

POSIX Regular Expression Parsing with Derivatives

Martin Sulzmann¹ and Kenny Zhuo Ming Lu²

Hochschule Karlsruhe - Technik und Wirtschaft
martin.sulzmann@hs-karlsruhe.de

Nanyang Polytechnic
luzhuomi@gmail.com

June 5, 2014

Regular Expression Matching is Ambiguous

- Regular Expression matching is ambiguous
- For example, there are two ways of matching “ab” with

$$(a + (b + ab))^*$$

Regular Expression Matching is Ambiguous

- Regular Expression matching is ambiguous
- For example, there are two ways of matching “ab” with

$$(a + (b + ab))^*$$

Regular Expression Matching is Ambiguous

- Regular Expression matching is ambiguous
- For example, there are two ways of matching “ab” with

$$(a + (b + ab))^*$$

Regular Expression Matching is Ambiguous

- Regular Expression matching is ambiguous
- Several disambiguation strategies exist
 - POSIX ✓
 - PCRE
 - ...

“Subpatterns should match the longest possible substrings, where sub-patterns that start earlier (to the left) in the regular expression take priority over ones starting later. Hence, higher-level subpatterns take priority over their lower-level component subpatterns. Matching an empty string is considered longer than no match at all.”

$$(a + (b + ab))^*$$

POSIX Implementation is Tricky

- We found bugs in existing NFA approaches
 - Google's re2
 - consider matching "abcd" with

$(a + b + c + d + ab + bc + cd + abc + bcd + abcd)^*$

POSIX Implementation is Tricky

- We found bugs in existing NFA approaches
 - Google's re2
 - consider matching "abcd" with

$$(a + b + c + d + ab + bc + cd + abc + bcd + abcd)^*$$

- re2 yields

$$(a + b + c + d + ab + bc + cd + abc + bcd + abcd)^*$$

POSIX Implementation is Tricky

- We found bugs in existing NFA approaches
 - Google's re2
 - consider matching "abcd" with

$$(a + b + c + d + ab + bc + cd + abc + bcd + abcd)^*$$

- re2 yields

$$(a + b + c + d + ab + bc + cd + abc + bcd + abcd)^*$$

- POSIX matching should yield

$$(a + b + c + d + ab + bc + cd + abc + bcd + abcd)^*$$

POSIX Implementation is Tricky

- We found bugs in existing NFA approaches
 - Google's re2
 - The default C POSIX shipped with Linux (GNU) and Mac OS (BSD)
 - consider matching "abcdefg" with

$$(a + bcdef + g + ab + c + d + e + efg + fg)^*$$

POSIX Implementation is Tricky

- We found bugs in existing NFA approaches
 - Google's re2
 - The default C POSIX shipped with Linux (GNU) and Mac OS (BSD)

- consider matching "abcdefg" with

$$(a + bcdef + g + ab + c + d + e + efg + fg)^*$$

- default C POSIX yields

$$(a + bcdef + g + ab + c + d + e + efg + fg)^*$$

POSIX Implementation is Tricky

- We found bugs in existing NFA approaches
 - Google's re2
 - The default C POSIX shipped with Linux (GNU) and Mac OS (BSD)

- consider matching "abcdefg" with

$$(a + bcdef + g + ab + c + d + e + efg + fg)^*$$

- default C POSIX yields

$$(a + bcdef + g + ab + c + d + e + efg + fg)^*$$

- POSIX matching should yield

$$(a + bcdef + g + ab + c + d + e + efg + fg)^*$$

POSIX Implementation is Tricky

- We found bugs in existing NFA approaches
 - Google's re2
 - The default C POSIX shipped with Linux (GNU) and Mac OS (BSD)
 - Many others
- Kuklewicz's TDFA addresses the POSIX matching
<http://hackage.haskell.org/package/regex-tdfa>
- We consider POSIX parsing using derivative (a DFA approach)

Parsing as Type Inhabitation Check

- Parsing “ab” with

$$(a + (b + ab))^*$$

produces

$$[Right\ Right\ (a, b)]$$

- Can be viewed as a type inhabitation check

$$\vdash [Right\ Right\ (a, b)] : (a + (b + ab))^*$$

i.e. Parse trees are proof terms of the type inhabitation check.
Detailed rules and definitions can be found in the paper.

Regex Parsing is also ambiguous

$$\frac{\frac{\frac{\vdash a : a \quad \vdash b : b}{\vdash (a, b) : ab}}{\vdash \text{Right } (a, b) : (a + ab)} \quad \frac{\vdash () : \epsilon}{\vdash \text{Right } () : (b + \epsilon)}}{\vdash (\text{Right } (a, b), \text{Right } ()) : (a + ab)(b + \epsilon)}$$

versus

$$\frac{\frac{\vdash a : a}{\vdash \text{Left } a : a + ab} \quad \frac{\vdash b : b}{\vdash \text{Left } b : (b + \epsilon)}}{\vdash (\text{Left } a, \text{Left } b) : (a + ab)(b + \epsilon)}$$

