the user is never presented with “incorrect” code.

Aims

This project aims to improve and extend on the REM tool by addressing key limitations and researching (and subsequently implementing) new approaches to refactoring.

1. **Fixing REM’s Issues:** Remove the dependency on (i) rustc and (ii) IntelliJ IDEA to make REM a standalone tool. Additionally, reduce file I/O overhead by optimising the tools internal data handling
2. **Platform-Agnostic Code Extraction:** Build a complete, end-to-end solution for method extraction using *Rust Analyzer* (RA) for the code extraction, wrapped in a single, user-friendly CLI tool written entirely in Rust.
3. **VSCode Extension Proof of Concept:** Develop a VS-Code extension to showcase how REM can be easily integrated into editors, with potential for extending to other environments.

¹See Rust Book - Lifetimes for more information

²Formal correctness proofs for refactorings are very rare, even for languages with fully formalised semantics [1]

4. **Exploring New Refactorings:** The big ticket item for this research project where the bulk of my time will be spent, investigating extract method refactoring for the far more advanced areas of *Generics*, *Asynchronous Code*, *Macros*, *Closures*, and *Const Functions*. *Generics* and *Asynchronous Code* present special challenges due to (a) the complex lifetimes and (b) the uncertainty of the code at compile time, and are what the majority of the research will focus on.

5. **Final Aim: Merging REM with RA:** An ultimate, ambitious goal for this project is to merge the REM toolchain with RA, Rust’s official language server protocol (LSP) implementation ³. This would require proving substantially more of the correctness of REM (although as previously noted, a full correctness proof is impossible), allowing its functionality to be available across any Rust-supported editor. If that is not feasible, then a well featured VSCode extension is a suitable alternative.

2. Work Completed to Date

Literature Review Progress

Research into Advanced Refactoring Methods

3. Coding Progress

Modifications to REM

The most important modification to the original toolchain has been decoupling it from the rust compiler. Previously, REM was heavily tied to rustc, the rust compiler. It relied compiler specific representations of files and their syntax trees. Because of this dependency, REM required a very specific buildchain, and could not be used outside of this build context. Additionally, when I took over the project, changes to external libraries meant that the toolchain was completely non-functional. Some work is still required to bring the toolchain off of nightly rust onto stable rust.

Another small modification has been changing the way REM handles data. Previously, the three main components of REM (controller, borrower and repairer) worked by each reading and writing to a file, and then passing the file to the next stage of the toolchain. By changing the way data is passed between the (now 5) modules of REM, a roughly 10-20% speedup has been achieved.

REM-Extract: Performing initial function extraction in Rust

REM-CLI: A comprehensive command line interface for REM

The REM-CLI provides a unified interface for developers to interact with the 5 separate modules of REM. It streamlines the process of performing method extraction by offering an intuitive command-line interface that handles everything from

³<https://rust-analyzer.github.io/>

code analysis to extraction. By consolidating the toolchain into a single, easy-to-use CLI, REM removes the need for users to manually manage multiple stages of the extraction process. It is also needed to integrate the functionality into a code editor.

REM-VSCode: A Visual Studio Code extension for REM

The REM-VSCode extension is a proof of concept that demonstrates how the REM tool can be easily integrated into any code editor. It provides a user-friendly interface within the far more popular Visual Studio Code editor. Because the entire backend of REM is now written in Rust, the extension is able to be independent of a specific platform, and thanks to the previous work on decoupling the tool from rustc, it no longer requires a very specific environment / buildchain to run. It is available on the VSCode Marketplace.

```
Test 99 | PASSED | FAILED | PASSED: return_to_parent in 1.82s
=====
Differences or compilation errors found for test 'return_to_parent':
Differences found between expected and output:
  fn foo() -> i64 {
    let n = 1;
    - let k = match fun_name(n) {
    + let k = match fun_name(&n) {
      Ok(value) => value,
      Err(value) => return value,
    };
    (n + k) as i64
  }

- fn fun_name(n: i32) -> Result<i32, i64> {
-   let m = n + 1;
+ fn fun_name(n: &i32) -> Result<i32, i64> {
+   let m = *n + 1;
+   return Err(1);
+   let k = 2;
+   Ok(k)
+ }

fn main() {
}
```

Figure 2. An example of the non-deterministic nature of the extraction algorithm. Red is expected output, green is the returned code

Come up with something else to go here

4. Future Work

5. Issues Encountered So Far

Non-Deterministic Nature of the Extraction Algorithm

The extraction algorithm used by RA is non-deterministic, meaning that for the same input code, the output of the extraction process can vary. This results in different versions of semantically equivalent but syntactically distinct code. In other words, while the extracted code maintains the same functionality, its structure and readability can change between runs. This poses a significant challenge for testing the tool, as each run may yield different outputs for the same input, making consistent verification difficult.

