

Certified Bit-Coded Regular Expression Parsing

Rodrigo Ribeiro
Universidade Federal de Ouro Preto
Ouro Preto, Minas Gerais — Brazil
rodrigo@decsi.ufop.br

André Du Bois
Universidade Federal de Pelotas
Pelotas, Rio Grande do Sul — Brazil
dubois@inf.ufpel.br

ABSTRACT

We describe the formalization of a regular expression (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 we prove that produced codes are equivalent to parse trees ensuring soundness and completeness w.r.t an inductive RE semantics. We include the certified algorithm in a tool developed by us, named *verigrep*, for regular expression based search in the style of the well known GNU *grep*. Practical experiments conducted with this tool are reported.

KEYWORDS

Certified algorithms, regular expressions, bit-codes, dependent types

ACM Reference format:

Rodrigo Ribeiro and André Du Bois. 2017. Certified Bit-Coded Regular Expression Parsing. In *Proceedings of SBLP 2017, Fortaleza, CE, Brazil, September 21–22, 2017*, 8 pages.
DOI: 10.1145/3125374.3125381

1 INTRODUCTION

Parsing is one of the most studied problems in computer science. It involves the process of checking if a string of symbols conforms to a given set of rules. Usually, parsing is preceded by the specification of rules in a formalism (e.g. a grammar) and, also, either the construction of data that makes evident which rules have been used to conclude that the string of symbols can be obtained from it or, otherwise, an indication of an error that represents the fact that the string of symbols cannot be generated.

In this work we are interested in the parsing problem for regular languages (RLs) [16], i.e. languages that can be recognized by non-deterministic finite automata and equivalent formalisms. Regular expressions (REs) are an algebraic and compact way of specifying RLs that are extensively used in lexical analyser generators [18] and string search utilities [14]. Since such tools are widely used and parsing is pervasive in computing, there is a growing interest on certified parsing algorithms [7, 9, 10]. This interest is motivated by the recent development of dependently typed languages. Such languages are powerful enough to express algorithmic properties as types, that are automatically checked by a compiler.

Permission to make digital or hard copies of all or part 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 components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SBLP 2017, Fortaleza, CE, Brazil

© 2017 ACM. 978-1-4503-5389-2/17/09...\$15.00

DOI: 10.1145/3125374.3125381

In a previous work by one the authors [19], we described the formalization of an algorithm for parsing RE's based on derivatives. RE derivatives were introduced by Brzozowski [6] as an alternative method to compute a finite state machine that is equivalent to a given RE and to perform RE-based parsing. Owens et. al [27] reintroduce this concept, since “derivatives have been lost in the sands of time” until his work on functional encoding of RE derivatives have renewed interest on its use for parsing [12, 23].

Nielsen et. al. [24] introduce algorithms to build a compact representation of RE parse trees as a sequence of bits, without explicitly constructing the parse tree first and describe algorithms that simulate a finite state machine that output the required bit list. Sulzmann et. al. [33] designed a RE derivative-based algorithm that incrementally builds a list of bits instead of a tree as parsing result. In both works no machine checked proof was provided and Sulzmann et. al. informal proof “has some unfillable gaps”, as pointed by Ausaf et. al. which proved that Sulzmann et. al. algorithm produces POSIX parse trees in Isabelle/HOL [3]. In this work we are interested in providing fully certified correctness proofs of Sulzmann et. al. derivative based parsing algorithm. The certified algorithm produces a bit sequence equivalent to traditional parse trees and we use it in a RE based search tool that has been developed by us, using the dependently typed language Agda [26]. Unlike Ausaf et. al., we do not provide mechanized proofs which ensure that our formalization produces POSIX parse trees. We leave such proof for further work.

More specifically, our contributions are:

- A formalization of bit-codes for RE parse trees, as presented in [24], and machine checkable proof of their properties.
- A formalization of bit-annotated regular expressions (BRE), its semantics and its soundness and completeness theorems w.r.t. standard RE semantics. We also relate produced bit sequences and RE semantics proving how they are related to each other.
- A formalization of Brzozowski derivatives for BRE and its soundness and completeness theorem w.r.t to BRE semantics.
- We build certified decision procedures for matching prefixes and substrings of BRE.

The rest of this paper is organized as follows. Section 2 presents a brief introduction to Agda. Section 3 describes the encoding of REs, its parse trees and their representation as bit-codes. In Section 3.2 we present BREs, its semantics and their relation with REs. In Section 4 we formalize a variant of Brzozowski derivatives for BREs, its properties and describe how to build a correct parsing algorithm from them. Section 5 comments on the use of the certified algorithm to build a tool for RE-based search and presents some experiments with this tool. Related work is discussed on Section 6. Section 7 concludes.

All the source code in this article has been formalized in Agda version 2.5.2 using Standard Library 0.13, but we do not present every detail. Proofs of some properties result in functions with a long pattern matching structure, that would distract the reader from understanding the high-level structure of the formalization. In such situations we give just proof sketches and point out where all details can be found in the source code.

All source code produced, including the \LaTeX source of this article, are available on-line [29].

2 AN OVERVIEW OF AGDA

Agda is a dependently-typed functional programming language based on Martin-Löf intuitionistic type theory [21]. Function types and an infinite hierarchy of types of types, $\text{Set } l$, where l is a natural number, are built-in. Everything else is a user-defined type. The type Set , also known as Set_0 , is the type of all “small” types, such as Bool , String and List Bool . The type Set_1 is the type of Set and “others like it”, such as $\text{Set} \rightarrow \text{Bool}$, $\text{String} \rightarrow \text{Set}$, and $\text{Set} \rightarrow \text{Set}$. We have that $\text{Set } l$ is an element of the type $\text{Set } (l + 1)$, for every $l \geq 0$. This stratification of types is used to keep Agda consistent as a logical theory [31].

An ordinary (non-dependent) function type is written $A \rightarrow B$ and a dependent one is written $(x : A) \rightarrow B$, where type B depends on x , or $\forall (x : A) \rightarrow B$. Agda allows the definition of *implicit parameters*, i.e. parameters whose values can be inferred from the context, by surrounding them in curly braces: $\forall \{x : A\} \rightarrow B$. To avoid clutter, we’ll omit implicit arguments from the source code presentation. The reader can safely assume that every free variable in a type is an implicit parameter.

