

Let's play with Regular Expressions

Sergey Vinokurov

serg.foo@gmail.com, @5ergv

2016-11-26

Outline

Introduction

Regular Expressions refresher

Glushkov Automaton Construction

Laziness allows matching strictly more expressive languages

Introduction

- ▶ This talk is based on “A Play on Regular Expressions”
Functional Pearl by Sebastian Fischer, Frank Huch and
Thomas Wilke

Introduction

- ▶ This talk is based on “A Play on Regular Expressions” Functional Pearl by Sebastian Fischer, Frank Huch and Thomas Wilke
- ▶ Regular expressions are usually taken for granted and are implemented in C

Why reimplement regular expressions in Haskell

- ▶ Some widespread implementations, e.g. Perl, Python, Java take exponential time on some inputs - we can do better

Why reimplement regular expressions in Haskell

- ▶ Some widespread implementations, e.g. Perl, Python, Java take exponential time on some inputs - we can do better
 - ▶ Not all implementations - grep and tcl, among others, are OK

Why reimplement regular expressions in Haskell

- ▶ Some widespread implementations, e.g. Perl, Python, Java take exponential time on some inputs - we can do better
 - ▶ Not all implementations - grep and tcl, among others, are OK
- ▶ Foreign libraries are harder to integrate with pure code

Why reimplement regular expressions in Haskell

- ▶ Some widespread implementations, e.g. Perl, Python, Java take exponential time on some inputs - we can do better
 - ▶ Not all implementations - grep and tcl, among others, are OK
- ▶ Foreign libraries are harder to integrate with pure code
 - ▶ Have a trie data structure, want all entries where prefix matches a regex. Can do this naively but will have to rerun matching process for each prefix

Why reimplement regular expressions in Haskell

- ▶ Some widespread implementations, e.g. Perl, Python, Java take exponential time on some inputs - we can do better
 - ▶ Not all implementations - grep and tcl, among others, are OK
- ▶ Foreign libraries are harder to integrate with pure code
 - ▶ Have a trie data structure, want all entries where prefix matches a regex. Can do this naively but will have to rerun matching process for each prefix
 - ▶ Compiling via GHCJS and can't link to C libraries

Why reimplement regular expressions in Haskell

- ▶ Some widespread implementations, e.g. Perl, Python, Java take exponential time on some inputs - we can do better
 - ▶ Not all implementations - grep and tcl, among others, are OK
- ▶ Foreign libraries are harder to integrate with pure code
 - ▶ Have a trie data structure, want all entries where prefix matches a regex. Can do this naively but will have to rerun matching process for each prefix
 - ▶ Compiling via GHCJS and can't link to C libraries
- ▶ Can do symbolic manipulation on regular expression AST - just for fun

Regular Expressions refresher

Definitions

Regular Expression are an expression language for describing sets of strings.

Regular Expressions refresher

Definitions

Regular Expression are an expression language for describing sets of strings.

Strings - sequences of characters. a, b, c - characters, $aaaab, c$ - strings

Regular Expressions refresher

Definitions

Regular Expression are an expression language for describing sets of strings.

Strings - sequences of characters. a, b, c - characters, $aaaab, c$ - strings

Symbol ε denotes empty string

Regular Expressions refresher

Definitions

Regular Expression are an expression language for describing sets of strings.

Strings - sequences of characters. a, b, c - characters, $aaaab, c$ - strings

Symbol ε denotes empty string

String concatenation is denoted with $\mathbin{\text{++}}$: $aaaaaa \mathbin{\text{++}} b = aaaaaab$

$$\forall x : x \mathbin{\text{++}} \varepsilon = \varepsilon \mathbin{\text{++}} x = x$$

Regular Expressions refresher

Definitions

Regular Expression are an expression language for describing sets of strings.

Strings - sequences of characters. a, b, c - characters, $aaaab, c$ - strings

Symbol ε denotes empty string

String concatenation is denoted with $\#$: $aaaaaa \# b = aaaaaab$

$$\forall x : x \# \varepsilon = \varepsilon \# x = x$$

Alphabet - set of all characters, denoted Σ . E.g. $\Sigma = \{a, b, c\}$.