The Key Idea

- We use a DFA construction via derivatives
- Brzozowski's Regular expression derivatives

$$\begin{aligned}\phi \setminus l &= \phi \\ \epsilon \setminus l &= \phi \\ l_1 \setminus l_2 &= \begin{cases} \epsilon & \text{if } l_1 == l_2 \\ \phi & \text{otherwise} \end{cases} \\ (r_1 + r_2) \setminus l &= r_1 \setminus l + r_2 \setminus l \\ (r_1 r_2) \setminus l &= \begin{cases} (r_1 \setminus l) r_2 + r_2 \setminus l & \text{if } \epsilon \in L(r_1) \\ (r_1 \setminus l) r_2 & \text{otherwise} \end{cases} \\ r^* \setminus l &= (r \setminus l) r^* \end{aligned}$$

The Key Idea

- 1 Apply Brzowski's Derivative to perform the matching by extraction pass

Matching by extraction: $r_0 \xrightarrow{l_1} r_1 \xrightarrow{l_2} \dots \xrightarrow{l_n} r_n$

$r \xrightarrow{l} r'$ is a short hand for $r \setminus l = r'$.

- 2 Build the empty proof term v_n of the final expression r_n
- 3 Reversing the previous pass to construct the parse trees by injection.

Parse trees by injection: $v_0 \xleftarrow{l_1} v_1 \xleftarrow{l_2} \dots \xleftarrow{l_n} v_n$

$v \xleftarrow{l} v'$ is short-hand for $v = inj_{r \setminus l} v'$ where $\vdash v : r$ and $\vdash v' : r \setminus l$

A POSIX Parsing Example

$(a + ab)(b + \epsilon)$

A POSIX Parsing Example

$$(a + ab)(b + \epsilon) \xrightarrow{a} (\epsilon + \epsilon b)(b + \epsilon)$$

A POSIX Parsing Example

$$(a + ab)(b + \epsilon) \xrightarrow{a} (\epsilon + \epsilon b)(b + \epsilon) \xrightarrow{b} (\phi + (\phi b + \epsilon))(b + \epsilon) + (\epsilon + \phi)$$

A POSIX Parsing Example

$$(a + ab)(b + \epsilon) \xrightarrow{a} (\epsilon + \epsilon b)(b + \epsilon) \xrightarrow{b} (\phi + (\phi b + \epsilon))(b + \epsilon) + (\epsilon + \phi)$$

*mkEps*_{($\phi + (\phi b + \epsilon)$)($b + \epsilon$) + ($\epsilon + \phi$)}

$$\begin{aligned}mkEps_{r^*} &= [] \\mkEps_{r_1 r_2} &= (mkEps_{r_1}, mkEps_{r_2}) \\mkEps_{r_1+r_2} \\&\quad | \epsilon \in L(r_1) = \textit{Left } mkEps_{r_1} \\&\quad | \epsilon \in L(r_2) = \textit{Right } mkEps_{r_2} \\mkEps_{\epsilon} &= ()\end{aligned}$$

$$mkEps_{(\phi+(\phi b+\epsilon))(b+\epsilon)+(\epsilon+\phi)} = \textit{Left } (\textit{Right } (\textit{Right } ()), \textit{Right } ())$$

- We have that $\vdash mkEps_r : r$ and $|mkEps_r| = \epsilon$ if $\epsilon \in L(r)$.
- Yields the empty POSIX parse tree.

A POSIX Parsing Example

$$(a + ab)(b + \epsilon) \xrightarrow{a} (\epsilon + \epsilon b)(b + \epsilon) \xrightarrow{b} (\phi + (\phi b + \epsilon))(b + \epsilon) + (\epsilon + \phi)$$

Left (Right (Right ()), Right ())

Brzozowski's Derivatives:

$$(r_1 r_2) \setminus l = \begin{cases} (r_1 \setminus l r_2) + r_2 \setminus l & \text{if } \epsilon \in L(r_1) \\ (r_1 \setminus l r_2) & \text{otherwise} \end{cases}$$

Derivative-Guided Proof Transformation

Brzozowski's Derivatives:

$$(r_1 r_2) \setminus l = \begin{cases} (r_1 \setminus l r_2) + r_2 \setminus l & \text{if } \epsilon \in L(r_1) \\ (r_1 \setminus l r_2) & \text{otherwise} \end{cases}$$

Proof Transformation:

$$\begin{aligned} \text{inj}_{(r_1 r_2) \setminus l} v \\ | \epsilon \in L(r_1) &= \text{case } v \text{ of} \\ &\quad \text{Left } (v_1, v_2) \rightarrow (\text{inj}_{r_1 \setminus l} v_1, v_2) \\ &\quad \text{Right } v_2 \rightarrow (\text{mkEps}_{r_1}, \text{inj}_{r_2 \setminus l} v_2) \\ | \text{otherwise} &= \text{case } v \text{ of} \\ &\quad (v_1, v_2) \rightarrow (\text{inj}_{r_1 \setminus l} v_1, v_2) \end{aligned}$$