As an example of Agda code, consider the following data type of length-indexed lists, also known as vectors.

```
data N : Set where
  zero : N
  suc : N → N

data Vec (A : Set) : N → Set where
  [] : Vec A zero
  _::_ : ∀ {n} → A → Vec A n → Vec A (suc n)
```

Constructor `[]` builds empty vectors. The cons-operator `(_::_)` inserts a new element in front of a vector of n elements (of type $\text{Vec } A \ n$) and returns a value of type $\text{Vec } A \ (\text{suc } n)$. The Vec datatype is an example of a dependent type, i.e. a type that uses a value (that denotes its length). The usefulness of dependent types can be illustrated with the definition of a safe list head function: `head` can be defined to accept only non-empty vectors, i.e. values of type $\text{Vec } A \ (\text{suc } n)$.

```
head : Vec A (suc n) → A
head (x :: xs) = x
```

In `head`’s definition, constructor `[]` is not used. The Agda type-checker can figure out, from `head`’s parameter type, that argument `[]` to `head` is not type-correct.

Thanks to the propositions-as-types principle¹ we can interpret types as logical formulas and terms as proofs. An example is the representation of equality as the following Agda type:

```
data ≡ {l} {A : Set l} (x : A) : A → Set where
  refl : x ≡ x
```

This type is called propositional equality. It defines that there is a unique evidence for equality, constructor `refl` (for reflexivity), that asserts that the only value equal to x is itself. Given a type P , type $\text{Dec } P$ is used to build proofs that P is a decidable proposition, i.e. that either P or not P holds. The decidable proposition type is defined as:

```
data Dec (P : Set) : Set where
  yes : P → Dec P
  no : ¬ P → Dec P
```

Constructor `yes` stores a proof that property P holds and constructor `no` an evidence that such proof is impossible. Some functions used in our formalization use this type. The type $\neg P$ is an abbreviation for $P \rightarrow \perp$, where \perp is a data type with no constructors (i.e. a data type for which it is not possible to construct a value, which corresponds to a false proposition).

A useful data type in dependently typed languages is the so-called dependent-product, or Σ -types, that generalizes cartesian products. A possible definition of dependent products in Agda is as follows:

```
data Σ {a b} (A : Set a) (B : A → Set b) : Set (a ⊔ b) where
  _-_- : (x : A) → B x → Σ A B
```

A value of type $\Sigma A B$ corresponds to a pair formed by a value x of type A and a value of type $B \ x$. Note that the second component type depends on the value of the first component. Under proposition-as-types interpretation, dependent products are equivalent to existential quantification, since values of dependent products are formed by a value $x : A$, which can be understood as a witness of existential quantification, and a value $B \ x$, which represents the proof that x holds for B . Agda standard library represents existential quantification as a dependent product in which the first component is an implicit parameter:

```
∃ : ∀ {a b} {A : Set a} → (A → Set b) → Set (a ⊔ b)
∃ = Σ _
```

In some functions involving dependent products we use wildcards “_” to avoid the need to explicitly write their values or even completely omit, since Agda type checker can infer them. Finally, we represent projections on dependent products as π_1 and π_2 which recover the first and second component of product types, respectively.

Dependently typed pattern matching is built by using the so-called `with` construct, that allows for matching intermediate values [22]. If the matched value has a dependent type, then its result can affect the form of other values. For example, consider the following code that defines a type for natural number parity. If the natural number is even, it can be represented as the sum of two equal natural numbers; if it is odd, it is equal to one plus the sum of two equal values. Pattern matching on a value of `Parity n` allows to discover if $n = j + j$ or $n = S(k + k)$, for some j and k in each branch of `with`. Note that the value of n is specialized accordingly, using information “learned” by the type-checker.

```
data Parity : N → Set where
  Even : ∀ {n : N} → Parity (n + n)
```

¹Also known as Curry-Howard “isomorphism” [31].

```

Odd : ∀ {n : ℕ} → Parity (S (n + n))
parity : (n : ℕ) → Parity n
parity = -- definition omitted
natToBin : ℕ → List Bool
natToBin zero = []
natToBin k with (parity k)
  natToBin (j + j) | Even = false :: natToBin j
  natToBin (suc (j + j)) | Odd = true :: natToBin j

```

For further information about Agda, see [26, 32].

3 REGULAR EXPRESSIONS

3.1 Standard Regular Expressions

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

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

where a is any symbol from the underlying alphabet. In our Agda formalization, we represent alphabet symbols using type `Char`.

Datatype `Regex` encodes RE syntax.

```

data Regex : Set where
  ∅ : Regex
  ε : Regex
  $ _ : Char → Regex
  _ • _ : Regex → Regex → Regex
  _ + _ : Regex → Regex → Regex
  _ ★ : Regex → Regex

```

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

The following datatype defines RE semantics inductively.

```

data _ ∈ [-] : List Char → Regex → Set where
  ε : [] ∈ [ε]
  $ _ : (c : Char) → [c] ∈ [$ c]
  _ • _ : xs ∈ [l] → ys ∈ [r] → (xs ++ ys) ∈ [l • r]
  _ + L_ : (r : Regex) → xs ∈ [l] → xs ∈ [l + r]
  _ + R_ : (l : Regex) → xs ∈ [r] → xs ∈ [l + r]
  _ ★ : xs ∈ [ε + (e • e ★)] → xs ∈ [e ★]

```

Constructor ϵ states that the empty string (denoted by the empty list $[]$) is in the language of RE ϵ .

For any single character a , the singleton string $[a]$ is in the RL for $\$ a$. Given membership proofs for REs l and r , $xs \in [l]$ and $ys \in [r]$, constructor \bullet can be used to build a proof for the concatenation of these REs. Constructor $+$ ($+$ L $+$ R) creates a membership proof for $l + r$ from a proof from l (r). Semantics for Kleene star is built using the following well known equivalence of REs: $e^* = \epsilon + e e^*$. Notice that we use the same constructor names both in RE syntax and in its semantics. Agda also allows overloading of data-constructors names.

Several inversion lemmas about RE parsing relation are necessary for derivative-based parsing formalization. They consist of pattern-matching on proofs of $_ \in [-]$ and are omitted for brevity.

