

From Ambiguous Regular Expressions to Deterministic Parsing Automata

Angelo Borsotti Luca Breveglieri Stefano Crespi Reghizzi
Angelo Morzenti

Politecnico di Milano - Piazza Leonardo Da Vinci n. 32 - I-20133 Milano - Italy

angelo.borsotti@mail.polimi.it luca.breveglieri@polimi.it
stefano.crespireghizzi@polimi.it angelo.morzenti@polimi.it

CIAA - 17-21 Aug. 2015 - Umeå - Sweden

downloads

paper: [full text](#) - [LNCS book](#)

SW tool: [short note](#) - [GitHub](#)

Table of Contents

- 1 Introduction
- 2 Basic Definitions
- 3 Local Testability
- 4 Parser Construction
- 5 Tree and Complexity
- 6 Disambiguation Criteria
- 7 Conclusion

Problem Statement and State of the Art

On the status of parsing Regular Expressions (*REs*)

Given a *RE*, it is well known that *checking* if a string belongs to the language $L(RE)$ can be done by a *DFA* (in real-time).

A well known method to build a recognizer *DFA* for a given *RE* is the Berry-Sethi algorithm (1986).

Parsing a *RE*, i.e., recognizing a string and *building* a syntax tree thereof, has been tackled (and solved) more recently:

- by a Thompson-based *NFA* recognizer (Dubè and Feeley, 2000)
- by the Brzozowski derivative method (Sulzmann and Lu, 2014)

On a different line, an inspiring work (Okui and Suzuki, 2010) resorts to a *NFA* that incorporates the metasymbolic structure of the *RE*, and that has a linear-time complexity for parsing a string.

Research Direction and Objectives

Objectives of the present research work

Develop a *DFA*-based method to (efficiently) parse a *RE*.

Generalize the classical Berry-Sethi algorithm (*BS*).

Obtain a Berry-Sethi Parser (*BSP*) that:

for non-ambiguous *REs*, efficiently recognizes a string, and outputs the unique syntax tree of the string

for ambiguous *REs*, represents all the syntax trees, and outputs one selected tree, e.g., *Greedy* or *POSIX*

Our method merges the metasymbolic structure of the *RE* and the Berry-Sethi approach, to construct a *DFA* that works both as a pure recognizer and as a parser able to build syntax trees.

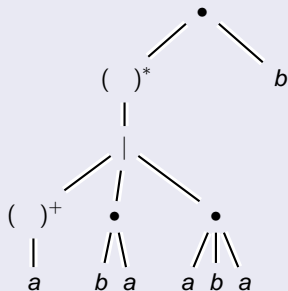
It is grounded in the theory of the locally testable languages.

Consider the following ambiguous *RE* e (it is the running example):

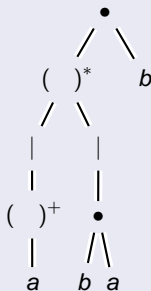
$$e = ((a)^+ \mid b \cdot a \mid a \cdot b \cdot a)^* \cdot b$$

Abstract Syntax Tree (AST) samples

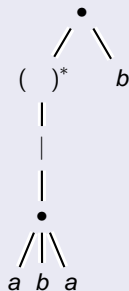
abstract syntax tree of the *RE* e



syntax trees of the string $abab \in L(e)$



Greedy



POSIX

Structural marking is hardly readable, a lighter marking helps.

Numerically marked RE \hat{e} with *partial* marking

$$\hat{e} = {}_1 \left({}_2 (a_3)_2^+ \mid b_4 a_5 \mid a_6 b_7 a_8 \right)_1^* b_9$$

The metasymbols ‘|’ and ‘.’ are left unmarked¹.

Numerically marked RE \hat{e} with *terminal* marking

$$\hat{e} = \left((a_3)^+ \mid b_4 a_5 \mid a_6 b_7 a_8 \right)^* b_9$$

This is used in the classical Berry-Sethi algorithm.

