

Code Similarity Detection Tool

Design Documentation

Kevin Cao

April 3, 2021

Contents

1	Introduction	1
1.1	Purpose	1
1.2	Scope	1
2	Conversion of Source Code	2
2.1	Programming Language Scope	2
2.2	Abstract Syntax Trees	3

1 Introduction

With the rising popularity of cryptocurrencies, it has become crucial that crypto-transactions are not only trustworthy, but also efficient and reliable. While Bitcoin continues to control a majority of the crypto-market, alternatives like Ethereum have also been growing in popularity.

For some of these alternatives, transactions are controlled and automated using “smart contracts”. These smart contracts are computer programs that explicitly detail the rules and regulations to an agreement between two or more parties. Once the terms of the agreement have been met, the contract will self-execute and complete the transaction as agreed upon.

1.1 Purpose

For many transactions, the corresponding smart contracts are derived from pre-existing popular and mature contracts. By detecting which parent contracts the newer contracts are derived from, the validation and verification process of the contracts can be expedited.

As such, a Code Similarity Detection Tool would serve the purpose of assisting in the process of determining the original source of a new contract.

1.2 Scope

This Software Design Document details a basic system that serves as a proof of concept for a Code Similarity Detection Tool. For the sake of simplicity, the tool will compare the similarity between code written in a simplified subset of the lisp-based programming language Racket. The system first translates the source code into an Abstract Syntax Tree (AST), and then converts the AST into De Bruijn notation before running the similarity algorithm. Lisp-based languages are best suited for this type of translation, but other languages can also be similarly converted, although the process is a bit more involved.

2 Conversion of Source Code

2.1 Programming Language Scope

As mentioned in the previous section, the tool will act on a simplified subset of the Racket programming language, which will be referred to as R_I . As such, it is best that this subset is defined explicitly.

The most basic component in R_I is an *atom*. Atoms include the basic primitive data values such as integers, booleans, and symbols. It also includes variables and empty lists.

However, most programs in R_I will not simply be made up of atoms — R_I is largely comprised of *exprs*, or expressions. An expression can be a simple atoms, an if-statement, a let-statement, or a function application. Figure 1 shows some examples of expressions.

Of course, in a functional programming language like R_I , it is imperative that function definitions are explicated. Functional definitions are in a class of their own, *defs*. Definitions detail function construction by first taking a variable name to describe the function, and then variable parameters, which are used in the expression body of the definition. Expressions can use the function name from the definition in functional applications, which are categorized under expressions.

Finally, an R_I program can be described as a list of function definitions followed by an expression statement.

```

atom ::= int | bool | sym | '()' | var
exp  ::= atom | (if exp exp exp) | (let ([var exp]) exp)
      | (+ exp exp) | (- exp exp) | (. exp exp) | (- exp)
      | (cons exp exp) | (eq? exp exp) | (and exp exp) | (or exp exp) | (not exp)
      | ...
def  ::= (define (var var ...) exp)
R1 ::= def ... exp

```

Figure 1: The concrete syntax of R_1 . “...” signifies a list of the previous token.

Compared to Racket, R_1 is a fairly barebones language—however, R_1 will be sufficient in demonstrating the core mechanic of the system.

2.2 Abstract Syntax Trees

Before we can begin the similarity algorithm, we first need to convert the source code into an Abstract Syntax Tree, which abstracts away the source code into the node objects of the tree. In short, we convert the concrete syntax detailed in Figure 1 into the abstract syntax detailed in Figure 2.

```

atom ::= (Int int) | (Bool bool) | (Sym sym) | (Var var)
op   ::= (Var var) | + | - | · | cons | eq? | and | or | not | ...
exp  ::= atom | (If exp exp exp) | (Let (Var var) exp exp) | (Apply op exp ...)
def  ::= (Def var (var ...) exp)
R1  ::= (Program (def ...) exp)

```

Figure 2: The abstract syntax of R_1 .

In Figure 2, the components of R_1 are converted into objects. You may also notice there is a new classification—*op*. Here, we have abstracted away functions to be either custom functions created in the program [(Var var)], or inbuilt functions in the language itself [+ , cons, etc.]. By doing so, we are able to abstract away all function applications to the (Apply op exp ...) object.

The rest of the Figure 2 is fairly self-explanatory, so we will not go into further detail here.