Regular Expressions refresher

Introducing Haskell ADT for regexps

Regular expressions are defined inductively

```
data Sigma = A | B | C
```

```
  deriving (Show, Eq, Ord, Enum, Bounded)
```

```
type Str a = [a]
```

```
data Regex a
```

```
  = Eps                -- Empty regex
```

```
  | Sym a              -- Singleton regex
```

```
  | Seq (Regex a) (Regex a) -- Sequence
```

```
  | Alt (Regex a) (Regex a) -- Alternatives
```

```
  | Rep (Regex a)      -- Repetition
```

```
  deriving (Show, Eq)
```


Regular Expressions refresher

Example

Consider regular expression $a^* \cdot b \cdot (c \mid \varepsilon)$

Regular Expressions refresher

Example

Consider regular expression $a^* \cdot b \cdot (c \mid \varepsilon)$

$$L(a^* \cdot b \cdot \varepsilon \mid c) = \{b, bc, ab, abc, aab, aabc, aaab, aaabc, \dots\}$$

In Haskell,

$re = Seq (Rep (Sym A)) (Seq (Sym B) (Alt Eps (Sym C)))$

Regular Expressions refresher

Example

Consider regular expression $a^* \cdot b \cdot (c \mid \varepsilon)$

$$L(a^* \cdot b \cdot \varepsilon \mid c) = \{b, bc, ab, abc, aab, aabc, aaab, aaabc, \dots\}$$

In Haskell,

$re = Seq (Rep (Sym A)) (Seq (Sym B) (Alt Eps (Sym C)))$

Or, with less parens

$re = Rep (Sym A) 'Seq' Sym B 'Seq' Alt Eps (Sym C)$

Regular Expressions refresher

Example

Consider regular expression $a^* \cdot b \cdot (c \mid \varepsilon)$

$$L(a^* \cdot b \cdot \varepsilon \mid c) = \{b, bc, ab, abc, aab, aabc, aaab, aaabc, \dots\}$$

In Haskell,

```
re = Seq (Rep (Sym A)) (Seq (Sym B) (Alt Eps (Sym C)))
```

Or, with less parens

```
re = Rep (Sym A) 'Seq' Sym B 'Seq' Alt Eps (Sym C)
```

Expected output

```
accept re [B] ⇒ True
```

```
accept re [B, C] ⇒ True
```

```
accept re [A, A, A, B] ⇒ True
```

```
accept re [A, A, A, C] ⇒ False
```

```
accept re [A, A, A, B, C] ⇒ True
```

Regular Expressions refresher

Defining regexp semantics

For any regex (A), notation $L(A)$ denotes all strings that regex matches.

Let's define regexp semantics in Haskell by defining regexp matching function:

$$\text{accept} :: \text{Eq } a \Rightarrow \text{Regex } a \rightarrow \text{Str } a \rightarrow \text{Bool}$$
$$\text{accept } \text{Eps} \quad s = \text{acceptEps } s$$
$$\text{accept } (\text{Sym } c) \quad s = \text{acceptSym } c \ s$$
$$\text{accept } (\text{Seq } x \ y) \quad s = \text{acceptSeq } x \ y \ s$$
$$\text{accept } (\text{Alt } x \ y) \quad s = \text{acceptAlt } x \ y \ s$$
$$\text{accept } (\text{Rep } x) \quad s = \text{acceptRep } x \ s$$

Regular Expressions refresher

Empty string regex

Empty string - ϵ is a regular expression that matches the empty string.

$$L(\epsilon) = \{\epsilon\}$$

acceptEps :: Str a → Bool

acceptEps s = null s

Regular Expressions refresher

Singleton character regex

Atoms - all symbols from Σ are regular expressions that match themselves.

$$\forall x \in \Sigma : L(x) = \{x\}$$

Haskell semantics:

```
acceptSym :: Eq a  $\Rightarrow$  a  $\rightarrow$  Str a  $\rightarrow$  Bool  
acceptSym c s = s == [c]
```

