

Towards certified virtual machine-based regular expression parsing

Thales Antônio Delfino
Universidade Federal de Ouro Preto
Ouro Preto, Minas Gerais, Brazil

Rodrigo Ribeiro
Universidade Federal de Ouro Preto
Ouro Preto, Minas Gerais, Brazil

Abstract

Regular expressions (REs) are pervasive in computing. We use REs in text editors, string search tools (like GNU-Grep) and lexical analysers generators. Most of these tools rely on converting regular expressions to its corresponding finite state machine or use REs derivatives for directly parse an input string. In this work, we want to investigate the suitability of another approach: instead of using derivatives or generate a finite state machine for a given RE, we will develop a virtual machine (VM) for parsing regular languages, in such a way that a RE is merely a program executed by the VM over the input string. We developed a prototype implementation in Haskell, test it using QuickCheck and provide proof sketches of its correctness with respect to RE standard inductive semantics.

CCS Concepts • **Theory of computation** → **Regular languages; Operational semantics;**

Keywords Regular Expressions, Parsing, Virtual Machines, Operational semantics

ACM Reference format:

Thales Antônio Delfino and Rodrigo Ribeiro. 2018. Towards certified virtual machine-based regular expression parsing. In *Proceedings of XXII BRAZILIAN SYMPOSIUM ON PROGRAMMING LANGUAGES, São Carlos, September 17–21, 2018 (SBLP2018)*, 8 pages. DOI: 10.475/123_4

1 Introduction

We name parsing the process of analyzing if a sequence of symbols matches a given set of rules. Such rules are usually specified in a formal notation, like a grammar. If a string can be obtained from those rules, we have success: we can build some evidence that the input is in the language described by the underlying formalism. Otherwise, we have a failure: no such evidence exists.

In this work, we focus on the parsing problem for regular expressions (REs), which are an algebraic and compact way of defining regular languages (RLs), i.e. languages that can be recognized by (non-)deterministic finite automata and equivalent formalisms. REs are widely used in string search tools, lexical analyser generators and XML schema languages [13]. Since RE parsing is pervasive in computing, its correctness is crucial and is the subject of study of several recent research works (e.g [3, 10, 20, 25]).

Approaches for RE parsing can use representations of finite state machines (e.g. [10]), derivatives (e.g. [20, 25]) or the so-called pointed RE's or its variants [3, 11]. Another approach for parsing is

based on the so-called parsing machines, which dates back to 70's with Knuth's work on top-down syntax analysis for context-free languages [17]. Recently, some works tried to revive the use of such machines for parsing: Cox [6] defined a VM for which a RE can be seen as “high-level programs” that can be compiled to a sequence of such VM instructions and Lua library LPEG [16] defines a VM whose instruction set can be used to compile Parser Expressions Grammars (PEGs) [12]. Such renewed research interest is motivated by the fact that is possible to include new features by just adding and implementing new machine instructions.

Since LPEG VM is designed with PEGs in mind, it is not appropriate for RE parsing, since the “star” operator for PEGs has a greedy semantics which differs from the conventional RE semantics for this operator. Also, Cox's work on VM-based RE parsing has problems. First, it is poorly specified: both the VM semantics and the RE compilation process are described only informally and no correctness guarantees is even mentioned. Second, it does not provide an evidence for matching, which could be used to characterize a disambiguation strategy, like Greedy [13] and POSIX [26]. To the best of our knowledge, no previous work has formally defined a VM for RE parsing that produces evidence (parse trees) for successful matches. The objective of this work is to give a first step in filling this gap. More specifically, we are interested in formally specify, implement and test the correctness of a VM based small-step semantics for RE parsing which produces bit-codes as a memory efficient representation of parse-trees. As pointed by [23], bit-codes are useful because they are not only smaller than the parse tree, but also smaller than the string being parsed and they can be combined with methods for text compression. We leave the task of proving that our VM follows a specific disambiguation strategy to future work.

Our contributions are:

- We present a small-step semantics for RE inspired by Thompson's NFA¹ construction [27]. The main novelty of this presentation is the use of data-type derivatives, a well-known concept in functional programming community, to represent the context in which the current RE being evaluated occur. We show informal proofs² that our semantics is sound and complete with respect to RE inductive semantics.
- We describe a prototype implementation of our semantics in Haskell and use QuickCheck [5] to test our semantics against a simple implementation of RE parsing, presented in [11], which we conjecture that is correct. Our test cases cover both accepted and rejected strings for randomly generated REs.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

SBLP2018, São Carlos

© 2018 Copyright held by the owner/author(s). 123-4567-24-567/08/06...\$15.00
DOI: 10.475/123_4

¹Non-deterministic finite automata.

²By “informal proofs” we mean proofs that are not mechanized in a proof-assistant. Due to space reasons, proofs of the relevant theorems are omitted from this version. Detailed proofs can be found in the accompanying technical report available on-line [7].

- We show how our proposed semantics can produce bit codes that denote parse trees [23] and test that such generated codes correspond to valid parsing evidence using QuickCheck.

We are aware that using automated testing is not sufficient to ensure correctness, but it can expose bugs before using more formal approaches, like formalizing our algorithm in a proof assistant. Such semantic prototyping step is crucial since it can avoid proof attempts that are doomed to fail due to incorrect definitions. The project's on-line repository [7] contains the partial Coq formalization of our semantics. Currently, we have formalized the semantics and its interpreter function. The Coq proof that the proposed small-step semantics is equivalent to the usual inductive RE semantics is under development.

The rest of this paper is organized as follows. Section 2 presents some background concepts on RE and data type derivatives that will be used in our semantics. Our operational semantics for RE parsing and its theoretical properties are described in Section 3. Our prototype implementation and the QuickCheck test suit used to validate it are presented in Section 4. Section 5 discuss related work and Section 6 concludes.

We assume that the reader knows the Haskell programming language, specially the list monad and how it can be used to model non-determinism. Good introductions to Haskell are available elsewhere [18]. All source code produced, including the literate Haskell source of this article (which can be preprocessed using `lhs2TeX` [19]), instructions on how to build it and reproduce the developed test suit are available on-line [7].