One way to look at RE parsing is to interpret parse trees as terms whose type is a RE [13, 24]. To ensure the consistency between a tree and its RE type, we use an indexed data-type, an usual approach in dependently typed languages [26].

```

data Tree : Regex → Set where
  ε : Tree ε
  $ _ : (c : Char) → Tree ($ c)
  inl : ∀ (r : Regex) (tl : Tree l) → Tree (l + r)
  inr : ∀ (l : Regex) (tr : Tree r) → Tree (l + r)
  _ • _ : ∀ (tl : Tree l) (tr : Tree r) → Tree (l • r)
  star[] : Tree (l ★)
  star- :: : Tree l → Tree (l ★) → Tree (l ★)

```

Each constructor of `Tree` specifies how to build a parse tree for a RE. As an example, value `inl ε (($ 'a') • ($ 'b'))` denotes a parse tree for RE $ab + \epsilon$.

The relation between RE semantics and its parse trees are formalized by functions `flat` and `unflat`. The `flat` function builds a membership proof from a given parse tree by recursion on its structure. At each step, `flat` returns a pair formed by a string and its RE membership proof. Matched strings can be recovered from parse trees by concatenating values in their leaves (that would be the empty string or a single character).

```

flat : Tree e → ∃ (λ xs → xs ∈ [e])
flat ε = [], ε
flat ($ c) = [c], ($ c)
flat (inl r t) with flat t
...| xs, prf = _, r +L prf
flat (inr l t) with flat t
...| xs, prf = _, l +R prf
flat (t • t') with flat t | flat t'
...| xs, prf | ys, prf' = , (prf • prf')
flat star[] = [], (_ +L ε) ★
flat (star- :: t t') with flat t | flat t'
...| xs, prf | ys, prf' = , (_ +R (prf • prf')) ★

```

Note that our definition of `flat` ensures, by construction, that the string produced by a parse tree is in its RE language. Such property is stated in Nielsen et. al. paper as a theorem (cf. Theorem 2.1 [24]).

Function `unflat` builds parse trees from RE membership proofs straightforwardly:

```

unflat : xs ∈ [e] → Tree e
unflat ε = ε
unflat ($ c) = $ c
unflat (prf • prf') = unflat prf • unflat prf'
unflat (r +L prf) = inl r (unflat prf)
unflat (l +R prf) = inr l (unflat prf)
unflat ((_ +L ε) ★) = star[]
unflat ((_ +R (prf • prf')) ★) = star- :: (unflat prf) (unflat prf')

```

The next lemmas state that functions `flat` and `unflat` are inverses and they are proved by induction on t and $xs \in [e]$, respectively.

LEMMA 1. *Let e be a RE and $t : \text{Tree } e$ a parse tree. Then $\text{unflat } (\pi_2 \text{ (flat } t)) \equiv t$.*

LEMMA 2. Let xs be a string and e a RE s.t. $xs \in [e]$. Then $\text{flat}(\text{unflat } prf) \equiv (xs, prf)$.

3.2 Bit-codes for RE Parse Trees

We follow the encoding of Nielsen et. al. [24] that uses bit marks to register which branch was chosen in a parse tree for a union operator, $+$, and the beginning and end of matches that correspond to a Kleene star. Note that no marking is needed for concatenation and single character since coding and decoding of bit sequences are directed by the underlining RE.

We represent bit sequences using lists whose elements are of type `Bit`.

data `Bit` : Set where

`0b 1b` : Bit

Not every bit list corresponds to a valid parse tree. Data type `lsCode` represents an inductive predicate that holds when a bit list bs denotes a valid parse tree for some RE e .

data `_lsCode_` : List Bit \rightarrow Regex \rightarrow Set where

`ε` : [] `lsCode` ϵ
`$ _` : (c : Char) \rightarrow [] `lsCode` ($\$ c$)
`inl` : xs `lsCode` $l \rightarrow (0_b :: xs)$ `lsCode` ($l + r$)
`inr` : xs `lsCode` $r \rightarrow (1_b :: xs)$ `lsCode` ($l + r$)
`_ . _` : xs `lsCode` $l \rightarrow ys$ `lsCode` $r \rightarrow (xs ++ ys)$ `lsCode` ($l \bullet r$)
`star[]` : [`1b`] `lsCode` ($l \star$)
`star-` : xs `lsCode` $l \rightarrow xss$ `lsCode` ($l \star$) $\rightarrow (0_b :: xs ++ xss)$ `lsCode` ($l \star$)

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

Function `code` builds a bit representation for a parse tree together with a proof that the produced bit string is a valid tree representation.

`code` : Tree $e \rightarrow \exists (\lambda bs \rightarrow bs$ `lsCode` $e)$
`code` $\epsilon = []$, ϵ
`code` ($\$ c$) = [], ($\$ c$)
`code` (`inl` r t) **with** `code` t
 \dots | ys , $pr = 0_b :: ys$, `inl` r pr
`code` (`inr` l t) **with** `code` t
 \dots | ys , $pr = 1_b :: ys$, `inr` l pr
`code` ($t \bullet t'$) **with** `code` t | `code` t'
 \dots | xs , pr | ys , $pr' = xs ++ ys$, $pr \bullet pr'$
`code` `star[]` = `1b` :: [], `star[]`
`code` (`star-` : t ts) **with** `code` t | `code` ts
 \dots | xs , pr | xss , $prs = (0_b :: xs ++ xss)$, `star-` : pr prs

Next we present function `decode` which generates a parse tree from its correspondent bit code.

`decode` : $\exists (\lambda bs \rightarrow bs$ `lsCode` $e) \rightarrow$ Tree e
`decode` ($_, \epsilon$) = ϵ
`decode` ($_, (\$ c)$) = $\$ c$

`decode` ($_, (\text{inl } r \text{ } pr)$) = `inl` r (`decode` ($_, pr$)
`decode` ($_, (\text{inr } l \text{ } pr)$) = `inr` l (`decode` ($_, pr$)
`decode` ($_, (pr \bullet pr')$) **with** `decode` ($_, pr$) | `decode` ($_, pr')$
 \dots | $bs1$, $pr1$ | $bs2$, $pr2 = pr1 \bullet pr2$
`decode` `star[]` = ($_{-} +L \epsilon$) \star
`decode` (`star-` : pr pr') **with** `decode` ($_, pr$) | `decode` ($_, pr')$
 \dots | $pr1$ | $pr2 = (_ +R (pr1 \bullet pr2)) \star$

As one might expect, decoding a bit sequence produced by the `code` function will produce the original parse tree.

THEOREM 1. Let t : Tree e for some RE e . Then, it holds that `decode` (`code` t) $\equiv t$.

Building bit codes from parse trees does not present any advantage, since after building parse trees we will need to traverse them in order to generate codes. A better strategy is to incrementally build bit-codes while parsing a given RE. In next sections, we describe the formalization of this approach, due to Sulzmann et. al. [33].

