

Certified Bit-Coded Regular Expression Parsing

Rodrigo Ribeiro¹ André Rauber Du Bois²

¹Departament of Computer Science
Universidade Federal de Ouro Preto

²Departament of Computer Science
Universidade Federal de Pelotas

September 21, 2017

Introduction

- ▶ Parsing is pervasive in computing
 - ▶ String search tools, lexical analysers...
 - ▶ Binary data files like images, videos ...
- ▶ Our focus: Regular Languages (RLs)
 - ▶ Languages denoted by Regular Expressions (REs) and equivalent formalisms

Introduction

- ▶ Is RE parsing a yes / no question?
 - ▶ No! Better to produce evidence: parse trees.
- ▶ Why use bit-codes for parse trees?
 - ▶ Memory compact representation of parsing evidence.
 - ▶ Easy serialization of parsing results.

Contributions

- ▶ We provide fully certified proofs of a derivative based algorithm that produces a bit representation of a parse tree.
- ▶ We mechanize results about the relation between RE parse trees and bit-codes.
- ▶ We provide sound and complete decision procedures for prefix and substring matching of RE.
- ▶ Coded included in a RE search tool developed by us — verigrep.
- ▶ All results formalized in Agda version 2.5.2.

Regular Expression Syntax

- Definition of RE over a finite alphabet Σ .

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

- Agda code

data `Regex` : `Set` **where**

`∅` : `Regex`

`ε` : `Regex`

`$ _` : `Char` → `Regex`

`_ • _` : `Regex` → `Regex` → `Regex`

`_ + _` : `Regex` → `Regex` → `Regex`

`_ ★` : `Regex` → `Regex`

Regular Expression Semantics

$$\frac{}{\epsilon \in \llbracket \epsilon \rrbracket}$$

$$\frac{a \in \Sigma}{a \in \llbracket a \rrbracket}$$

$$\frac{s \in \llbracket e \rrbracket \quad s' \in \llbracket e' \rrbracket}{ss' \in \llbracket ee' \rrbracket}$$

$$\frac{s \in \llbracket e \rrbracket}{s \in \llbracket e + e' \rrbracket}$$

$$\frac{s' \in \llbracket e' \rrbracket}{s' \in \llbracket e + e' \rrbracket}$$

$$\frac{s \in \llbracket \epsilon + ee^* \rrbracket}{s \in \llbracket e^* \rrbracket}$$

Regular Expression Semantics — Agda code

```
data _ ∈  $\llbracket$ _  $\rrbracket$  : List Char → Regex → Set where  
   $\epsilon$       :  $\llbracket \rrbracket \in \llbracket \epsilon \rrbracket$   
   $\$$ _      : ( $c$  : Char) →  $\$ c \in \llbracket \$ c \rrbracket$   
  _ • _   :  $s \in \llbracket l \rrbracket \rightarrow$   
             $s' \in \llbracket r \rrbracket \rightarrow$   
             $(s ++ s') \in \llbracket l \bullet r \rrbracket$   
  _ + L _ :  $s \in \llbracket l \rrbracket \rightarrow s \in \llbracket l + r \rrbracket$   
  -- some code omitted...
```

Parse trees for REs

- ▶ We interpret RE as types and parse tree as terms.
- ▶ Informally:
 - ▶ leafs: empty string and character.
 - ▶ concatenation: pair of parse trees.
 - ▶ choice: just the branch of chosen RE.
 - ▶ Kleene star: list of parse trees.

Parse trees for RE — Example

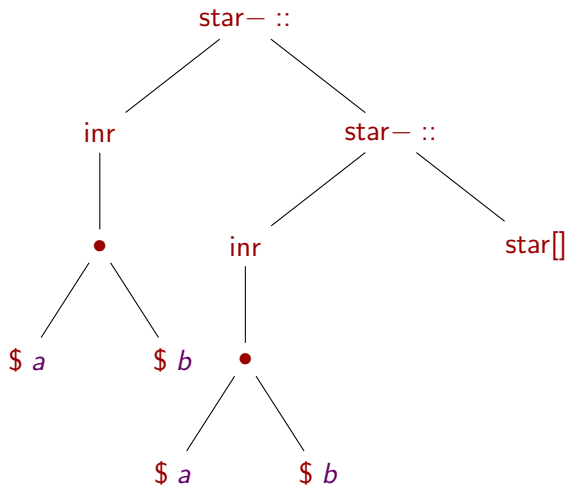


Figure: Parse tree for RE: $(c + ab)^*$ and the string $w = abab$.