2 Background

2.1 Regular expressions: syntax and semantics

REs are defined with respect to a given alphabet. Formally, the following context-free grammar defines RE syntax:

$$e ::= \emptyset \mid \epsilon \mid a \mid ee \mid e + e \mid e^*$$

Meta-variable e will denote an arbitrary RE and a an arbitrary alphabet symbol. As usual, all meta-variables can appear primed or subscripted. In our Haskell implementation, we represent alphabet symbols using type `Char`.

data `Regex` = $\emptyset \mid \epsilon \mid \text{Chr Char} \mid \text{Regex} \bullet \text{Regex}$
 $\mid \text{Regex} + \text{Regex} \mid \text{Star Regex}$

Constructors \emptyset and ϵ denote respectively the empty set (\emptyset) and the empty string (ϵ) REs. Alphabet symbols are constructed by using the `Chr` constructor. Bigger REs are built using concatenation (\bullet), union ($+$) and Kleene star (`Star`).

Following common practice [20, 24, 25], we adopt an inductive characterization of RE membership semantics. We let judgement $s \in \llbracket e \rrbracket$ denote that string s is in the language denoted by RE e .

Rule *Eps* states that the empty string (denoted by the ϵ) is in the language of RE ϵ .

For any single character a , the singleton string a is in the RL for `Chr a`. Given membership proofs for REs e and e' , $s \in \llbracket e \rrbracket$ and $s' \in \llbracket e' \rrbracket$, rule *Cat* can be used to build a proof for the concatenation of these REs. Rule *Left* (*Right*) creates a membership proof for $e + e'$ from a proof for e (e'). Semantics for Kleene star is built using the following well known equivalence of REs: $e^* = \epsilon + e e^*$.

We say that a RE e is *problematic* if $e = e'^*$ and $\epsilon \in \llbracket e' \rrbracket$ [13]. In this work, we limit our attention to non-problematic RE's. Our

$$\begin{array}{c} \frac{}{\epsilon \in \llbracket \epsilon \rrbracket} \{Eps\} \qquad \frac{a \in \Sigma}{a \in \llbracket a \rrbracket} \{Chr\} \\[10pt] \frac{s \in \llbracket e \rrbracket}{s \in \llbracket e + e' \rrbracket} \{Left\} \qquad \frac{s' \in \llbracket e' \rrbracket}{s' \in \llbracket e + e' \rrbracket} \{Right\} \\[10pt] \frac{}{\epsilon \in \llbracket e^* \rrbracket} \{StarBase\} \qquad \frac{s \in \llbracket e \rrbracket \quad s' \in \llbracket e^* \rrbracket}{ss' \in \llbracket e^* \rrbracket} \{StarRec\} \\[10pt] \frac{s \in \llbracket e \rrbracket \quad s' \in \llbracket e' \rrbracket}{ss' \in \llbracket ee' \rrbracket} \{Cat\} \end{array}$$

Figure 1. RE inductive semantics.

results can be extended to problematic REs without providing any new insight [13, 23].

2.2 RE parsing and bit-coded parse trees

RE parsing. One way to represent parsing evidence is to build a tree that denotes a RE membership proof. Following [13, 23], we let parse trees be terms whose type is underlying RE.

data `Tree` = $() \mid \text{Chr Char} \mid \text{Tree} \bullet \text{Tree} \mid \text{InL Tree}$
 $\mid \text{InR Tree} \mid \text{List [Tree]}$

Constructor $()$ denotes a tree for RE ϵ and `Chr` is a tree for a single character RE. Trees for concatenations are pairs, constructors `InL` and `InR` denotes trees for the left and right component of a choice operator. Finally, a tree for RE e^* is a list of trees for RE e . This informal relation is specified by the following inductive relation between parse trees and RE. We let $\vdash t : e$ denote that t is a parse tree for RE e .

$$\begin{array}{c} \frac{}{\vdash () : \epsilon} \qquad \frac{}{\vdash \text{Chr } a : a} \qquad \frac{\vdash t : e}{\vdash \text{InL } t : e + e'} \\[10pt] \frac{\vdash t' : e'}{\vdash \text{InR } t' : e + e'} \qquad \frac{\vdash t : e \quad \vdash t' : e'}{\vdash t \bullet t' : ee'} \qquad \frac{\forall t. t \in ts \rightarrow \vdash t : e}{\vdash \text{List } ts : e^*} \end{array}$$

Figure 2. Parse tree typing relation.

The relation between RE semantics and its parse trees are formalized using the function `flat`, which builds the string stored in a given parse tree. The Haskell implementation of `flat` is immediate.

`flat` :: `Tree` → `String`

`flat ()` = ""

`flat (Chr c)` = [c]

`flat (t • t')` = `flat t` # `flat t'`

`flat (InL t)` = `flat t`

`flat (InR t)` = `flat t`

`flat (List ts)` = `concatMap flat ts`

The next theorem, which relates parse trees and RE semantics, can be proved by an easy induction on the RE semantics derivation.

Theorem 1. *For all s and e , if $s \in \llbracket e \rrbracket$ then exists a tree t such that `flat t` = s and $\vdash t : e$.*

Bit-coded parse trees. Nielsen et. al. [23] proposed the use of bit-marks to register which branch was chosen in a parse tree for union operator, $+$, and to delimit different matches done by Kleene star expression. Evidently, not all bit sequences correspond to valid parse trees. Ribeiro et. al. [25] showed an inductively defined relation between valid bit-codes and RE, accordingly to the encoding proposed by [23]. We let the judgement $bs \triangleright e$ denote that the sequence of bits bs corresponds to a parse-tree for RE e .

$$\begin{array}{c}
\frac{}{[] \triangleright \epsilon} \quad \frac{}{[] \triangleright a} \quad \frac{bs \triangleright e}{0_b : bs \triangleright e + e'} \\
\frac{bs \triangleright e'}{1_b : bs \triangleright e + e'} \quad \frac{bs \triangleright e \quad bs' \triangleright e'}{bs \# bs' \triangleright ee'} \quad \frac{}{[1_b] \triangleright e^*} \\
\frac{bs \triangleright e \quad bss \triangleright e^*}{0_b : bs \# bss \triangleright e^*}
\end{array}$$

Figure 3. Typing relation for bit-codes.