3.3 Bit-annotated Regular Expressions

Intuitively, BRE just attach a list of bits bs to every non-empty regular expression. Type `BitRegex` defines the syntax of BRE, which just associates a list of bits to a BRE.

data `BitRegex` : Set where

`empty` : `BitRegex`
`eps` : List Bit \rightarrow `BitRegex`
`char` : List Bit \rightarrow (c : Char) \rightarrow `BitRegex`
`choice` : List Bit \rightarrow `BitRegex` \rightarrow `BitRegex` \rightarrow `BitRegex`
`cat` : List Bit \rightarrow `BitRegex` \rightarrow `BitRegex` \rightarrow `BitRegex`
`star` : List Bit \rightarrow `BitRegex` \rightarrow `BitRegex`

Conversion between `Regex` and `BitRegex` types is done by functions `internalize` and `erase`. First, we define an auxiliary function, `fuse`, which attaches a bit code to the top-most position of a BRE.

`fuse` : List Bit \rightarrow `BitRegex` \rightarrow `BitRegex`
`fuse` bs `empty` = `empty`
`fuse` bs (`eps` x) = `eps` ($bs ++ x$)
`fuse` bs (`char` x c) = `char` ($bs ++ x$) c
`fuse` bs (`choice` x e e') = `choice` ($bs ++ x$) e e'
`fuse` bs (`cat` x e e') = `cat` ($bs ++ x$) e e'
`fuse` bs (`star` x e) = `star` ($bs ++ x$) e

Function `internalize` converts a standard RE into BRE by inserting empty bit lists on it.

`internalize` : Regex \rightarrow `BitRegex`
`internalize` \emptyset = `empty`
`internalize` ϵ = `eps` []
`internalize` ($\$ x$) = `char` [] x
`internalize` ($e \bullet e'$) = `cat` [] (`internalize` e) (`internalize` e')
`internalize` ($e + e'$)
= `choice` [] (`fuse` [`0b`] (`internalize` e)
(`fuse` [`1b`] (`internalize` e'))
`internalize` ($e \star$) = `star` [] (`internalize` e)

Function `erase` is straightforwardly defined by recursion.

```

erase : BitRegex → Regex
erase empty = 0
erase (eps x) = ε
erase (char x c) = $ c
erase (choice x e e') = erase e + (erase e')
erase (cat x e e') = erase e • (erase e')
erase (star x e) = (erase e) ★
    
```

The relation between these functions is expressed by the following lemmas, which are proved by induction on the structure of e .

LEMMA 3. For all e and bs , $\text{erase (fuse } bs \text{ (internalize } e)) \equiv e$.

LEMMA 4. Let $e : \text{Regex}$. Then, $\text{erase (internalize } e) \equiv e$

Next, we define an alternative inductive semantics for BRE and show that it is sound and complete w.r.t. standard RE semantics. A value of type $xs \in \langle e \rangle$ states that string xs is in the language denoted by $\text{BitRegex } e$.

```

data _ ∈ ⟨_⟩ : List Char → BitRegex → Set where
  eps : (bs : List Bit) → [] ∈ ⟨ eps bs ⟩
  char : (bs : List Bit) (c : Char) → [ c ] ∈ ⟨ char bs c ⟩
  inl : (r : BitRegex) bs → xs ∈ ⟨ l ⟩ → xs ∈ ⟨ choice bs l r ⟩
  inr : (l : BitRegex) bs → xs ∈ ⟨ r ⟩ → xs ∈ ⟨ choice bs l r ⟩
  cat : (bs : List Bit) → xs ∈ ⟨ l ⟩ → ys ∈ ⟨ r ⟩ →
        (xs ++ ys) ∈ ⟨ cat bs l r ⟩
  star[] : (bs : List Bit) → [] ∈ ⟨ star bs l ⟩
  star- :: bs → (x :: xs) ∈ ⟨ l ⟩ → xss ∈ ⟨ star [] l ⟩ →
        (x :: xs ++ xss) ∈ ⟨ star bs l ⟩
    
```

The relation between BitRegex and Regex semantic is specified by the following theorem, proved by induction on the derivation of $xs \in [e]$.

THEOREM 2. Let $e : \text{Regex}$ and $xs : \text{List Char}$. Then, $xs \in [e]$, if and only if, $xs \in \langle \text{internalize } e \rangle$.

4 DERIVATIVES AND PARSING

Formally, the derivative of a formal language $L \subseteq \Sigma^*$ with respect to a symbol $a \in \Sigma$ is the language formed by suffixes of L words without the prefix a .

An algorithm for computing the derivative of a language represented by a RE as another RE is due to Brzozowski [6]. It relies on a function (called v) that determines if some RE accepts or not the empty string (by returning ϵ or 0 , respectively):

$$\begin{aligned}
 v(\emptyset) &= 0 \\
 v(\epsilon) &= \epsilon \\
 v(a) &= 0 \\
 v(e e') &= \begin{cases} \epsilon & \text{if } v(e) = v(e') = \epsilon \\ 0 & \text{otherwise} \end{cases} \\
 v(e + e') &= \begin{cases} \epsilon & \text{if } v(e) = \epsilon \text{ or } v(e') = \epsilon \\ 0 & \text{otherwise} \end{cases} \\
 v(e^*) &= \epsilon
 \end{aligned}$$

Decidability of $v(e)$ is proved by function $v[e]$, which is defined by induction over the structure of the input BRE e and returns a proof that the empty string is accepted or not, using Agda type of decidable propositions, $\text{Dec } P$.

```

v[_] : ∀ (e : BitRegex) → Dec ([ ] ∈ ⟨ e ⟩)
v[ empty ] = no (λ ())
v[ eps bs ] = yes eps
v[ char bs x ] = no (λ ())
v[ cat bs e e' ] with v[ e ] | v[ e' ]
v[ cat bs e e' ] | yes pr | (yes pr1)
  = yes (cat bs pr pr1)
v[ cat bs e e' ] | yes pr | (no ¬pr1)
  = no (¬pr1 ∘ π₂ ∘ [ ] ∈ • – invert)
v[ cat bs e e' ] | no ¬pr | pr1
  = no (¬pr ∘ π₁ ∘ [ ] ∈ • – invert)
v[ choice bs e e' ] with v[ e ] | v[ e' ]
v[ choice bs e e' ] | yes pr | pr1 = yes (e' +L pr)
v[ choice bs e e' ] | no ¬pr | (yes pr1) = yes (e +R pr1)
v[ choice bs e e' ] | no ¬pr | (no ¬pr1)
  = no ([ ¬pr , ¬pr1 ] ∘ e + – invert)
v[ star bs e ] = yes star[]
    
```

Function $v[e]$ is an immediate translation of $v(e)$, as defined by Brzozowski, to Agda code and it uses several inversion lemmas about BRE semantics. Lemma $[] \in \bullet - \text{invert}$ states that if the empty string is in the language of $\text{cat } bs \ l \ r$ (where l and r are arbitrary BRE's) then the empty string belongs to l and r 's languages. Lemma $e + - \text{invert}$ is defined similarly for choice.