Numerically partially marked LSTs of the ambiguous string *abab*

$$\text{Greedy: } {}_1 \left({}_2 \left(\begin{matrix} a \\ 3 \end{matrix} \right)_2 \begin{matrix} b & a \\ 4 & 5 \end{matrix} \right)_1 b_9 \qquad \text{POSIX: } {}_1 \left(\begin{matrix} a & b & a \\ 6 & 7 & 8 \end{matrix} \right)_1 b_9$$

A SW tool should use the structural representation of the LST.

¹They can be easily inferred from their neighbours.

Use square brackets '[' and ']' as the new metasymbolic parentheses.

Reintroduction rules of the metasymbols - from \hat{e} to \check{e}

parenthesis pair	reintroduction rule
${}_h(\hat{e})_h$	${}_h([\hat{e}])_h$
${}_h(\hat{e})_h^*$	$\left[{}_h([\hat{e}]^+)_h \mid 1_h \right]$
${}_h(\hat{e})_h^+$	${}_h([\hat{e}]^+)_h$

Reintroduction example - REs: e , \hat{e} (marked) and \check{e} (reparenthesized)

$$\begin{aligned}
 e &= \left((a)^+ \mid ba \mid aba \right)^* b & \hat{e} &= {}_1 \left({}_2 (a_3)_2^+ \mid b_4 a_5 \mid a_6 b_7 a_8 \right)_1^* b_9 \\
 \check{e} &= \left[{}_1 \left(\left[{}_2 ([a_3]^+)_2 \mid b_4 a_5 \mid a_6 b_7 a_8 \right]^+ \right)_1 \mid 1_1 \right] b_9 \\
 abab \in L(e) & \iff & & {}_1 \left({}_2 \left(\begin{matrix} a \\ 3 \end{matrix} \right)_2 \begin{matrix} b & a \\ 4 & 5 \end{matrix} \right)_1 \begin{matrix} b \\ 9 \end{matrix} \in L(\check{e})
 \end{aligned}$$

Alphabets

Σ	$= \{ a, b, \dots \}$	terminals
$\widehat{\Sigma}$	$= \{ a_1, a_2, \dots \}$	marked terminals
\widehat{M}	$= \{ 1_h^a, h(,)_h \mid 1 \leq h \leq \text{a bound} \}$	marked metasympols
$\widehat{\Omega}$	$= \widehat{\Sigma} \cup \widehat{M}$	all marked symbols

^aSymbol 1_h represents the empty string ε at position h .

Languages

original: $L(e) \subseteq \Sigma^*$ reparenthetized: $L(\check{e}) \subseteq \widehat{\Omega}^+$

Important observation

Given a *RE* e , it is evident that the reparenthetized language $L(\check{e})$ is the language of all and only the *LSTs* of the strings generated by e .

Let a, b be terminals, x, y strings (possibly empty) and L a language.

Classical local sets of a (generic) language L

$$\begin{aligned} Ini(L) &= \{ a \mid ax \in L \} & Fin(L) &= \{ b \mid xb \in L \} \\ Dig(L) &= \{ ab \mid xaby \in L \} & Fol(L, a) &= \{ b \mid ab \in Dig(L) \} \end{aligned}$$

Locally testable language L

$$L = Ini \cdot \Sigma^* \cap \Sigma^* \cdot Fin \cap \{ \text{strings with only digrams} \in Dig \}$$

A locally testable language is characterized by its Ini , Fin and Dig sets.

The local sets of a regular language are easily computable².

Important observation

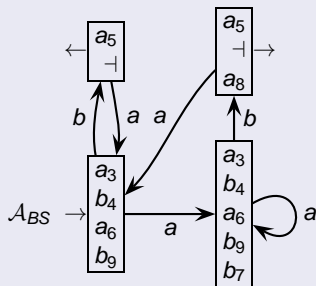
Since the reparenthesized $RE \check{e}$ is *linear* (no repeated symbols), the (reparenthesized) language $L(\check{e}) \subseteq \widehat{\Omega}^+$ of the *LSTs* is *locally testable* !