The empty string and single character RE are both represented by empty bit lists. Codes for RE ee' are built by concatenating codes of e and e' . In RE union operator, $+$, the bit 0_b marks that the parse tree for $e + e'$ is built from e 's and bit 1_b that it is built from e' 's. For the Kleene star, we use bit 1_b to denote the parse tree for the empty string and bit 0_b to begin matchings of e in a parse tree for e^* .

The relation between a bit-code and its underlying parse tree can be defined using functions `code` and `decode`. Type `Code` used in `code` and `decode` definition is just a synonym for `[Bit]`. Function `code` has an immediate definition by recursion on the structure of parse tree.

```

code :: Tree → Regex → Code
code (InL t) (e + _) = 0b : code t e
code (InR t') (_, + e') = 1b : code t' e'
code (List ts) (Star e) = 0b : codeList ts e
code (t • t') (e • e') = code t e # code t' e'
code _ _ = []

codeList :: [Tree] → Regex → Code
codeList ts e = foldr (λt ac → 0b : code t e # ac) [1b] ts

```

To define function `decode`, we need to keep track of the remaining bits to be processed to finish tree construction. This task is done by an auxiliar definition, `dec`.

```

dec :: Regex → Code → Maybe (Tree, Code)
dec ε bs = return ((), bs)
dec (Chr c) bs = return (Chr c, bs)
dec (e + _) (0b : bs) = do
  (t, bs1) ← dec e bs
  return (InL t, bs1)
dec (_, + e') (1b : bs) = do
  (t', bs1) ← dec e' bs
  return (InR t', bs1)
dec (e • e') bs = do
  (t, bs1) ← dec e bs
  (t', bs') ← dec e' bs1

```

```

return (t • t', bs')
dec (Star e) bs = do
  (ts, bs') ← decodeList e bs
  return (List ts, bs')
dec _ _ = fail "invalid bit code"

```

For single character and empty string REs, its decoding consists in just building the tree and leaving the input bit-coded untouched. We build a left tree (using `InL`) for $e + e'$ if the code starts with bit 0_b . A parse tree using constructor `InR` is built whenever we find bit 1_b for a union RE. Building a tree for concatenation is done by sequencing the processing of codes for left component of concatenation and starting the processing of right component with the remaining bits from the processing of the left RE.

```

decodeList :: Regex → Code → Maybe ([Tree], Code)
decodeList _ [] = fail "fail decodeList"
decodeList _ (1b : bs) = return ([], bs)
decodeList e (0b : bs) = do
  (t, be) ← dec e bs
  (ts, bs') ← decodeList e be
  return (t : ts, bs')

```

Function `decodeList` generate a list of parse trees consuming the bit 0_b used as a separator, and bit 1_b which finish the list of parsing results for star operator.

Finally, using `dec`, the definition of `decode` is immediate.

```

decode :: Regex → Code → Maybe Tree
decode e bs
  = case dec e bs of
    Just (t, []) → Just t
    _ → Nothing

```

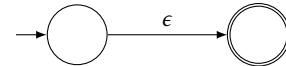
The relation between codes and its correspondent parse trees are specified by the next theorem.

Theorem 2. Let t be a parse tree such that $\vdash t : e$, for some RE e . Then $(\text{code } t) \triangleright e$ and $\text{decode } e (\text{code } t) = \text{Just } t$.

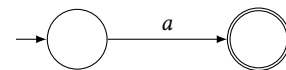
Next, we review Thompson NFA construction which is similar to the proposed semantics for RE parsing developed in Section 3.

2.3 Thompson NFA construction

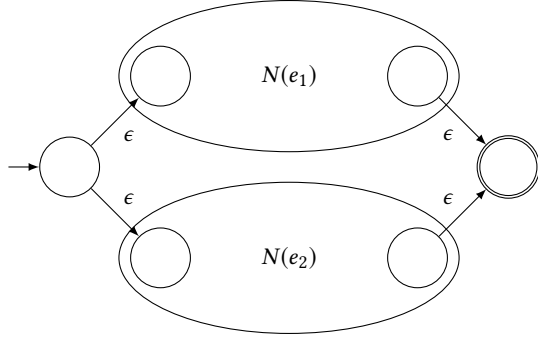
The Thompson NFA construction is a classical algorithm for building an equivalent NFA with ϵ -transitions by induction over the structure of an input RE. We follow a presentation given in [2] where $N(e)$ denotes the NFA equivalent to RE e . The construction proceeds as follows. If $e = \epsilon$, we can build the following NFA equivalent to e .



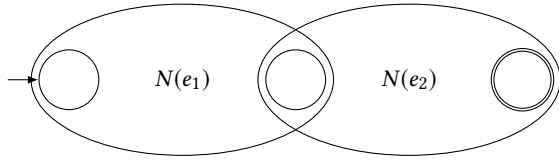
If $e = a$, for $a \in \Sigma$, we can make a NFA with a single transition consuming a :



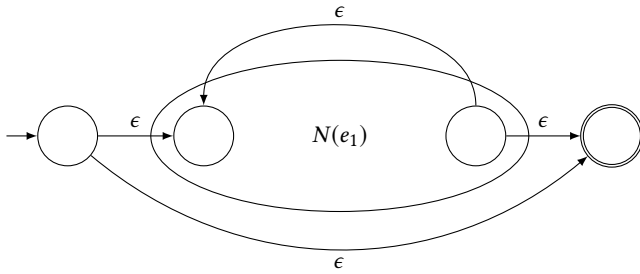
When $e = e_1 + e_2$, we let $N(e_1)$ be the NFA for e_1 and $N(e_2)$ the NFA for e_2 . The NFA for $e_1 + e_2$ is built by adding a new initial and accepting state which can be combined with $N(e_1)$ and $N(e_2)$ using ϵ -transitions as shown in the next picture.