4.1 Derivatives for Bit-annotated REs and its Properties

Following Sulzmann [33], we define a function mkEps that builds the bit code for a nullable BRE, (i.e. a BRE s.t. $[] \in \langle e \rangle$).

```

mkEps : [ ] ∈ ⟨ t ⟩ → List Bit
mkEps (eps bs) = bs
mkEps (inl br bs pr) = bs ++ mkEps pr
mkEps (inr bl bs pr) = bs ++ mkEps pr
mkEps (cat bs pr pr') = bs ++ mkEps pr ++ mkEps pr'
mkEps (star[] bs) = bs ++ [ 1b ]
mkEps (star- :: bs pr pr' x) = bs ++ [ 1b ]
    
```

Next we define the derivative operation on BREs in Agda. The difference between this definition and standard Brzozowski derivatives [6] is that the former inserts parse tree information in terms of bit annotations. For example, consider $\text{cat } bs \ l \ r$ where $[] \in \langle l \rangle$, additional parse information is built from a nullability test result using functions mkEps and fuse . For the Kleene star operation we record the start of a new iteration fusing $[0_b]$ and we mark a start of a new matching iteration by attaching the empty list in $\text{star } [] \ e$.

```

∂[_,_] : BitRegex → Char → BitRegex
∂[ empty , c ] = empty
∂[ eps bs , c ] = eps bs
∂[ char bs c , c' ] with c ≐ c'
...| yes refl = eps bs
...| no prf = empty
∂[ cat bs e e' , c ] with v[ e ]
∂[ cat bs e e' , c ] | yes pr
  = choice bs (cat bs ∂[ e , c ] e') (fuse (mkEps pr) ∂[ e' , c ])
    
```

$$\begin{aligned}
&\partial[\text{cat } bs \ e \ e', c] \mid \text{no } \neg pr = \text{cat } bs \ \partial[e, c] \ e' \\
&\partial[\text{choice } bs \ e \ e', c] = \text{choice } bs \ (\partial[e, c]) \ (\partial[e', c]) \\
&\partial[\text{star } bs \ e, c] \\
&\quad = \text{cat } bs \ (\text{fuse } [0_b] \ \partial[e, c]) \ (\text{star } [] \ e)
\end{aligned}$$

From this definition we prove the following important properties of the derivative operation: soundness of $\partial[-, -]$ ensures that if a string xs is in the language of $\partial[e, x]$, then $(x :: xs) \in \langle e \rangle$ holds. Completeness ensures that the other direction of the implication holds. Both are proved by induction on the structure of e .

THEOREM 3 (DERIVATIVE OPERATION SOUNDNESS). *For all BREs e , all strings xs and all symbols x , if $xs \in \langle \partial[e, x] \rangle$ holds then $(x :: xs) \in \langle e \rangle$ holds.*

THEOREM 4 (DERIVATIVE OPERATION COMPLETENESS). *For all BREs e , all strings xs and all symbols x , if $(x :: xs) \in \langle e \rangle$ holds then $xs \in \langle \partial[e, x] \rangle$ holds.*

4.2 Parsing

RE parsing involves determining which prefixes and substrings of the input string match a given RE. For this, we define datatypes that represent the fact that a given BRE matches a prefix or a substring of a given string.

We say that BRE e matches a prefix of string xs if there exist strings ys and zs such that $xs \equiv ys ++ zs$ and $ys \in \langle e \rangle$. Definition of `IsPrefix` datatype encode this concept. Datatype `IsSubstr` specifies when a BRE e matches a substring in xs : there must exist strings ys , zs and ws such that $xs \equiv ys ++ zs ++ ws$ and $zs \in \langle e \rangle$ hold.

data `IsPrefix` ($xs : \text{List Char}$) ($e : \text{BitRegex}$) : `Set` **where**
`Prefix` : $xs \equiv ys ++ zs \rightarrow ys \in \langle e \rangle \rightarrow \text{IsPrefix } xs \ e$

data `IsSubstr` ($xs : \text{List Char}$) ($e : \text{BitRegex}$) : `Set` **where**
`Substr` : $xs \equiv ys ++ zs ++ ws \rightarrow zs \in \langle e \rangle \rightarrow \text{IsSubstr } xs \ e$

Using these datatypes we can define the following relevant properties of prefixes and substrings that are used to fill proof obligations present in decidability tests. All these lemmas are direct consequences of prefix and substring definitions.

LEMMA 5 (LEMMA $\neg \text{IsPrefix}$). *For all BREs e , if $[] \in \langle e \rangle$ does not hold then neither does `IsPrefix [] e`.*

LEMMA 6 (LEMMA $\neg \text{IsPrefix-}$::). *For all BREs e and all strings xs , if $[] \in \langle e \rangle$ and `IsPrefix xs e` do not hold then neither does `IsPrefix (x :: xs) e`.*

LEMMA 7 (LEMMA $\neg \text{IsSubstr}$). *For all BREs e , if `IsPrefix [] e` does not hold then neither does `IsSubstr [] e`.*

LEMMA 8 (LEMMA $\neg \text{IsSubstr-}$::). *For all strings xs , all symbols x and all BREs e , if `IsPrefix (x :: xs) e` and `IsSubstr xs e` do not hold then neither does `IsSubstr (x :: xs) e`.*

Function `IsPrefixDec` decides if a given BRE e matches a prefix in xs by induction on the structure of xs , using Lemmas 5, 6, decidable emptiness test `v[-]` and Theorem 3. Intuitively, `IsPrefixDec` first checks if current RE e accepts the empty string. In this case, `[]` is returned as a prefix. Otherwise, it verifies, for each symbol x , whether BRE $\partial[e, x]$ matches a prefix of the input string. If

this is the case, a prefix including x is built from a recursive call to `IsPrefixDec` or if no prefix is matched a proof of such impossibility is constructed using lemma $\neg \text{IsPrefix-}$::.