²See e.g., Berstel and Pin (1996).

$$\hat{e} = ((a_3)^+ \mid b_4 a_5 \mid a_6 b_7 a_8)^* b_9 \dashv$$

Classical Berry-Sethi recognizer \mathcal{A}_{BS}

Initials	a_3	b_4	a_6	b_9
a_3	a_3	b_4	a_6	b_9
b_4	a_5			
a_5	a_3	b_4	a_6	b_9
a_6	b_7			
b_7	a_8			
a_8	a_3	b_4	a_6	b_9
b_9	\dashv			



It holds: $L(\mathcal{A}_{BS}) = L(e)$ (Berry and Sethi, 1986).

Let $\mu \in \widehat{M}^*$ be a (possibly empty) sequence of marked *metasymbols*, without repetitions³. Let $a_h, b_k \in \widehat{\Sigma}$ be marked *terminal symbols*.

Extended set of the initial terminals of the *RE* e - Ini_E

$$Ini_E(e) = \{ \mu a_h \mid Ini(\mu a_h) \subseteq Ini(\check{e}) \wedge Dig(\mu a_h) \subseteq Dig(\check{e}) \}$$

Extended set of the terminals following a_h in the *RE* e - Fol_E

$$Fol_E(e, a_h) = \{ \mu b_k \mid Dig(a_h \mu b_k) \subseteq Dig(\check{e}) \}$$

Intuitively, such sets represent the acyclic metasymbolic paths in the *RE* to reach the initial terminals or between any two adjacent ones.

The extended local sets of a *RE* e can be easily computed (see paper), starting from the classical local sets of the reparenthetized *RE* \check{e} .

³So the number of μ prefixes is finite.

$$\hat{e} = {}_1 ({}_2 (a_3)_2^+ \mid b_4 a_5 \mid a_6 b_7 a_8)_1^* b_9 \dashv$$

Extended sets of the initials and followers of the RE e

Initials	${}_1 ({}_2 (a_3$	${}_1 (b_4$	${}_1 (a_6$	${}_1 b_9$
a_3	a_3	$)_2 {}_2 (a_3$	$)_2 b_4$	$)_2 a_6$
b_4	a_5			$)_2)_1 b_9$
a_5	${}_2 (a_3$	b_4	a_6	$)_1 b_9$
a_6	b_7			
b_7	a_8			
a_8	${}_2 (a_3$	b_4	a_6	$)_1 b_9$
b_9	\dashv			

Item - a fragment of a syntax tree

$$\text{item id: } \langle \mu a_h, \overbrace{\{\text{item id.s}\}}^{\text{links to predecessors}} \rangle \quad \mu \in \widehat{M}^* \quad a_h \in \widehat{\Sigma} \cup \dashv$$

The prefix μ must not contain any pair of repeated metasymbols.

The states of the *BSP DFA* \mathcal{A}_{BSP} are sets (non-empty) of items.

Sketch of the *BSP* construction algorithm

Input: extended local sets Ini_E and Fol_E of a RE e .

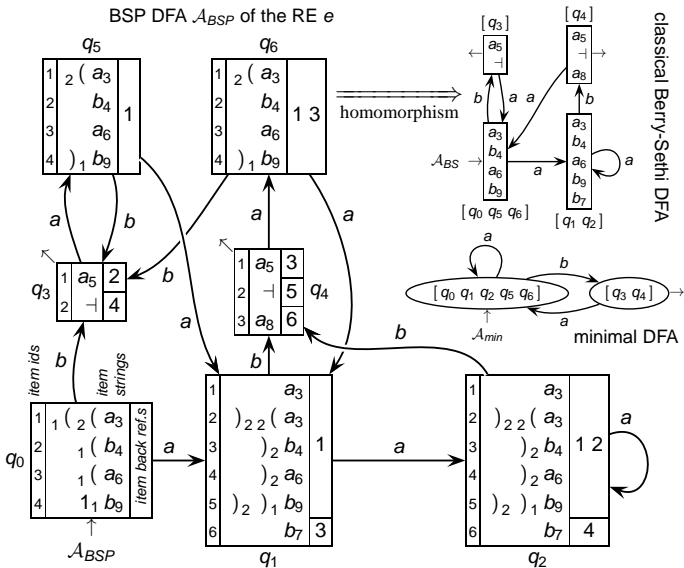
Output: state-transition graph of the *DFA* \mathcal{A}_{BSP} of e .