The NFA for the concatenation $e = e_1e_2$ is built from the NFAs $N(e_1)$ and $N(e_2)$. The accepting state of $N(e_1e_2)$ will be the accepting state from $N(e_2)$ and the starting state of $N(e_1)$ will be the initial state of $N(e_1)$.



Finally, for the Kleene star operator, we built a NFA for the RE e , add a new starting and accepting states and the necessary ϵ transitions, as shown below.



Originally, Thompson formulate its construction as a IBM 7094 program [27]. Next we reformulate it as a small-step operational semantics using contexts, modelled as data-type derivatives for RE, which is the subject of the next section.

2.4 Data-type derivatives

The usage of evaluation contexts is standard in reduction semantics [9]. Contexts for evaluating a RE during the parse of a string s can be defined by the following context-free syntax:

$$E[] \rightarrow E[] + e \mid e + E[] \mid E[] e \mid e E[] \mid \star$$

The semantics of a $E[]$ context is a RE with a hole that needs to be “filled” to form a RE. We have two cases for union and concatenation denoting that the hole could be the left or the right component of such operators. Since the Kleene star has only a recursive occurrence, it is denoted just as a “mark” in context syntax.

Having defined our semantics (Figure 4), we have noticed that our RE context syntax is exactly the data type for *one-hole contexts*, known as derivative of an algebraic data type. Derivatives were introduced by McBride and its coworkers [21] as a generalization of Huet’s zippers for a large class of algebraic data types [1]. RE contexts are implemented by the following Haskell data-type:

```
data Hole = InChoiceL Regex | InChoiceR Regex
          | InCatL Regex | InCatR Regex | InStar
```

Constructor **InChoiceL** store the right component of a union RE (similarly for **InChoiceR**). We need to store contexts for union because such information is used to allow backtracking in case of failure. Constructors **InCatL** and **InCatR** store the right (left) component of a concatenation and they are used to store the next subexpressions that needed to be evaluated during input string parsing. Finally, **InStar** marks that we are currently processing an expression with a Kleene star operator.

3 Proposed semantics

In this section we present the definition of an operational semantics for RE parsing which is equivalent to executing the Thompson’s construction NFA over the input string. Observe that, the inductive semantics for RE (Figure 1) can be understood as a big-step operational semantics for RE, since it ignores many details on how should we proceed to match an input [24].

The semantics is defined as a binary relation between *configurations*, which are 5-uples $\langle d, e, c, b, s \rangle$ where:

- d is a direction, which specifies if the semantics is starting (denoted by B) or finishing (F) the processing of the current expression e .
- e is the current expression being evaluated;
- c is a context in which e occurs. Contexts are just a list of **Hole** type in our implementation.
- b is a bit-code for the current parsing result, in reverse order.
- s is the input string currently being processed.

Notation $\langle d, e, c, b, s \rangle \rightarrow \langle d', e', c', b', s' \rangle$ denotes that from configuration $\langle d, e, c, b, s \rangle$ we can give a step leading to a new state $\langle d', e', c', b', s' \rangle$ using the rules specified in Figure 4.

The rules of the semantics can be divided in two groups: starting rules and finishing rules. Starting rules deal with configurations with a begin (B) direction and denote that we are beginning the parsing for its RE e . Finishing rules use the context to decide how the parsing for some expression should end. Intuitively, starting rules correspond to transitions entering a sub-automata of Thompson NFA and finishing rules to transitions exiting a sub-automata.

The meaning of each starting rule is as follows. Rule $\{Eps\}$ specifies that we can mark a state as finished if it consists of a starting configuration with RE ϵ . We can finish any configuration for **Chr** a if it is starting with current string with a leading a . Whenever we have a starting configuration with a choice RE, $e_1 + e_2$, we can non-deterministically choose if input string s can be processed by e_1 (rule $Left_B$) or e_2 (rule $Right_B$). For beginning configurations with concatenation, we parse input string using each of its components sequentially. Finally, for starting configurations with a Kleene star operator, e^* , we can either start the processing of e or finish the processing for e^* . In all recursive cases for RE, we insert context information in the third component of the resulting configuration in order to decide how the machine should step after finishing the execution of the RE currently on focus.

Rule (Cat_{EL}) applies to any configuration which is finishing with a left concatenation context $(E[]e')$. In such situation, rule specifies that a computation should continue with e' and push the context $e E[]$. We end the computation for a concatenation, whenever we find a context $e E[]$ in the context component (rule (Cat_{ER})). Finishing a computation for choice consists in just popping its

$$\begin{array}{c}
\frac{}{\langle B, \epsilon, c, b, s \rangle \rightarrow \langle F, \epsilon, c, b, s \rangle} (Eps) \quad \frac{}{\langle B, a, c, b, a : s \rangle \rightarrow \langle F, a, c, b, s \rangle} (Chr) \quad \frac{b' = 0_b : b \quad c' = E[] + e' : c}{\langle B, e + e', c, b, s \rangle \rightarrow \langle B, e, c', b', s \rangle} (Left_B) \\
\frac{b' = 1_b : b \quad c' = e + E[] : c}{\langle B, e + e', c, b, s \rangle \rightarrow \langle B, e', c', b', s \rangle} (Right_B) \quad \frac{c' = E[] e' : c}{\langle B, ee', c, b, s \rangle \rightarrow \langle B, e, c', b, s \rangle} (Cat_B) \quad \frac{}{\langle B, e^*, c, b, s \rangle \rightarrow \langle B, e, \star : c, 0_b : b, s \rangle} (Star_1) \\
\frac{}{\langle B, e^*, c, b, s \rangle \rightarrow \langle F, e^*, c, 1_b : b, s \rangle} (Star_2) \quad \frac{c' = eE[] : c}{\langle F, e, E[] e' : c, b, s \rangle \rightarrow \langle B, e', c', b, s \rangle} (Cat_{EL}) \quad \frac{}{\langle F, e', eE[] : c, b, s \rangle \rightarrow \langle F, ee', c, b, s \rangle} (Cat_{ER}) \\
\frac{c = E[] + e' : c'}{\langle F, e, c, b, s \rangle \rightarrow \langle F, e + e', c', 0_b : b, s \rangle} (Left_E) \quad \frac{c = e + E[] : c'}{\langle F, e, c, b, s \rangle \rightarrow \langle F, e + e', c', 1_b : b, s \rangle} (Right_E) \\
\frac{}{\langle F, e, \star : c, b, s \rangle \rightarrow \langle B, e, \star : c, 0_b : b, s \rangle} (Star_{E1}) \quad \frac{}{\langle F, e, \star : c, b, s \rangle \rightarrow \langle F, e^*, c, 1_b : b, s \rangle} (Star_{E2})
\end{array}$$