Regular Expressions refresher

Concatenation

Concatenation - if A and B are regular expressions, so is $A \cdot B$

$$L(A \cdot B) = \{x \# y \mid x \in L(A), y \in L(B)\}$$

Haskell semantics:

```
acceptSeq :: Eq a  $\Rightarrow$  Regex a  $\rightarrow$  Regex a  $\rightarrow$  Str a  $\rightarrow$  Bool  
acceptSeq r1 r2 s =  
  or [ accept r1 s1 && accept r2 s2 | (s1, s2)  $\leftarrow$  split s ]
```

```
split :: Str a  $\rightarrow$  [(Str a, Str a)]  
split xs = map ( $\lambda(x, y) \rightarrow$  (reverse x, y)) (go [] xs)
```

where

```
go p []      = [(p, [])]  
go p (c : cs) = (p, c : cs) : go (c : p) cs
```

```
split "abc"  $\Rightarrow$  [("", "abc"), ("a", "bc"), ("ab", "c"), ("abc", "")]
```


Regular Expressions refresher

Alternatives

Alternative - if A and B are regular expressions, so is A | B

$$L(A \mid B) = L(A) \cup L(B)$$

Haskell semantics:

```
acceptAlt :: Eq a => Regex a -> Regex a -> Str a -> Bool  
acceptAlt r1 r2 s = accept r1 s || accept r2 s
```

Regular Expressions refresher

Repetition

Repetition - also known as Kleene star, if A is a regular expression, so is A^*

$$L(A^*) = \{\varepsilon, A, A \cdot A, A \cdot A \cdot A, \dots\}$$

Haskell semantics:

```
acceptRep :: Eq a => Regex a -> Str a -> Bool  
acceptRep r s = or [and [accept r p | p <- ps] | ps <- parts s]
```

```
parts :: Str a -> [[Str a]]  
parts []      = [[]]  
parts [c]     = [[[c]]]  
parts (c : cs) =  
    concat [[(c : p) : ps, [c] : p : ps] | p : ps <- parts cs]
```

```
parts "abcd" =>
```

```
[[ "abcd" ], [ "a", "bcd" ], [ "ab", "cd" ], [ "a", "b", "cd" ], [ "abc", "d" ], [ "a", "bc", "d" ], [ "a", "b", "c", "d" ]]
```

Regular Expressions refresher

What is not a regular expression

- ▶ Some popular implementations allow “regexps” with backreferences, like $(a|b|c)\backslash 1$

Regular Expressions refresher

What is not a regular expression

- ▶ Some popular implementations allow “regexps” with backreferences, like $(a|b|c)\backslash 1$
- ▶ $(.*)\backslash 1$ describes non-regular language. Therefore it cannot be converted to finite automaton

Regular Expressions refresher

What is not a regular expression

- ▶ Some popular implementations allow “regexps” with backreferences, like $(a|b|c)\backslash 1$
- ▶ $(.*)\backslash 1$ describes non-regular language. Therefore it cannot be converted to finite automaton
- ▶ Generally, we want an automaton since it's very fast

Regular Expressions refresher

What is not a regular expression

- ▶ Some popular implementations allow “regexps” with backreferences, like $(a|b|c)\backslash 1$
- ▶ $(.*)\backslash 1$ describes non-regular language. Therefore it cannot be converted to finite automaton
- ▶ Generally, we want an automaton since it's very fast
- ▶ We don't include backreferences in our regexps. Regexps without them have enough expressive power

Regular Expressions refresher

Drawbacks

- ▶ Although usefull for defining semantics, *accept* function is not appropriate for practical applications.

Regular Expressions refresher

Drawbacks

- ▶ Although useful for defining semantics, *accept* function is not appropriate for practical applications.
- ▶ Reliance on function *parts* means exponential worst-case runtime in length of string being matched

Glushkov Automaton Construction

Glushkov proposed Algorithm for generating Nondeterministic Finite Automaton, NFA, from regular expression.

NFA is a good way to run regexps because

- ▶ Size of NFA is linear in size of regular expression

Glushkov Automaton Construction

Glushkov proposed Algorithm for generating Nondeterministic Finite Automaton, NFA, from regular expression.