Set the initial state of \mathcal{A}_{BSP} to the set $Ini_E(e \dashv)$.

For each not yet examined state $p \in \mathcal{A}_{BSP}$, for each unmarked terminal $a \in \Sigma$
 for all the (source) items $id: \langle \mu a_h, \{\dots\} \rangle \in p$ such that $unmark(a_h) = a$
 create^a a new state q with all the (target) items in the sets $Fol_E(e \dashv, a_h)$,
 link such targets to their sources id , and create a new transition $p \xrightarrow{a} q$.

Set the final states of \mathcal{A}_{BSP} to those that contain \dashv .

^aUnless it was already created.



Equivalence of the *BSP* and *BS DFAs*

For a *RE* e , the two *DFAs* \mathcal{A}_{BSP} (parser) and \mathcal{A}_{BS} (classical recognizer) of e are equivalent, i.e., both recognize the same language $L(e)$.

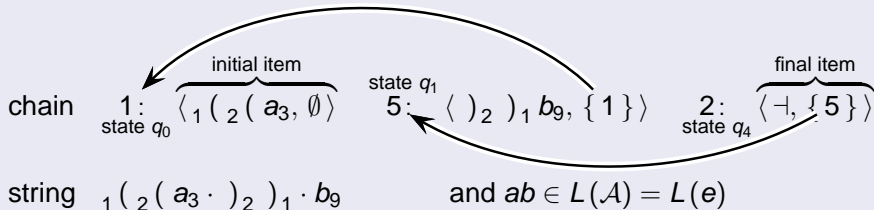
Sketch of proof

The classical recognizer \mathcal{A}_{BS} is a homomorphic image of the new parser \mathcal{A}_{BSP} . By canceling the item metasymbols, as well as the references to the predecessor items, the *DFA* \mathcal{A}_{BSP} (after merging the states that have turned identical) becomes exactly the *DFA* \mathcal{A}_{BS} .

Further observations

In general neither \mathcal{A}_{BSP} nor \mathcal{A}_{BS} are minimal. The sets of backward references in a state of \mathcal{A}_{BSP} are either disjoint or coincident. This helps us optimize the parsing. The time complexity of the *BSP* construction is (at least) exponential in the size of the alphabet $\hat{\Omega}$.

Item chain and string - sample from the running example

Extended language of the *BSP*

Given a *BSP* \mathcal{A} , define the *extended language* $\hat{L}(\mathcal{A}) \subseteq \hat{\Omega}^+$ to be the set of the strings (over marked symbols) of all the item chains in \mathcal{A} .

Infinite ambiguity - Kleene star (or cross) with a nullable argument

A string is infinitely ambiguous if and only if it has a *LST* that contains a substring of type $1_h \mu 1_h$, for some mark h and metasymbols $\mu \in \widehat{M}^*$.

Main theorem - The *BSP* extended language is the set of *RE LSTs*

Let e be a *RE* and \mathcal{A} the *BSP* of e . The extended language $\widehat{L}(\mathcal{A})$ of \mathcal{A} coincides with the language $L(\check{e})$ of the *LSTs* of e , removing from $L(\check{e})$ the strings that contain a substring of type $1_h \mu 1_h$ (infinite ambiguity).

Sketch of proof

Language $\widehat{L}(\mathcal{A})$ is locally testable and defined by the extended local sets^a Ini_E and Fol_E of e , characterized by the classical local sets Ini , Fin and Dig of \check{e} . Language $L(\check{e})$ is also locally testable and defined by the same classical local sets. Thus $\widehat{L}(\mathcal{A})$ and $L(\check{e})$ coincide (see note).

^aTheir definitions rule out any pair of repeated metasymbols, e.g., $1_h \mu 1_h$.