Figure 4. Small-step semantics for RE parsing.

correspondent context, as done by rules $(Left_E)$ and $(Right_E)$. For the Kleene star operator, we can either finish the computation by popping the contexts and adding the corresponding 1_b to end its matching list or restart with RE e for another matching over the input string.

The starting state of the semantics is given by the configuration $\langle B, e, [], [], s \rangle$ and accepting configurations are $\langle F, e', [], bs, [] \rangle$, for some RE e' and code bs . Following common practice, we let \rightarrow^* denote the reflexive, transitive closure of the small-step semantics defined in Figure 4. We say that a string s is accepted by RE e if $\langle B, e, [], [], s \rangle \rightarrow^* \langle F, e', [], bs, [] \rangle$. The next theorem asserts that our semantics is sound and complete with respect to RE inductive semantics (Figure 1).

Theorem 3. *For all strings s and non-problematic REs $e, s \in \llbracket e \rrbracket$ if, and only if, $\langle B, e, [], [], s \rangle \rightarrow^* \langle F, e', [], bs, [] \rangle$ and $\langle F, e', [], bs, [] \rangle$ is an accepting configuration.*

4 Implementation details

In order to implement the small-step semantics of Figure 4, we need to represent configurations. We use type **Conf** to denote configurations and directions are represented by type **Dir**, where **Begin** denote the starting and **End** the finishing direction.

data **Dir** = **Begin** | **End**

type **Conf** = (**Dir**, **Regex**, [**Hole**], **Code**, **String**)

Function **finish** tests if a configuration is an accepting one.

finish :: **Conf** → **Bool**

finish (**End**, $_$, $_$, $_$, $_$) = **True**

finish $_$ = **False**

The small-step semantics is implemented by function **next**, which returns a list of configurations that can be reached from a given input configuration. We will begin by explaining the equations that code the set of starting rules from the small-step semantics. The first alternative

next :: **Conf** → [**Conf**]

next (**Begin**, ϵ , ctx , bs , s) = [**(End**, ϵ , ctx , bs , s)]

implements rule (Eps) , which finishes a starting **Conf** with an ϵ . Rule (Chr) is implemented by the following equation

next (**Begin**, **Chr** c , ctx , bs , $a : s$)
 | $a \equiv c = [(\text{End}, **Chr** c , ctx , bs , s)]
 | **otherwise** = []$

which consumes input character a if it matches RE **Chr** c , otherwise it fails by returning an empty list. For a choice expression, we can use two distinct rules: one for parsing the input using its left component and another rule for the right. Since both union and Kleene star introduce non-determinism in RE parsing, we can easily model this using the list monad, by return a list of possible resulting configurations.

next (**Begin**, $e + e'$, ctx , bs , s)
 = [**(Begin**, e , **InChoiceL** e' : ctx , 0_b : bs , s)
 , **(Begin**, e' , **InChoiceR** e : ctx , 1_b : bs , s)]

Concatenation just sequences the computation of each of its composing RE.

next (**Begin**, $e \bullet e'$, ctx , bs , s)
 = [**(Begin**, e , **InCatL** e' : ctx , bs , s)]

For a starting configuration with Kleene star operator, **Star** e , we can proceed in two ways: by beginning the parsing of RE e or by finishing the computation for **Star** e over the input.

next (**Begin**, **Star** e , ctx , bs , s)
 = [**(Begin**, e , **InStar** : ctx , 0_b : bs , s)
 , **(End**, (**Star** e), ctx , 1_b : bs , s)]

The remaining equations of **next** deal with operational semantics finishing rules. The equation below implements rule (Cat_{EL}) which specifies that an ended computation for the left component of a concatenation should continue with its right component.

next (**End**, e , **InCatL** e' : ctx , bs , s)
 = [**(Begin**, e' , **InCatR** e : ctx , bs , s)]

Whenever we are in a finishing configuration with a right concatenation context, (**InCatR** e), we end the parsing of the input for the whole concatenation RE.


```
next (End, e', InCatR e : ctx, bs, s)
  = [(End, e • e', ctx, bs, s)]
```

Next equations implement the rules that finish configurations for the union, by committing to its first successful branch.

```
next (End, e, InChoiceL e' : ctx, bs, s)
  = [(End, e + e', ctx, 0b : bs, s)]
next (End, e', InChoiceR e : ctx, bs, s)
  = [(End, e + e', ctx, 1b : bs, s)]
```

Equations for Kleene star implement rules ($Star_{E1}$) and ($Star_{E2}$) which allows ending or add one more match for an RE e .

```
next (End, e, InStar : ctx, bs, s)
  = [(Begin, e, InStar : ctx, 0b : bs, s)
    , (End, (Star e), ctx, 1b : bs, s)]
```

Finally, stuck states on the semantics are properly handled by the following equation which turns them all into a failure (empty list).

```
next _ = []
```

The reflexive-transitive closure of the semantics is implemented by function `steps`, which computes the trace of all states needed to determine if a string can be parsed by the RE e .

```
steps :: [Conf] → [Conf]
steps [] = []
steps cs = steps [c' | c ← cs, c' ← next c] # cs
```

Finally, the function for parsing a string using an input RE is implemented as follows:

```
vmAccept :: String → Regex → (Bool, Code)
vmAccept s e = let r = [c | c ← steps initcfg, finish c]
  in if null r then (False, []) else (True, bitcode (head r))
  where
    initcfg = [(Begin, e, [], [], s)]
    bitcode (→, →, →, bs, →) = reverse bs
```

Function `vmAccept` returns a pair formed by a boolean and the bit-code produced during the parsing of an input string and RE. Observe that, we need to reverse the bit-codes, since they are built in reverse order.