NFA is a good way to run regexps because

- ▶ Size of NFA is linear in size of regular expression
- ▶ Given NFA with m states we can run it on string of n characters in $O(m \cdot n)$ time

Glushkov Automaton Construction

Glushkov proposed Algorithm for generating Nondeterministic Finite Automaton, NFA, from regular expression.

NFA is a good way to run regexps because

- ▶ Size of NFA is linear in size of regular expression
- ▶ Given NFA with m states we can run it on string of n characters in $O(m \cdot n)$ time
- ▶ DFA is simpler than NFA but can have $\Omega(2^m)$ states of states for an NFA with m states.

Glushkov Automaton Construction

Example

We're going to fuse translation of regexp to NFA and execution to NFA.

Glushkov Automaton Construction

Example

We're going to fuse translation of regexp to NFA and execution to NFA.

Glushkov observation: every NFA states corresponds to a position within regexp. Position is a place where a symbol occurs.

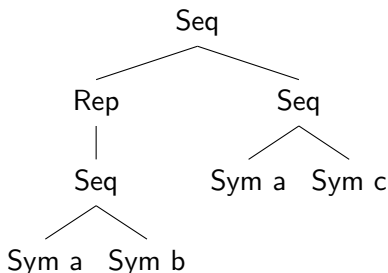
Glushkov Automaton Construction

Example

We're going to fuse translation of regexp to NFA and execution to NFA.

Glushkov observation: every NFA states corresponds to a position within regexp. Position is a place where a symbol occurs.

Let's match $(a \cdot b)^* \cdot a \cdot c$ against "ababac"



Glushkov Automaton Construction

Example - step 1

Let's match $(a \cdot b)^* \cdot a \cdot c$ against "ababac"

Glushkov Automaton Construction

Example - step 1

Let's match $(a \cdot b)^* \cdot a \cdot c$ against "ababac"

"a|babac"

Glushkov Automaton Construction

Example - step 1

Let's match $(a \cdot b)^* \cdot a \cdot c$ against "ababac"

"ababac"

$(\text{a} \cdot b)^* \cdot a \cdot c$

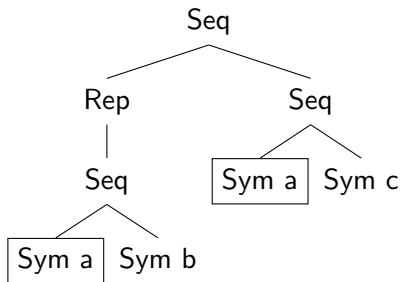
Glushkov Automaton Construction

Example - step 1

Let's match $(a \cdot b)^* \cdot a \cdot c$ against "ababac"

"a|babac"

$(a \cdot b)^* \cdot a \cdot c$



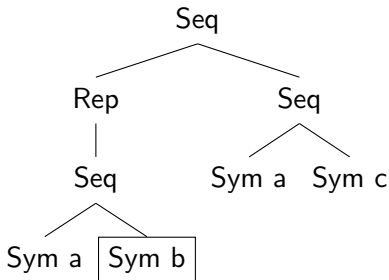
Glushkov Automaton Construction

Example - step 2

Let's match $(a \cdot b)^* \cdot a \cdot c$ against "ababac"

"a**b**abac"

$(a \cdot \boxed{b})^* \cdot a \cdot c$



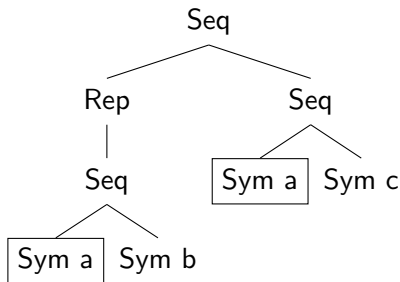
Glushkov Automaton Construction

Example - step 3

Let's match $(a \cdot b)^* \cdot a \cdot c$ against "ababac"

"ab**a**bac"