`IsPrefixDec` : $\forall xs \ e \rightarrow \text{Dec } (\text{IsPrefix } xs \ e)$
`IsPrefixDec [] e with v[e]`
`IsPrefixDec [] e | yes p = yes (Prefix [] [] refl p)`
`IsPrefixDec [] e | no $\neg p$ = no ($\neg \text{IsPrefix } \neg p$)`
`IsPrefixDec (x :: xs) e with v[e]`
`IsPrefixDec (x :: xs) e | yes p = yes (Prefix [] (x :: xs) refl p)`
`IsPrefixDec (x :: xs) e | no $\neg p$ with IsPrefixDec xs ($\partial[e, x]$)`
`IsPrefixDec (x :: xs) e | no $\neg p$ | (yes (Prefix ys zs eq wit))`
`= yes (Prefix (x :: ys) zs (cong (- :: - x) eq) (∂ - sound wit))`
`IsPrefixDec (x :: xs) e | no $\neg p$ | (no $\neg p$)`
`= no ($\neg \text{IsPrefix-}$:: $\neg p$ $\neg p$)`

Function `IsSubstrDec` is also defined by induction on the structure of the input string e , using `IsPrefixDec` to check whether it is possible to match a prefix of e . In this case, a substring is built from this prefix. If there's no such prefix, a recursive call is made to check if there is a substring match, returning such substring or a proof that it does not exist.

`IsSubstrDec` : $\forall xs \ e \rightarrow \text{Dec } (\text{IsSubstr } xs \ e)$
`IsSubstrDec [] e with v[e]`
`IsSubstrDec [] e | yes p = yes (Substr [] [] [] refl p)`
`IsSubstrDec [] e | no $\neg p$ = no ($\neg \text{IsSubstr } (\neg \text{IsPrefix } \neg p)$)`
`IsSubstrDec (x :: xs) e with IsPrefixDec (x :: xs) e`
`IsSubstrDec (x :: xs) e | yes (Prefix ys zs eq wit)`
`= yes (Substr [] ys zs eq wit)`
`IsSubstrDec (x :: xs) e | no $\neg p$ with IsSubstrDec xs e`
`IsSubstrDec (x :: xs) e | no $\neg p$ | (yes (Substr ys zs ws eq wit))`
`= yes (Substr (x :: ys) zs ws (cong (- :: - x) eq) wit)`
`IsSubstrDec (x :: xs) e | no $\neg p_1$ | (no $\neg p$)`
`= no ($\neg \text{IsSubstr-}$:: $\neg p_1$ $\neg p$)`

5 IMPLEMENTATION DETAILS AND EXPERIMENTS

We include the formalized algorithm in a tool for RE parsing developed by us, named `verigrep`, in the style of GNU Grep [14]. We have built a simple parser combinator library for parsing RE syntax, using the Agda Standard Library and its support for calling Haskell functions through its foreign function interface.

Experimentation with our tool involved a comparison of its performance with GNU Grep [14] (`grep`), Google regular expression library `re2` [28] and Haskell RE parsing algorithms `haskell-regex`, described in [12]. The experiments consider three distinct algorithms implemented in `verigrep`: RE parsing using Brzozowski derivatives, Antimirov partial derivatives and bit-coded parsing as described in this work. We run RE parsing experiments on a machine with a Intel Core I7 1.7 GHz, 8GB RAM running Mac OS X 10.12.3; the results were collected and the median of several test runs was computed.

We use the same experiments as those used in [33]; these consist of parsing files containing thousands of occurrences of symbol a , using the RE $(a + b + ab)^*$; and parsing files containing thousands

of occurrences of ab , using the same RE. Results are presented in Figures 1 and 2, respectively.

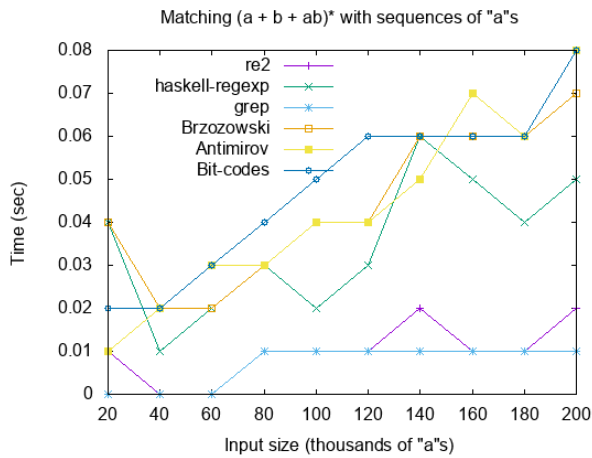


Figure 1: Results of experiment 1.

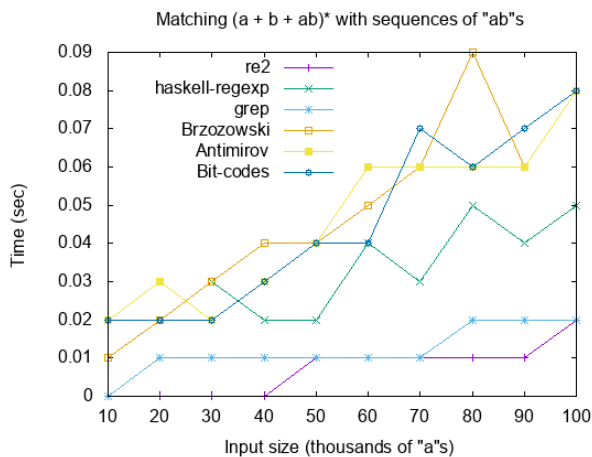


Figure 2: Results of experiment 2.

Our tool behaves poorly when compared with all other options considered. The cause of this inefficiency needs further investigation, since the algorithm formalized uses the POSIX disambiguation strategy, which avoids quotienting the result of derivative operations w.r.t. ACUI axioms as usual Brzozowski derivatives. The main reason behind POSIX and greedy disambiguation strategies in derivative based parsing is to improve efficiency by eliminating simplification steps on derivatives result [33]. We leave the proof that the formalized algorithm indeed produces POSIX parse trees for future work.