4.1 Test suite

An overview of QuickCheck. Our tests are implemented using QuickCheck [5], a library that allows the testing of properties expressed as Haskell functions. Such verification is done by generating random values of the desired type, instantiating the relevant property with them, and checking it directly by evaluating it to a boolean. This process continues until a counterexample is found or a specified number of cases are tested with success. The library provides generators for several standard library data types and combinators to build new generators for user-defined types.

As an example of a custom generator, consider the task of generating a random alpha-numeric character. To implement such generator, `genChar`, we use QuickCheck function `suchThat` which generates a random value with satisfies a predicate passed as argument (in example, we use `isAlphaNum` which is true whenever we pass an alpha-numeric character to it), using an random generator taken as input.

```
genChar :: Gen Char
genChar = suchThat (arbitrary :: Gen Char) isAlphaNum
```

Test case generators. In order to test the correctness of our semantics, we needed to build generators for REs and for strings. We develop functions to randomly generate strings accepted and rejected for a RE, using the QuickCheck library.

Generation of random RE is done by function `sizedRegex` with takes a depth limit to restrict the size of the generated RE. Whenever the input depth limit is less or equal to 1, we can only build a ϵ or a single character RE. The definition of `sizedRegex` uses QuickCheck function `frequency`, which receives a list of pairs formed by a weight and a random generator and produce, as result, a generator which uses such frequency distribution. In `sizedRegex` implementation we give a higher weight to generate characters and equal distributions to build concatenation, union or star.

```
sizedRegex :: Int → Gen Regex
sizedRegex n
  | n ≤ 1 = frequency [(10, return ε), (90, Chr ($) genChar)]
  | otherwise = frequency [(10, return ε), (30, Chr ($) genChar)
    , (20, (•) ($) sizedRegex n2 (★) sizedRegex n2)
    , (20, (+) ($) sizedRegex n2 (★) sizedRegex n2)
    , (20, Star ($) suchThat (sizedRegex n2) (not ◦ nullable))]
  where n2 = div n 2
```

Given an RE e , we can generate a random string s such that $s \in \llbracket e \rrbracket$ using the next definition. We generate strings by choosing randomly between branches of a union or by repeating n times a string s which is accepted by e , whenever we have e^* (function `randomMatches`).

```
randomMatch :: Regex → Gen String
randomMatch ε = return ""
randomMatch (Chr c) = return [c]
randomMatch (e • e') = liftM2 (++) (randomMatch e)
  (randomMatch e')
randomMatch (e + e') = oneof [randomMatch e, randomMatch e']
randomMatch (Star e) = do
  n ← choose (0, 3) :: Gen Int
  randomMatches n e
randomMatches :: Int → Regex → Gen String
randomMatches m e'
  | m ≤ 0 = return []
  | otherwise = liftM2 (++) (randomMatch e')
    (randomMatches (m - 1) e')
```

The algorithm for generating random strings that aren't accepted by a RE is similarly defined and omitted for brevity.

Properties considered. In order to verify if the defined semantics is correct, we need to check the following properties:

- Our semantics accepts only and all the strings in the language described by the input RE: we test this property by generating random strings that should be accepted and strings that must be rejected by a random RE.
- Our semantics generates valid parsing evidence: the bit-codes produced as result, have the following properties: 1) the bit-codes can be parsed into a valid parse tree t for the

random produced RE e , i.e. $\vdash t : e$ holds ; 2) `flat` $t = s$ and 3) `code` $e = \text{bs}$.

In addition to coding / decoding of parse trees, we need a function which checks if a tree is indeed a parsing evidence for some RE e . Function `tc` takes, as arguments, a parse tree t and a RE e and verifies if t is an evidence for e .

```
tc :: Tree → Regex → Bool
tc () ε = True
tc (Chr c) (Chr c') = c == c'
tc (t • t') (e • e') = tc t e ∧ tc t' e'
tc (InL t) (e + _) = tc t e
tc (InR t') (_ + e') = tc t' e'
tc (List ts) (Star e) = all (flip tc e) ts
```

Function `tc` is a implementation of parsing tree typing relation, as specified by the following result.

Theorem 4. *For all tree t and RE e , $\vdash t : e$ if, and only if, `tc` $t e = \text{True}$.*

Code coverage results. After running thousands of well-succeeded tests, we gain a high degree of confidence in the correctness of our semantics, however, it is important to measure how much of our code is covered by the test suite. We use the Haskell Program Coverage tool (HPC) [14] to generate statistics about the execution of our tests. Code coverage results are presented in Figure 5.

Top Level Definitions		Alternatives		Expressions	
%	covered / total	%	covered / total	%	covered / total
100%	3/3	100%	10/10	100%	74/74
100%	4/4	100%	18/18	97%	163/167
-	0/0	-	0/0	-	0/0
100%	7/7	100%	21/21	100%	173/173
100%	7/7	100%	25/25	100%	142/142
100%	21/21	100%	74/74	99%	552/556

Figure 5. Code coverage results

Our test suite give us almost 100% of code coverage, which provides a strong evidence that our semantics is indeed correct. All top level definitions and function alternatives are actually executed by the test cases and just two expressions are marked as non-executed by HPC.

5 Related work

Ierusalimsky [16] proposed the use of Parsing Expression Grammars (PEGs) as a basis for pattern matching. He argued that pure REs is a weak formalism for pattern-matching tasks: many interesting patterns either are difficult to describe or cannot be described by REs. He also said that the inherent non-determinism of REs does not fit the need to capture specific parts of a match. Following this proposal, he presented LPEG, a pattern-matching tool based on PEGs for the Lua language. He argued that LPEG unifies the ease of use of pattern-matching tools with the full expressive power of PEGs. He also presented a parsing machine (PM) that allows an implementation of PEGs for pattern matching. In [22], Medeiros et. al. presents informal correctness proofs of LPEG PM. While such proofs represent a important step towards the correctness of LPEG, there is no guarantee that LPEG implementation follows its specification.