Derivative-Guided Proof Transformation

Brzozowski's Derivatives:

$$(r_1 r_2) \setminus l = \begin{cases} (r_1 \setminus l r_2) + r_2 \setminus l & \text{if } \epsilon \in L(r_1) \\ (r_1 \setminus l r_2) & \text{otherwise} \end{cases}$$

Proof Transformation:

$$\begin{aligned} \text{inj}_{(r_1 r_2) \setminus l} v \\ | \epsilon \in L(r_1) &= \text{case } v \text{ of} \\ &\quad \text{Left } (v_1, v_2) \rightarrow (\text{inj}_{r_1 \setminus l} v_1, v_2) \\ &\quad \text{Right } v_2 \rightarrow (\text{mkEps}_{r_1}, \text{inj}_{r_2 \setminus l} v_2) \\ | \text{otherwise} &= \text{case } v \text{ of} \\ &\quad (v_1, v_2) \rightarrow (\text{inj}_{r_1 \setminus l} v_1, v_2) \end{aligned}$$

Derivative-Guided Proof Transformation

- If $\vdash v : r \setminus l$ then $\vdash inj_{r \setminus l} v : r$.
- Injection preserves POSIX parse trees.
- Given

$$\begin{aligned} inj_{r_1 r_2 \setminus l} v \\ \quad | \epsilon \in L(r_1) &= \text{case } v \text{ of} \\ \quad \quad \text{Left } (v_1, v_2) &\rightarrow (inj_{r_1 \setminus l} v_1, v_2) \\ \quad \quad \text{Right } v_2 &\rightarrow (mkEps_{r_1}, inj_{r_2 \setminus l} v_2) \\ \quad \text{otherwise} &= \text{case } v \text{ of} \\ \quad \quad (v_1, v_2) &\rightarrow (inj_{r_1 \setminus l} v_1, v_2) \\ inj_{r_1+r_2 \setminus l} (\text{Left } v_1) | \epsilon \in L(r_1) &= \text{Left } inj_{r_1 \setminus l} v_1 \\ inj_{r_1+r_2 \setminus l} (\text{Right } v_2) | \epsilon \in L(r_2) &= \text{Right } inj_{r_2 \setminus l} v_2 \end{aligned}$$

we have

$$\begin{aligned} inj_{((\epsilon+\epsilon b)(b+\epsilon)) \setminus b} (\text{Left } (\text{Right } (\text{Right } ()), \text{Right } ())) &\longrightarrow \\ inj_{(\epsilon+\epsilon b) \setminus b} (\text{Right } (\text{Right } ()), \text{Right } ()) &\longrightarrow \\ (\text{Right } (inj_{\epsilon b \setminus b} (\text{Right } ())), \text{Right } ()) &\longrightarrow^* \\ (\text{Right } ((), b), \text{Right } ()) & \end{aligned}$$

A POSIX Parsing Example

$$\begin{array}{ccc} (a + ab)(b + \epsilon) \xrightarrow{a} & (\epsilon + \epsilon b)(b + \epsilon) \xrightarrow{b} & (\phi + (\phi b + \epsilon))(b + \epsilon) + (\epsilon + \phi) \\ & \text{Right } ((), b), \text{Right } () \xleftarrow{b} & \text{Left } (\text{Right } (\text{Right } ()), \text{Right } ()) \\ & \parallel & \parallel \\ & v_2 & v_3 \end{array}$$

A POSIX Parsing Example

$$\begin{array}{ccccc} (a + ab)(b + \epsilon) & \xrightarrow{a} & (\epsilon + \epsilon b)(b + \epsilon) & \xrightarrow{b} & (\phi + (\phi b + \epsilon))(b + \epsilon) + (\epsilon + \phi) \\ inj_{((a+ab)(b+\epsilon)) \setminus a} v_2 & \xleftarrow{a} & (Right ((), b), Right ()) & \xleftarrow{a} & Left (Right (Right ()), Right ()) \\ & & \parallel & & \parallel \\ & & v_2 & & v_3 \end{array}$$

A POSIX Parsing Example

$$\begin{array}{ccc} (a + ab)(b + \epsilon) \xrightarrow{a} & (\epsilon + \epsilon b)(b + \epsilon) \xrightarrow{b} & (\phi + (\phi b + \epsilon))(b + \epsilon) + (\epsilon + \phi) \\ (Right(a, b), Right()) \xleftarrow{a} & (Right((), b), Right()) \xleftarrow{b} & Left(Right(Right()), Right()) \\ & \parallel & \parallel \\ & v_2 & v_3 \end{array}$$

- Implemented in Haskell
- Apply simplification to ensure the finiteness of the derivatives
- BitCoding Forward Parsing for efficiency
- Compilation by constructing DFA from the derivatives
- Fine-tuned (sub-matching) version to keep the last match of Kleene's star
- Competitive performance

Our Contributions and Conclusion

- Formally define POSIX parsing by view regular expressions as types and parse tree as values.
- Relate parsing to the more specific submatching problem.
- Present a method to compute POSIX parse trees based on Brzozowski's derivatives and verify its correctness
- Built an optimized versions for parsing as well as submatching.

Thank you. Questions are welcome!