6 RELATED WORK

Parsing with derivatives. Recently, derivative-based parsing has received a lot of attention. Owens et al. were the first to present

a functional encoding of RE derivatives and use it to parsing and DFA building. They use derivatives to build scanner generators for ML and Scheme [27]; no formal proof of correctness was presented.

Might et al. [23] report on the use of derivatives for parsing not only RLs but also context-free ones. He uses derivatives to handle context-free grammars (CFG) and develops an equational theory for compaction that allows for efficient CFG parsing using derivatives. Implementation of derivatives for CFGs are described by using the Racket programming language [8]. However, Might et al. do not present formal proofs related to the use of derivatives for CFGs.

Fischer et al. describe an algorithm for RE-based parsing based on weighted automata in Haskell [12]. The paper describes the design evolution of such algorithm as a dialog between three persons. Their implementation has a competitive performance when compared with Google’s RE library [28]. This work also does not consider formal proofs of RE parsing.

An algorithm for POSIX RE parsing is described in [33]. 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.

Certified parsing algorithms. Certified algorithms for parsing also received attention recently. Firsov et al. describe a certified algorithm for RE parsing by converting an input RE to an equivalent NFA represented as a boolean matrix [9]. A matrix library based on some “block” operations [20] is developed and used Agda formalization of NFA-based parsing. Compared to our work, a NFA-based formalization requires a lot more infrastructure (such as a Matrix library). No experiments with the certified algorithm were reported.

Firsov describes an Agda formalization of a parsing algorithm that deals with any CFG (CYK algorithm) [11]. Bernardy et al. describe a formalization of another CFG parsing algorithm in Agda [4]: Valiant’s algorithm [34], which reduces CFG parsing to boolean matrix multiplication. In both works, no experiment with formalized parsing algorithms were reported.

A certified LR(1) CFG validator is described in [17]. The formalized checking procedure verifies if CFG and an automaton match. They proved soundness and completeness of the validator in the Coq proof assistant [5]. Termination of the LR(1) automaton interpreter is ensured by imposing a natural number bound on allowed recursive calls.

Formalization of a parser combinator library was the subject of Danielsson’s work [7]. He built a library of parser combinators using coinduction and provides correctness proofs of such combinators.

Almeida et al. [1] describe a Coq formalization of partial derivatives and its equivalence with automata. Partial derivatives were introduced by Antimirov [2] as an alternative to Brzozowski derivatives, since it avoids quotient resulting REs with respect to ACUI axioms. Almeida et al. motivation is to use such formalization as a basis for a decision procedure for RE equivalence.

Ridge [30] describes a formalization, in the HOL4 theorem prover, of a combinator parsing library. A parser generator for such combinators is described and a proof that generated parsers are sound and complete is presented. According to Ridge, preliminary results shows that parsers built using his generator are faster than those created by the Happy parser generator [15].

Ausaf et. al. [3] describe a formalization, in Isabelle/HOL [25], of the POSIX matching algorithm proposed by Sulzmann et al. [33]. They give a constructive characterization of what a POSIX matching is and prove that such matching is unique for a given RE and string. No experiments with the verified algorithm are reported.

7 CONCLUSION

We have given a complete formalization of a bit-coded derivative-based parsing for REs in Agda. To the best of our knowledge, this is the first work that presents a complete verification of a bit-code based parsing algorithm and uses it in a tool for RE-based search.

As future work, we intend to continue the development of verigrep by certifying greedy and POSIX disambiguation strategies and finite state machine based algorithms for parsing RE.

ACKNOWLEDGMENTS

We thank the anonymous reviewers for their comments. This work was conducted during a scholarship supported by PNPd at Universidade Federal de Pelotas, RS - Brazil. Financed by CAPES, Brazilian Federal Agency for Support and Evaluation of Graduate Education within the Ministry of Education of Brazil.