In [24], Rathnayake and Thielecke formalized a VM implementation for RE matching using operational semantics. Specifically, they

derived a series of abstract machines, moving from the abstract definition of matching to realistic machines. First, a continuation is added to the operational semantics to describe what remains to be matched after the current expression. Next, they represented the expression as a data structure using pointers, which enables redundant searches to be eliminated via testing for pointer equality. Although their work has some similarities with ours (a VM-based parsing of REs), they did not present any evidence or proofs that their VM is correct.

Fischer, Huch and Wilke [11] developed a Haskell program for matching REs. The program is purely functional and it is overloaded over arbitrary semirings, which solves the matching problem and supports other applications like computing leftmost longest matchings or the number of matchings. Their program can also be used for parsing every context-free language by taking advantage of laziness. Their developed program is based on an old technique to turn REs into finite automata, which makes it efficient compared to other similar approaches. One advantage of their implementation over our proposal is that their approach works with context-free languages, not only with REs purely. However, they did not present any correctness proofs of their Haskell code.

Cox [6] said that viewing RE matching as executing a special machine makes it possible to add new features just by the inclusion of new machine instructions. He presented two different ways to implement a VM that executes a RE that has been compiled into byte-codes: a recursive and a non-recursive backtracking implementation, both in C programming language. Cox's work on VM-based RE parsing is poorly specified: both the VM semantics and the RE compilation process are described only informally and no correctness guarantees is even mentioned.

Frisch and Cardelli [13] studied the theoretical problem of matching a flat sequence against a type (RE): the result of the process is a structured value of a given type. Their contributions were in noticing that: (1) A disambiguated result of parsing can be presented as a data structure that does not contain ambiguities. (2) There are problematic cases in parsing values of star types that need to be disambiguated. (3) The disambiguation strategy used in XDuce and CDuce (two XML-oriented functional languages) pattern matching can be characterized mathematically by what they call greedy RE matching. (4) There is a linear time algorithm for the greedy matching. Their approach is different since they want to axiomatize abstractly the disambiguation policy, without providing an explicit matching algorithm. They identify three notions of problematic words, REs, and values (which represent the ways to match words), relate these three notions, and propose matching algorithms to deal with the problematic case.

Ribeiro and Du Bois [25] described the formalization of a RE parsing algorithm that produces a bit representation of its parse tree in the dependently typed language Agda. The algorithm computes bit-codes using Brzozowski derivatives and they proved that the produced codes are equivalent to parse trees ensuring soundness and completeness with respect to an inductive RE semantics. They included the certified algorithm in a tool developed by themselves, named verigrep, for RE-based search in the style of GNU grep. While the authors provide formal proofs, their tool show a bad performance when compared with other approaches to RE parsing. Besides, they did not prove that their algorithm follows some disambiguation policy, like POSIX or greedy.

Nielsen and Henglein [23] showed how to generate a compact *bit-coded* representation of a parse tree for a given RE efficiently, without explicitly constructing the parse tree first, by simplifying the DFA-based parsing algorithm of Dubé and Feeley [8] to emit a bit representation without explicitly materializing the parse tree itself. They also showed that Frisch and Cardelli's greedy RE parsing algorithm [13] can be straightforwardly modified to produce bit codings directly. They implemented both solutions as well as a backtracking parser and performed benchmark experiments to measure their performance. They argued that bit codings are interesting in their own right since they are typically not only smaller than the parse tree, but also smaller than the string being parsed and can be combined with other techniques for improved text compression. As others related works, the authors did not present a formal verification of their implementations.

An algorithm for POSIX RE parsing is described in [26]. The main idea of the article is to adapt derivative parsing to construct parse trees incrementally to solve both matching and submatching for REs. In order to improve the efficiency of the proposed algorithm, Sulzmann et al. use a bit encoded representation of RE parse trees. Textual proofs of correctness of the proposed algorithm are presented in an appendix.

6 Conclusion

In this work, we presented a small-step operational semantics for a virtual machine for RE parsing inspired on Thompson's NFA construction. Our semantics produces, as parsing evidence, bit-codes which can be used to characterize which disambiguation strategy is followed by the semantics. We use data-type derivatives to represent evaluation contexts for RE. Such contexts are used to decide how to finish the execution of the RE on focus. We have developed a prototype implementation of our semantics in Haskell and use QuickCheck to verify its relevant properties with respect to a simple implementation of RE parsing by Fisher et. al. [11].

Currently, we have a formalized interpreter of our semantics in Coq proof assistant [4] available at project's on-line repository [7]. We are working on formalizing the equivalence between the proposed semantics and the standard RE inductive semantics.

As future work we intend to use our verified semantics to build a certified tool for RE parsing, work on proofs that the semantics follow a specific disambiguation strategy and investigate how other algorithms (e.g. the Glushkov construction [15]) for converting a RE into a finite state machine could be expressed in terms of an operational semantics.

Acknowledgements

We would like to thank Prof. *Leonardo Vieira, Samuel Feitosa* and the anonymous reviewers for their valuable suggestions and comments on early versions of this paper.