Parse trees for REs

```
data Tree : Regex → Set where  
   $\epsilon$  : Tree  $\epsilon$   
   $\$_{\_}$  : ( $c$  : Char) → Tree ( $\$ c$ )  
  inl : Tree  $l$  → Tree ( $l + r$ )  
  inr : Tree  $r$  → Tree ( $l + r$ )  
   $\_ \bullet \_$  : Tree  $l$  → Tree  $r$  → Tree ( $l \bullet r$ )  
  star[] : Tree ( $l \star$ )  
  star- :: : Tree  $l$  → Tree ( $l \star$ ) → Tree ( $l \star$ )
```

Relating parse trees and RE semantics

- ▶ Using function *flat*.
- ▶ Property: Let t be a parse tree for a RE e and a string s . Then, $flat(t) = s$ and $s \in \llbracket e \rrbracket$.

$$\begin{aligned} flat(\epsilon) &= [] \\ flat(\$ a) &= [a] \\ flat(inl\ t) &= flat(t) \\ flat(inr\ t) &= flat(t) \\ flat(t \bullet t') &= flat(t) ++ flat(t') \\ flat(star[]) &= [] \\ flat(star- :: t\ ts) &= flat(t) ++ flat(ts) \end{aligned}$$

Relating parse trees and RE semantics

- flat type ensure its correctness property!

```
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')
-- some code omitted
```

Bit-codes for parse trees

- ▶ Bit-codes mark...
 - ▶ which branch of choice was chosen during parsing.
 - ▶ matchings done by the Kleene star operator.
- ▶ Predicate relating bit-codes to its RE.

data `_IsCode_` : `List Bit` \rightarrow `Regex` \rightarrow `Set` **where**

ϵ : `[] IsCode` ϵ

$\$_$: $(c : \text{Char}) \rightarrow [] \text{IsCode } (\$ c)$

`inl` : $xs \text{ IsCode } l \rightarrow (0_b :: xs) \text{ IsCode } (l + r)$

`inr` : $xs \text{ IsCode } r \rightarrow (1_b :: xs) \text{ IsCode } (l + r)$

$_ \bullet _$: $xs \text{ IsCode } l \rightarrow ys \text{ IsCode } r \rightarrow (xs ++ ys) \text{ IsCode } (l \bullet r)$

`star[]` : $[1_b] \text{ IsCode } (l \star)$

`star-` :: $xs \text{ IsCode } l \rightarrow xss \text{ IsCode } (l \star) \rightarrow$
 $(0_b :: xs ++ xss) \text{ IsCode } (l \star)$

How to relate bit-codes and parse trees?

- Function `code` builds bit-codes for parse trees.

```
code : Tree e → ∃ (λ bs → bs lsCode e)
code ($ c) = [], ($ c)
code (inl r t) with code t
...| ys , pr = 0b :: ys , inl r pr
code (inr l t) with code t
...| ys , pr = 1b :: ys , inr l pr
code star[] = 1b :: [], star[]
code (star- :: t ts) with code t | code ts
...| xs , pr | xss , prs = (0b :: xs ++ xss) , star- :: pr prs
-- some code omitted
```

How to relate bit-codes and parse-trees?