For example in Figure 2, a test of function extraction, RA may choose to pass variables by reference to the extracted function, dereferencing them inside the function body. In another instance, it might pass the same variables by ownership, altering the code's appearance but not its behavior. While both versions of the code will likely compile to identical machine code, the readability and clarity of the resulting code can vary.

Acknowledgments

I would like to thank my supervisor, Alex Potanin, for his guidance and support. Additionally, I would like to thank Sasha Pak for listening in on all of our meetings and providing valuable feedback. Finally, I would like to thank Sewen Thy, the original developer of REM, for his help in understanding the project and his assistance in getting me started.

Source Code

The source code for this project can be found at the following repositories:

- REM-cli
- REM-extract — (Crates.io)
- REM-vscode
- REM-utils (Crates.io)
- REM-controller — (Crates.io)
- REM-borrower — (Crates.io)
- REM-repairer — (Crates.io)
- REM-constraint — (Crates.io)

References

- [1] Sewen Thy, Andreea Costea, Kiran Gopinathan, and Ilya Sergey. Adventure of a lifetime: Extract method refactoring for rust. *Proc. ACM Program. Lang.*, 7(OOPSLA2), October 2023.
- [2] Sewen Thy. Borrowing without sorrowing: Implementing extract method refactoring for rust. 2023.
- [3] Aaron Matsakis, Nicholas Turon. The rust rfc book: Statements and expressions. 2020.
- [4] Aaron Matsakis, Nicholas Turon. The rust rfc book: Lifetime elision. 2020.
- [5] Sorin Lerner, Todd Millstein, and Craig Chambers. Automatically proving the correctness of compiler optimizations. *SIGPLAN Not.*, 38(5):220–231, May 2003.

Appendices

A1. Non Local Controlflow Repair [1]

Algorithm 1: FIXNONLOCALCONTROL

Input : an extracted function EF , an introduced function call expression E (i.e., $EF(\dots)$) in the caller
Output: a list of patches PS to apply to the refactored file

```

1  $PS \leftarrow []$ 
2  $R \leftarrow$  collect return statements in  $EF$ 
3  $B, C \leftarrow$  collect top-level break and continue statements in  $EF$ 
4 if  $R \cup B \cup C \neq \emptyset$  then
5    $RTY \leftarrow \text{BUILDRETURNType}(R, B, C)$ 
6    $PS \leftarrow \text{UPDATEReturnType}(EF, RTY) :: PS$ 
7   for  $l_r \in R$  do  $PS \leftarrow (l_r, \text{return } e \rightsquigarrow \text{return Ret}(e)) :: PS$ 
8   for  $l_b \in B$  do  $PS \leftarrow (l_b, \text{break} \rightsquigarrow \text{return Break}) :: PS$ 
9   for  $l_c \in C$  do  $PS \leftarrow (l_c, \text{continue} \rightsquigarrow \text{return Continue}) :: PS$ 
10   $l_E \leftarrow$  find location of the final expression of  $EF$ 
11   $PS \leftarrow (l_E, E \rightsquigarrow \text{Ok}(E)) :: PS$ 
12   $\overline{CS} \leftarrow \text{BUILDCASESFORReturnType}(RTY)$ 
13   $l_{\text{caller}} \leftarrow$  location of  $E$ 
14   $PS \leftarrow (l_{\text{caller}}, E \rightsquigarrow \text{match } E \text{ with } \overline{CS}) :: PS$ 
15 return  $PS$ 

```

A2. Ownership and Borrowing Repair [1]