References

- [1] Michael Gordon Abbott, Thorsten Altenkirch, Neil Ghani, and Conor McBride. 2003. Derivatives of Containers. In *Typed Lambda Calculi and Applications, 6th International Conference, TLCA 2003, Valencia, Spain, June 10-12, 2003, Proceedings. (Lecture Notes in Computer Science)*, Martin Hofmann (Ed.), Vol. 2701. Springer, 16–30. https://doi.org/10.1007/3-540-44904-3_2
- [2] Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. 1986. *Compilers: Principles, Techniques, and Tools*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA.
- [3] Andrea Asperti, Claudio Sacerdoti Coen, and Enrico Tassi. 2010. Regular Expressions, au point. *CoRR* abs/1010.2604 (2010). [arXiv:1010.2604](http://arxiv.org/abs/1010.2604) <http://arxiv.org/abs/1010.2604>
- [4] Yves Bertot and Pierre Castran. 2010. *Interactive Theorem Proving and Program Development: Coq'Art The Calculus of Inductive Constructions* (1st ed.). Springer Publishing Company, Incorporated.
- [5] Koen Claessen and John Hughes. 2000. QuickCheck: A Lightweight Tool for Random Testing of Haskell Programs. In *Proceedings of the Fifth ACM SIGPLAN International Conference on Functional Programming (ICFP '00)*. ACM, New York, NY, USA, 268–279.
- [6] Russ Cox. 2009. Regular Expression Matching: the Virtual Machine Approach. (2009). <https://swtch.com/>
- [7] Thales Delfino and Rodrigo Ribeiro. 2018. Towards certified virtual machine-based regular expression parsing — On-line repository. <https://github.com/thalesad/regexvm>. (2018).
- [8] Danny Dubé and Marc Feeley. 2000. Efficiently Building a Parse Tree from a Regular Expression. *Acta Inf.* 37, 2 (Oct. 2000), 121–144. <https://doi.org/10.1007/s002360000037>
- [9] Matthias Felleisen, Robert Bruce Findler, and Matthew Flatt. 2009. *Semantics Engineering with PLT Redex* (1st ed.). The MIT Press.
- [10] Denis Firsov and Tarmo Uustalu. 2013. Certified Parsing of Regular Languages. In *Proceedings of the Third International Conference on Certified Programs and Proofs - Volume 8307*. Springer-Verlag New York, Inc., New York, NY, USA, 98–113. https://doi.org/10.1007/978-3-319-03545-1_7
- [11] Sebastian Fischer, Frank Huch, and Thomas Wilke. 2010. A play on regular expressions. *ACM SIGPLAN Notices* 45, 9 (2010), 357. <https://doi.org/10.1145/1932681.1863594>
- [12] Bryan Ford. 2004. Parsing Expression Grammars: A Recognition-based Syntactic Foundation. In *Proceedings of the 31st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '04)*. ACM, New York, NY, USA, 111–122. <https://doi.org/10.1145/964001.964011>
- [13] Alain Frisch and Luca Cardelli. 2004. Greedy Regular Expression Matching. *ICALP 2004 - International Colloquium on Automata, Languages and Programming* 3142 (2004), 618–629.
- [14] Andy Gill and Colin Runciman. 2007. Haskell Program Coverage. In *Proceedings of the ACM SIGPLAN Workshop on Haskell Workshop (Haskell '07)*. ACM, New York, NY, USA, 1–12. <https://doi.org/10.1145/1291201.1291203>
- [15] Victor M. Glushkov. 1961. The Abstract Theory of Automata. *Russian Mathematical Surveys* 16, 5 (1961), 1–53.
- [16] Roberto Ierusalimsky. 2009. A text patternmatching tool based on parsing expression grammars. *Software - Practice and Experience* (2009). <https://doi.org/10.1002/spe.892>
- [17] Donald E. Knuth. 1971. Top-down Syntax Analysis. *Acta Inf.* 1, 2 (June 1971), 79–110. <https://doi.org/10.1007/BF00289517>
- [18] Miran Lipovaca. 2011. *Learn You a Haskell for Great Good!: A Beginner's Guide* (1st ed.). No Starch Press, San Francisco, CA, USA.
- [19] A. Loh. [n. d.]. Typesetting Haskell and more with lhs2TeX. ([n. d.]). <http://www.cs.uu.nl/~>
- [20] Raul Lopes, Rodrigo Ribeiro, and Carlos Camarão. 2016. Certified Derivative-Based Parsing of Regular Expressions. In *Programming Languages — Lecture Notes in Computer Science 9889*. Springer, 95–109.
- [21] Conor McBride. 2008. Clowns to the left of me, jokers to the right (pearl): dissecting data structures. In *Proceedings of the 35th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2008, San Francisco, California, USA, January 7-12, 2008*. 287–295. <https://doi.org/10.1145/1328438.1328474>
- [22] Sérgio Medeiros and Roberto Ierusalimsky. 2008. A parsing machine for PEGs. *Proceedings of the 2008 symposium on Dynamic languages - DLS '08* (2008), 1–12. <https://doi.org/10.1145/1408681.1408683>
- [23] Lasse Nielsen and Fritz Henglein. 2011. Bit-coded Regular Expression Parsing. In *Language and Automata Theory and Applications*, Adrian-Horia Dediu, Shunsuke Inenaga, and Carlos Martín-Vide (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 402–413.
- [24] Asiri Rathnayake and Hayo Thielecke. 2011. Regular Expression Matching and Operational Semantics. *Electronic Proceedings in Theoretical Computer Science* 62, Sos (2011), 31–45. <https://doi.org/10.4204/EPTCS.62.3> [arXiv:1108.3126](https://arxiv.org/abs/1108.3126)
- [25] Rodrigo Ribeiro and André Du Bois. 2017. Certified Bit-Coded Regular Expression Parsing. *Proceedings of the 21st Brazilian Symposium on Programming Languages - SBLP 2017* (2017), 1–8. <http://dl.acm.org/citation.cfm?doid=3125374.3125381>
- [26] Martin Sulzmann and Kenny Zhuo Ming Lu. 2014. POSIX Regular Expression Parsing with Derivatives. In *Functional and Logic Programming - 12th International Symposium, FLOPS 2014, Kanazawa, Japan, June 4-6, 2014. Proceedings (Lecture Notes in Computer Science)*, Michael Codish and Eijiro Sumii (Eds.), Vol. 8475. Springer, 203–220. https://doi.org/10.1007/978-3-319-07151-0_13
- [27] Ken Thompson. 1968. Programming Techniques: Regular Expression Search Algorithm. *Commun. ACM* 11, 6 (June 1968), 419–422.