$(\boxed{a} \cdot b)^* \cdot \boxed{a} \cdot c$



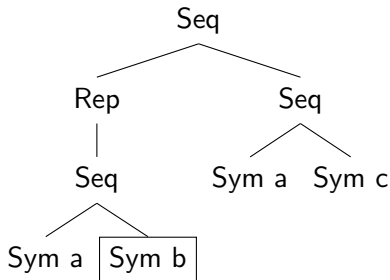
Glushkov Automaton Construction

Example - step 4

Let's match $(a \cdot b)^* \cdot a \cdot c$ against "ababac"

"aba**b**ac"

$(a \cdot \boxed{b})^* \cdot a \cdot c$



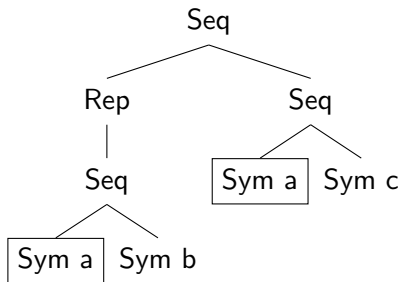
Glushkov Automaton Construction

Example - step 5

Let's match $(a \cdot b)^* \cdot a \cdot c$ against "ababac"

"ababac"

$(\boxed{a} \cdot b)^* \cdot \boxed{a} \cdot c$



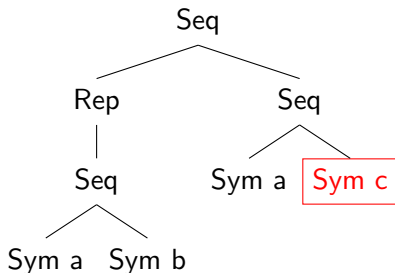
Glushkov Automaton Construction

Example - step 6

Let's match $(a \cdot b)^* \cdot a \cdot c$ against "ababac"

"ababac"

$(a \cdot b)^* \cdot a \cdot$ c



Glushkov Automaton Construction

Haskell implementation

```
data Regex' a
  = Eps'                                -- Empty regex
  | Sym' Bool a                          -- Singleton regex
  | Seq' (Regex' a) (Regex' a)         -- Sequence
  | Alt' (Regex' a) (Regex' a)         -- Alternatives
  | Rep' (Regex' a)                      -- Repetition
deriving (Show, Eq)

-- Convenient constructor
sym' :: a → Regex' a
sym' c = Sym' False c
```


Glushkov Automaton Construction

Moving mark

$shift :: Eq\ a \Rightarrow Bool \rightarrow Regex'\ a \rightarrow a \rightarrow Regex'\ a$
 $shift\ _ \ Eps' \quad \quad _ = Eps'$
 $shift\ m\ (Sym'\ _ \ c)\ c' = Sym'\ (m\ \&\&\ c == c')\ c$
 $shift\ m\ (Alt'\ p\ q)\ c = Alt'\ (shift\ m\ p\ c)\ (shift\ m\ q\ c)$
 $shift\ m\ (Seq'\ p\ q)\ c =$
 $\quad Seq'\ (shift\ m\ p\ c)\ (shift\ (m\ \&\&\ empty\ p\ ||\ final\ p)\ q\ c)$
 $shift\ m\ (Rep'\ p)\ c = Rep'\ (shift\ (m\ ||\ final\ p)\ p\ c)$

$empty :: Regex'\ a \rightarrow Bool$
 $empty\ Eps' \quad \quad = True$
 $empty\ (Sym'\ _ \ _) = False$
 $empty\ (Alt'\ p\ q) = empty\ p\ ||\ empty\ q$
 $empty\ (Seq'\ p\ q) = empty\ p\ \&\&\ empty\ q$
 $empty\ (Rep'\ p)\ \quad = True$

Glushkov Automaton Construction

Checking if mark is in the final position

$final :: Regex' a \rightarrow Bool$

$final\ Eps' = False$

$final\ (Sym'\ m\ _) = m$

$final\ (Alt'\ p\ q) = final\ p \parallel final\ q$

$final\ (Seq'\ p\ q) = final\ p \ \&\&\ empty\ q \parallel final\ q$

$final\ (Rep'\ p) = final\ p$

Glushkov Automaton Construction

Matching

$accept' :: Eq\ a \Rightarrow Regex'\ a \rightarrow [a] \rightarrow Bool$

$accept' re [] = empty\ re$

$accept' re (c : cs) = final\ (foldl\ (shift\ False)\ (shift\ True\ re\ c)\ cs)$

Glushkov Automaton Construction

Matching

$accept' :: Eq\ a \Rightarrow Regex'\ a \rightarrow [a] \rightarrow Bool$
 $accept' re [] = empty\ re$
 $accept' re (c : cs) = final\ (foldl\ (shift\ False)\ (shift\ True\ re\ c)\ cs)$