How to check the ambiguity of a RE

A RE is non-ambiguous if and only if every item of its BSP has at most one predecessor and every state of its BSP has at most one final item.

The same holds for a string, by restricting to the string state sequence.

How to parse a non-ambiguous string $w \in L(e)$ - two phases

First recognize w by using the $BSP \mathcal{A}$ of e and find the state sequence s of w (s is unique as \mathcal{A} is a DFA). **Second** follow back the unique item chain of w in s . The string of such a chain is the unique LST of w .

Parsing complexity (non-ambiguous case)

The time complexity of parsing a non-ambiguous string $w \in \Sigma^*$ and building its LST is linear in the string length $|w|$.

Notice a string state sequence contains the whole string tree forest⁴.

Disambiguation criteria - *Greedy* and *POSIX*

To select a prior *LST* for an ambiguous string, *Greedy*^a makes the *RE* leftmost choices, *POSIX* (an *IEEE* standard) makes the *RE* longest matches (both criteria are formalized by Frisch and Cardelli, 2004).

^aDiverging interpretations exist, e.g., Frisch and Cardelli vs *Java.regex*.

How to parse an ambiguous string $w \in L(e)$ - three phases (diff.s red)

First phase as before. **Second** follow back **all** the item chains of w in s , and tag the touched items. **Third** navigate such tagged items from the initial ones, and at each bifurcation point make the prior *LST* choice.

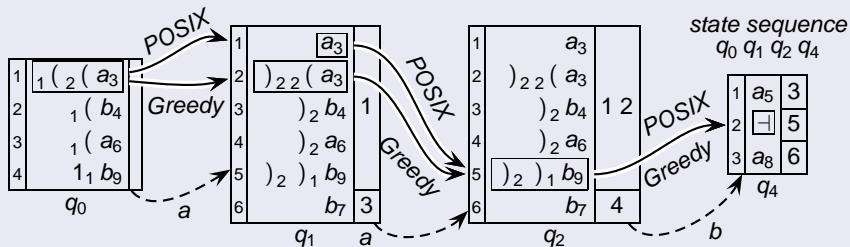
Parsing complexity (ambiguous case)

The time complexity of parsing an ambiguous string $w \in \Sigma^*$ and building its *Greedy* or *POSIX LST* is linear in the string length $|w|$.

⁴In a packed form, with common tree parts shared.

$$\hat{e} = {}_1 ({}_2 (a_3)_2^+ | b_4 a_5 | a_6 b_7 a_8)_1^* b_9 \dashv$$

Parsing the ambiguous string $aab \in L(e)$



Here *Greedy* is computed according to the class *Java.regex*.

About the *POSIX* criterion

In the third parsing phase, for *POSIX* we use the Okui and Suzuki algorithm (2010), which dynamically bookkeeps information while navigating the items, and reads it back to make the prior choice.

Current results

A det. parser (*BSP*) for any *RE*, even ambiguous, such that:
it generalizes the Berry-Sethi pure recognizer
disambiguation criteria can be easily plugged in
parsing a string (building its *LST*) is linear-time

The *BSP* also provides a simple test of the *RE* ambiguity.

A demonstrative *SW* tool (in *JavaScript*) is available ([GitHub](#)).

Future research

Complexity analytical evaluation vs other approaches.

Experimental performance evaluation (on benchmarks).

Ongoing research . . .

Experimental performance evaluation (using *Java*).

Promising figures (speed-up vs Okui-Suzuki *NFA*).

Three significant references



G. Berry and R. Sethi.

From regular expressions to deterministic automata.

Theor. Comput. Sci., 48(1):117–126, 1986.



A. Frisch and L. Cardelli.

Greedy regular expression matching.

In J. Díaz, J. Karhumäki, A. Lepistö, and D. Sannella, editors, *ICALP*, volume 3142 of *LNCS*, pages 618–629. Springer, 2004.



S. Okui and T. Suzuki.

Disambiguation in regular expression matching via position automata with augmented transitions.

In M. Domaratzki and K. Salomaa, editors, *CIAA*, volume 6482 of *LNCS*, pages 231–240. Springer, 2010.