- ▶ Function `decode` parses a bit string for w.r.t. a RE.
- ▶ Property: forall t , $\text{decode}(\text{code } t) \equiv t$

```
decode :  $\exists (\lambda bs \rightarrow bs \text{ isCode } e) \rightarrow \text{Tree } e$ 
decode (–, ($ c)) = $ c
decode (–, (inl r pr)) = inl r (decode (–, pr))
decode (–, (inr l pr)) = inr l (decode (–, pr))
decode star[] = (– +L  $\epsilon$ )  $\star$ 
decode (star– :: pr pr') with decode (–, pr) | decode (–, pr')
...| pr1 | pr2 = (– +R (pr1 • pr2))  $\star$ 
-- some code omitted
```

Bit-codes for RE parse trees

- ▶ Building a parse tree for compute bit-codes is expensive.
- ▶ Better idea: build bit-codes during parsing, instead of computing parse trees.
- ▶ How? Just attach bit-codes to RE.

data BitRegex : Set **where**

empty : BitRegex

eps : List Bit \rightarrow BitRegex

char : List Bit \rightarrow 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

Relating REs and BREs

- Function `internalize` converts a RE into a BRE.

`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)

Relating REs and BREs

- ▶ Function `erase` converts a BRE into a RE.
- ▶ Property: for all e , $\text{erase}(\text{internalize } e) \equiv e$.

`erase` : `BitRegex` \rightarrow `Regex`

`erase empty` = \emptyset

`erase (eps x)` = ϵ

`erase (char x c)` = $\$ c$

`erase (choice x e e')` = $\text{erase } e + (\text{erase } e')$

`erase (cat x e e')` = $\text{erase } e \bullet (\text{erase } e')$

`erase (star x e)` = $(\text{erase } e) \star$

Semantics of BREs

- ▶ Same as RE semantics, but includes the bit-codes.
- ▶ Property: $s \in \llbracket e \rrbracket$ iff $s \in \langle \text{internalize } e \rangle$.

```
data  $_ \in \langle \_ \rangle : \text{List Char} \rightarrow \text{BitRegex} \rightarrow \text{Set where}$   
   $\text{eps} : [] \in \langle \text{eps } bs \rangle$   
   $\text{char} : (c : \text{Char}) \rightarrow [c] \in \langle \text{char } bs\ c \rangle$   
   $\text{inl} : s \in \langle l \rangle \rightarrow s \in \langle \text{choice } bs\ l\ r \rangle$   
   $\text{inr} : s \in \langle r \rangle \rightarrow s \in \langle \text{choice } bs\ l\ r \rangle$   
   $\text{cat} : s \in \langle l \rangle \rightarrow s' \in \langle r \rangle \rightarrow (s\ ++\ s') \in \langle \text{cat } bs\ l\ r \rangle$   
   $-- \text{ some code omitted}$ 
```

Derivatives in a nutshell

- ▶ What is the derivative of an (B)RE?
 - ▶ Derivatives are defined w.r.t. an alphabet symbol.
- ▶ The derivative of a (B)RE e w.r.t. a , $\partial(e, a)$, is another (B)RE that denotes all strings in e language with the leading a removed.

$$\partial_a(e) = \{s \mid as \in \llbracket e \rrbracket\}$$

- ▶ Property: $s \in \langle \partial[e, x] \rangle$ holds iff $(x :: s) \in \langle e \rangle$ holds.

Derivatives for BREs

- Follows the definition of Brzozowski's derivatives.

```
 $\partial[-, -] : \text{BitRegex} \rightarrow \text{Char} \rightarrow \text{BitRegex}$   
 $\partial[\text{eps } bs, c] = \text{eps } bs$   
 $\partial[\text{cat } bs \ e \ e', c] \text{ with } \nu[e]$   
 $\partial[\text{cat } bs \ e \ e', c] \mid \text{yes } pr$   
    = choice  $bs$  (cat  $bs$   $\partial[e, c] \ e'$ ) (fuse (mkEps  $pr$ )  $\partial[e', c]$ )  
 $\partial[\text{cat } bs \ e \ e', c] \mid \text{no } pr = \text{cat } bs \ \partial[e, c] \ e'$   
 $\partial[\text{star } bs \ e, c]$   
    = cat  $bs$  (fuse [0b]  $\partial[e, c]$ ) (star []  $e$ )  
-- some code omitted
```

Nullability test for BREs

- ▶ Checks if $\epsilon \in \llbracket e \rrbracket$, for some (B)RE e .
- ▶ Agda code: decision procedure for $[] \in \langle e \rangle$.