Let's try our previous test regexp $a^* \cdot b \cdot (\epsilon \mid c)$

$re = Rep'\ (sym'\ A)\ 'Seq'\ sym'\ B\ 'Seq'\ Alt'\ Eps'\ (sym'\ C)$

Expected output

$accept'\ re\ [B] \Rightarrow True$

$accept'\ re\ [B, C] \Rightarrow True$

$accept'\ re\ [A, A, A, B] \Rightarrow True$

$accept'\ re\ [A, A, A, C] \Rightarrow False$

$accept'\ re\ [A, A, A, B, C] \Rightarrow True$

Glushkov Automaton Construction

Improving efficiency

It's costly to recompute *final* and *empty* every time. We can memoize them in expressions

```
data RegexMemo a
  = EpsMemo           -- Empty regex
  | SymMemo Bool a    -- Singleton regex
  | SeqMemo (R a) (R a) -- Sequence
  | AltMemo (R a) (R a) -- Alternatives
  | RepMemo (R a)      -- Repetition
deriving (Show, Eq)

data R a = R
  { rEmpty :: Bool
  , rFinal  :: Bool
  , rRegex  :: RegexMemo a
  } deriving (Show, Eq)
```

Glushkov Automaton Construction

Improving efficiency

Having defined *RegexMemo* and *R*, we can define smart constructors and we'll only need to replace all explicit constructors with smart constructors in function *shift* to get the benefit.

E.g.

$$reEps :: R\ a$$
$$reEps = R$$
$$\begin{aligned} &\{ rEmpty = True \\ &\quad , rFinal = False \\ &\quad , rRegex = EpsMemo \\ &\quad \} \end{aligned}$$
$$reAlt :: R\ a \rightarrow R\ a \rightarrow R\ a$$
$$reAlt\ p\ q = R$$
$$\begin{aligned} &\{ rEmpty = rEmpty\ p \parallel rEmpty\ q \\ &\quad , rFinal = rFinal\ p \parallel rFinal\ q \\ &\quad , rRegex = AltMemo\ p\ q \\ &\quad \} \end{aligned}$$

One cool laziness trick

Laziness allows matching strictly more expressive languages. Even some context-sensitive ones

To make our regex lazy we'll have to add another field that tracks, whether the expression was active. If it wasn't, we can avoid forcing it. This would allow us to operate infinite regexps, which can match some context-free languages and even some context-sensitive ones.

One cool laziness trick

Laziness allows matching strictly more expressive languages. Even some context-sensitive ones

To make our regex lazy we'll have to add another field that tracks, whether the expression was active. If it wasn't, we can avoid forcing it. This would allow us to operate infinite regexps, which can match some context-free languages and even some context-sensitive ones.

Let's match language of parentheses $\{\epsilon, (), (()), ()(), (())(), \dots\}$

Our alphabet is $\Sigma = \{ (,) \}$

```
data Parens = LParen | RParen deriving (Show, Eq)
```

Now make regex that matches this language

```
let re = (symL LParen 'seqL' re 'seqL' symL RParen) 'seqL' re in re
```


One cool laziness trick

Sample run

```
let re = (symL LParen 'seqL' re 'seqL' symL RParen) 'seqL' re in re
```

One cool laziness trick

Sample run

let *re* = (*symL LParen 'seqL' re 'seqL' symL RParen*) *'seqL' re* **in** *re*

Try it out

acceptLazy re [LParen, RParen] ⇒ True

acceptLazy re [LParen, RParen, LParen, RParen] ⇒ True

acceptLazy re [LParen, RParen, LParen, RParen, RParen] ⇒ False

Thank you!

Thank you!

Questions?