Algorithm 2: FIXOWNERSHIPANDBORROWING

Input : the extracted function EF , the expression E of the call to EF , original function F
Output: a set of patches PS

```

1  $Aliases \leftarrow$  alias analysis on  $F$  /* maps variables to their aliases */
2  $Mut \leftarrow$  COLLECTMUTABILITYCONSTRAINTS( $EF, Aliases$ )
3  $Own \leftarrow$  COLLECTOWNERSHIPCONSTRAINTS( $EF, Aliases, F$ )
4  $PS \leftarrow []$ 
5 for  $param \in EF.params$  do /* derive patches for the signature of  $EF$  */
6    $v, \tau, l \leftarrow param.var, param.type, param.loc$ 
7   if UNSAT( $Mut \cup Own, v$ ) then raise RefactorError
8   if LUB( $Mut \cup Own, v$ ) =  $\langle mut, ref \rangle$  then  $PS \leftarrow (l, v : \tau \rightsquigarrow v : \&mut \tau) :: PS$ 
9   if LUB( $Mut \cup Own, v$ ) =  $\langle imm, ref \rangle$  then  $PS \leftarrow (l, v : \tau \rightsquigarrow v : \&\tau) :: PS$ 
10 for  $param \in EF.params$  do /* derive the patches for the body of  $EF$  */
11   if LUB( $Mut \cup Own, param.var$ ) =  $\langle \_, ref \rangle$  then
12      $Exps \leftarrow$  collect from  $EF.body$  all the occurrences of  $param.var$ 
13     for  $e \in Exps$  do  $PS \leftarrow (e.loc, e \rightsquigarrow (* e)) :: PS$ 
14 for  $arg \in E.args$  do /* derive patches for the call to  $EF$  */
15    $v, e, l \leftarrow arg.var, arg.exp, arg.loc$ 
16   if LUB( $Mut \cup Own, v$ ) =  $\langle mut, ref \rangle$  then  $PS \leftarrow (l, e \rightsquigarrow \&mut e) :: PS$ 
17   if LUB( $Mut \cup Own, v$ ) =  $\langle imm, ref \rangle$  then  $PS \leftarrow (l, e \rightsquigarrow \&e) :: PS$ 

```

Algorithm 3: COLLECTMUTABILITYCONSTRAINTS

Input : extracted function EF , an alias map $Aliases$
Output: a set Mut of mutability constraints

```

1  $MV \leftarrow$  collect all the variables in  $EF$  that are part of an  $lvalue$  expression
2  $MV \leftarrow$  add to  $MV$  all the variables in the body of  $EF$  that are function call arguments with mutable requirements
3  $MV \leftarrow$  add to  $MV$  all the variables in  $EF$  that are mutably borrowed
4  $Mut \leftarrow \{ imm <: p.var \mid p \in EF.params \wedge \forall v' \in Aliases(p.var) : v' \notin MV \} \cup$ 
5    $\{ mut <: p.var \mid p \in EF.params \wedge \exists v' \in Aliases(p.var) : v' \in MV \}$ 

```

Algorithm 4: COLLECTOWNERSHIPCONSTRAINTS

Input : extracted function EF , an alias map $Aliases$, original caller function F
Output: a set $Ownership$ of ownership constraints

```

1  $FV \leftarrow$  free variables in  $F$  in the code snippet after the call to  $EF$ 
2  $PBVV \leftarrow$  collect all vars in  $EF.params$  declared as pass-by-value
3  $Borrows \leftarrow PBVV \cap \{ p.var \mid p \in EF.params \wedge \exists v' \in Aliases(p.var) : v' \in FV \}$ 
4  $Own \leftarrow$  collect all the vars in  $EF$  which are moved into or out of
5  $Ownership \leftarrow \{ v <: ref \mid v \in Borrows \} \cup \{ own <: v \mid v \in Own \}$ 

```

A3. Lifetime Repair [1]

Algorithm 5: FIXLIFETIMES

Input : a cargo manifest file *CARGO_MANIFEST* for the whole project, extracted function *EF*

Output: patched extracted function *EF'*

```

1 EF' ← clone EF
2 EF' ← update EF' by annotating each borrow in EF'.params and EF'.ret with a fresh lifetime where none exists
3 EF' ← update EF' by adding the freshly introduced lifetimes to the list of lifetime parameters in EF'.sig
4 Loop
5   err ← (cargo check CARGO_MANIFEST).errors
6   if err =  $\emptyset$  then break                                     /* refactoring is completed */
7   suggestions ← collect lifetime bounds suggestions from err
8   if suggestions =  $\emptyset$  then raise RefactorError               /* refactoring failed */
9   EF' ← apply suggestions to EF'
  // readability optimisations:
10 EF' ← collapse the cycles in the where clause of EF'.sig
11 EF' ← apply elision rules

```