$$\nu(\emptyset) = \emptyset$$

$$\nu(\epsilon) = \epsilon$$

$$\nu(a) = \emptyset$$

$$\nu(e \ e') = \begin{cases} \epsilon & \text{if } \nu(e) = \nu(e') = \epsilon \\ \emptyset & \text{otherwise} \end{cases}$$

$$\nu(e + e') = \begin{cases} \epsilon & \text{if } \nu(e) = \epsilon \text{ or } \nu(e') = \epsilon \\ \emptyset & \text{otherwise} \end{cases}$$

$$\nu(e^*) = \epsilon$$

Parsing with derivatives

- ▶ RE-based text search tools parse prefixes and substrings.
- ▶ Types `IsPrefix` and `IsSubstr` are proofs that a string is a prefix / substring of an input BRE.
- ▶ Parsing algorithms defined as proofs of decidability of `IsPrefix` and `IsSubstr`. Proof by induction on the input string, using properties of derivative operation.

```
data IsPrefix (a : List Char) (e : BitRegex) : Set where
```

```
  Prefix : a ≡ b ++ c → b ∈ ⟨ e ⟩ → IsPrefix a e
```

```
data IsSubstr (a : List Char) (e : BitRegex) : Set where
```

```
  Substr : a ≡ b ++ c ++ d → c ∈ ⟨ e ⟩ → IsSubstr a e
```

Experimental Results

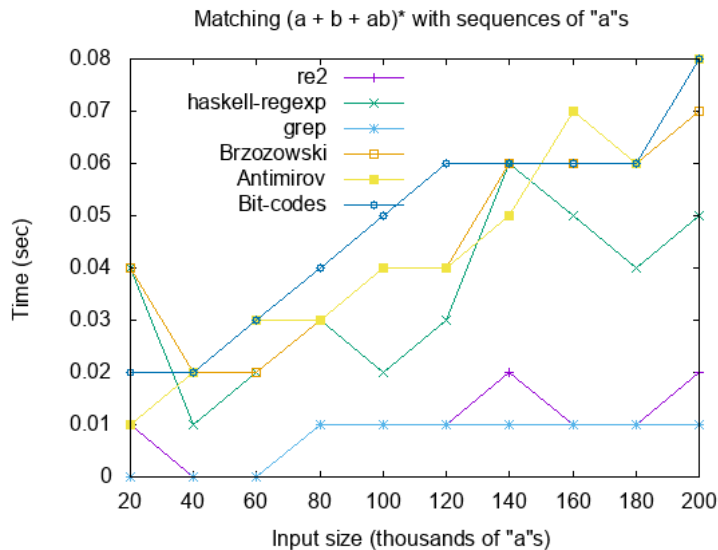


Figure: Results of experiment 1.

Future Work

- ▶ How to improve efficiency?
 - ▶ How intrinsic verification affects generated code efficiency?
 - ▶ Currently porting code to use extrinsic verification (proofs separated from program code).
 - ▶ Experiment with alternative formalization: BRE semantics defined by **erase**: $s \in \langle e \rangle = s \in \llbracket \text{erase } e \rrbracket$.
- ▶ How to measure memory consumption, without compiler support?
 - ▶ No profiling support in Agda compiler.
 - ▶ Agda compiles to Haskell, but there's no direct correspondence between Agda source code and Haskell generated code.

Conclusion

- ▶ We build a certified algorithm for BRE parsing in Agda.
- ▶ We certify several previous results about bit-coded parse trees and their relationship with RE semantics.
- ▶ Algorithm included in verigrep tool for RE text search.

Relating REs and BREs

- Function `fuse` attach a bit-string into a BRE.

```
fuse : List Bit → BitRegex → 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
```

Building bit-codes for ϵ

```
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' eq) = bs ++ mkEps pr ++ mkEps pr'
mkEps (star[] bs) = bs ++ [ 1b ]
mkEps (star- :: bs pr pr' x) = bs ++ [ 1b ]
```

Relating parse trees and RE semantics

- ▶ `unflat` builds parse trees out of RE semantics evidence.
- ▶ Functions `flat` and `unflat` are inverses.

```
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)
-- some code omitted
```