REFERENCES

- [1] José Bacelar Almeida, Nelma Moreira, David Pereira, and Simão Melo de Sousa. 2010. Partial Derivative Automata Formalized in Coq. In *Implementation and Application of Automata - 15th International Conference, CIAA 2010, Winnipeg, MB, Canada, August 12-15, 2010. Revised Selected Papers (Lecture Notes in Computer Science)*, Michael Domaratzki and Kai Salomaa (Eds.), Vol. 6482. Springer, 59–68. DOI: http://dx.doi.org/10.1007/978-3-642-18098-9_7
- [2] Valentin Antimirov. 1996. Partial derivatives of regular expressions and finite automaton constructions. *Theoretical Computer Science* 155, 2 (1996), 291 – 319. DOI: [http://dx.doi.org/10.1016/0304-3975\(95\)00182-4](http://dx.doi.org/10.1016/0304-3975(95)00182-4)
- [3] Fahad Ausaf, Roy Dyckhoff, and Christian Urban. 2016. POSIX Lexing with Derivatives of Regular Expressions (Proof Pearl). In *Interactive Theorem Proving - 7th International Conference, ITP 2016, Nancy, France, August 22-25, 2016, Proceedings (Lecture Notes in Computer Science)*, Jasmin Christian Blanchette and Stephan Merz (Eds.), Vol. 9807. Springer, 69–86. DOI: http://dx.doi.org/10.1007/978-3-319-43144-4_5
- [4] Jean-Philippe Bernardy and Patrik Jansson. 2016. Certified Context-Free Parsing: A formalisation of Valiant’s Algorithm in Agda. *CoRR abs/1601.07724* (2016). <http://arxiv.org/abs/1601.07724>
- [5] 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.
- [6] Janusz A. Brzozowski. 1964. Derivatives of Regular Expressions. *J. ACM* 11, 4 (Oct. 1964), 481–494. DOI: <http://dx.doi.org/10.1145/321239.321249>
- [7] Nils Anders Danielsson. 2010. Total Parser Combinators. *SIGPLAN Not.* 45, 9 (Sept. 2010), 285–296. DOI: <http://dx.doi.org/10.1145/1932681.1863585>
- [8] Matthias Felleisen, M.D. Barski Conrad, David Van Horn, and Eight Students of Northeastern University. 2013. *Realm of Racket: Learn to Program, One Game at a Time!* No Starch Press, San Francisco, CA, USA.
- [9] Denis Firsov and Tarmo Uustalu. 2013. Certified Parsing of Regular Languages. In *Certified Programs and Proofs - Third International Conference, CPP 2013, Melbourne, VIC, Australia, December 11-13, 2013, Proceedings (Lecture Notes in Computer Science)*, Georges Gonthier and Michael Norrish (Eds.), Vol. 8307. Springer, 98–113. DOI: http://dx.doi.org/10.1007/978-3-319-03545-1_7
- [10] Denis Firsov and Tarmo Uustalu. 2014. Certified CYK parsing of context-free languages. *J. Log. Algebr. Meth. Program.* 83, 5-6 (2014), 459–468. DOI: <http://dx.doi.org/10.1016/j.jlmp.2014.09.002>
- [11] Denis Firsov and Tarmo Uustalu. 2014. Certified (CYK) parsing of context-free languages. *Journal of Logical and Algebraic Methods in Programming* 83, 5-6 (2014), 459 – 468. DOI: <http://dx.doi.org/10.1016/j.jlmp.2014.09.002>
- [12] Sebastian Fischer, Frank Huch, and Thomas Wilke. 2010. A Play on Regular Expressions: Functional Pearl. In *Proceedings of the 15th ACM SIGPLAN International Conference on Functional Programming (ICFP ’10)*. ACM, New York, NY, USA, 357–368. DOI: <http://dx.doi.org/10.1145/1863543.1863594>
- [13] Alain Frisch and Luca Cardelli. 2004. Greedy Regular Expression Matching. In *Automata, Languages and Programming: 31st International Colloquium, ICALP 2004, Turku, Finland, July 12-16, 2004. Proceedings (Lecture Notes in Computer Science)*, Josep Diaz, Juhani Karhumäki, Arto Lepistö, and Donald Sannella (Eds.), Vol. 3142. Springer, 618–629. DOI: http://dx.doi.org/10.1007/978-3-540-27836-8_53
- [14] Grep 2017. GNU Grep home page. <https://www.gnu.org/software/grep/>. (2017).
- [15] Happy 2001. Happy: The parser generator for Haskell. <http://www.haskell.org/happy/>. (2001).
- [16] John E. Hopcroft, Rajeev Motwani, Rotwani, and Jeffrey D. Ullman. 2000. *Introduction to Automata Theory, Languages and Computability* (2nd ed.). Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA.
- [17] Jacques-Henri Jourdan, François Pottier, and Xavier Leroy. 2012. Validating LR(1) Parsers. In *Proceedings of the 21st European Conference on Programming Languages and Systems (ESOP’12)*. Springer-Verlag, Berlin, Heidelberg, 397–416. DOI: http://dx.doi.org/10.1007/978-3-642-28869-2_20
- [18] M. E. Lesk and E. Schmidt. 1990. UNIX Vol. II. W. B. Saunders Company, Philadelphia, PA, USA, Chapter Lex&Mdash; Lexical Analyzer Generator, 375–387. <http://dl.acm.org/citation.cfm?id=107172.107193>
- [19] Raul Lopes, Rodrigo Ribeiro, and Carlos Camarão. 2015. Certified Derivative-Based Parsing of Regular Expressions. In *Programming Languages – Lecture Notes in Computer Science 9889*. Springer, 95–109.
- [20] Hugo Daniel Macedo and José Nuno Oliveira. 2013. Typing linear algebra: A biproduct-oriented approach. *CoRR abs/1312.4818* (2013). <http://arxiv.org/abs/1312.4818>
- [21] Per Martin-Löf. 1998. An intuitionistic theory of types. In *Twenty-five years of constructive type theory (Venice, 1995)*. Oxford Logic Guides, Vol. 36. Oxford Univ. Press, New York, 127–172.
- [22] Conor McBride and James McKinna. 2004. The View from the Left. *J. Funct. Program.* 14, 1 (Jan. 2004), 69–111. DOI: <http://dx.doi.org/10.1017/S0956796803004829>
- [23] Matthew Might, David Darais, and Daniel Spiewak. 2011. Parsing with Derivatives: A Functional Pearl. *SIGPLAN Not.* 46, 9 (Sept. 2011), 189–195. DOI: <http://dx.doi.org/10.1145/2034574.2034801>
- [24] Lasse Nielsen and Fritz Henglein. 2011. *Bit-coded Regular Expression Parsing*. Springer Berlin Heidelberg, Berlin, Heidelberg, 402–413. DOI: http://dx.doi.org/10.1007/978-3-642-21254-3_32
- [25] Tobias Nipkow, Markus Wenzel, and Lawrence C. Paulson. 2002. *Isabelle/HOL: A Proof Assistant for Higher-order Logic*. Springer-Verlag, Berlin, Heidelberg.
- [26] Ulf Norell. 2009. Independently Typed Programming in Agda. In *Proceedings of the 4th International Workshop on Types in Language Design and Implementation (TLDI’09)*. ACM, New York, NY, USA, 1–2. DOI: <http://dx.doi.org/10.1145/1481861.1481862>
- [27] Scott Owens, John Reppy, and Aaron Turon. 2009. Regular-expression Derivatives Re-examined. *J. Funct. Program.* 19, 2 (March 2009), 173–190. DOI: <http://dx.doi.org/10.1017/S0956796808007090>
- [28] re2 2016. Google Regular Expression Library - re2. <https://github.com/google/re2>. (2016).
- [29] Rodrigo Ribeiro, Raul Lopes, and Carlos Camarão. 2017. Certified Derivative Based Parsing of Regular Expressions – On-line repository. <https://github.com/rodrigoriibeiro/regex>. (2017).
- [30] Tom Ridge. 2011. Simple, Functional, Sound and Complete Parsing for All Context-free Grammars. In *Proceedings of the First International Conference on Certified Programs and Proofs (CPP’11)*. Springer-Verlag, Berlin, Heidelberg, 103–118. DOI: http://dx.doi.org/10.1007/978-3-642-25379-9_10
- [31] Morten Heine Sørensen and Pawel Urzyczyn. 2006. *Lectures on the Curry-Howard Isomorphism, Volume 149 (Studies in Logic and the Foundations of Mathematics)*. Elsevier Science Inc., New York, NY, USA.
- [32] Aaron Stump. 2016. *Verified Functional Programming in Agda*. Association for Computing Machinery and Morgan; Claypool, New York, NY, USA.
- [33] 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. DOI: http://dx.doi.org/10.1007/978-3-319-07151-0_13
- [34] Leslie G. Valiant. 1975. General Context-free Recognition in Less Than Cubic Time. *J. Comput. Syst. Sci.* 10, 2 (April 1975), 308–315. DOI: [http://dx.doi.org/10.1016/S0022-0000\(75\)80046-8](http://dx.doi.org/10.1016/S0022-0000(75)80046-8)