

# Deterministic Regular Expressions in Linear Time

reasoning about deterministic regular expressions  
without building the automata

Benoit Groz<sup>1,3</sup>, Sebastian Maneth<sup>2,3</sup>, Slawek Staworko<sup>1,3</sup>

<sup>1</sup>University of Lille, Mostrare INRIA

<sup>2</sup>UNSW, NICTA

<sup>3</sup>Équipe associée Transduce

PODS, May, 2012

# Context

DTD and XML Schema use regular expressions to specify which sequence of children may appear below a node.

```
<!ELEMENT book (author,chapter*,index?)>
```

Constraint: regular expression must be **deterministic**.

We provide new algorithms to:

- Check if a regular expression is deterministic.
- Decide the membership problem for deterministic regular expressions.

# Outline

- 1 Glushkov relations: First, Last, Follow ... and determinism
- 2 Problem statement
- 3 Structure of the expression
- 4 Algorithms to test membership

# Outline

- 1 Glushkov relations: First, Last, Follow ... and determinism
- 2 Problem statement
- 3 Structure of the expression
- 4 Algorithms to test membership

# Structure of regular expressions

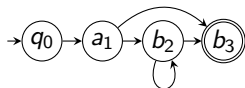
$ab^*b$

$abb^*$

# Structure of regular expressions

$a_1 b_2^* b_3$

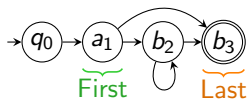
$a_1 b_2 b_3^*$



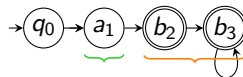
$b_3$  follows  $a_1$ ,  $b_2$  follows  $a_1 \dots$

# Structure of regular expressions

$a_1 b_2^* b_3$

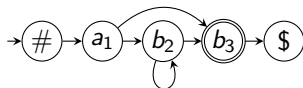


$a_1 b_2 b_3^*$

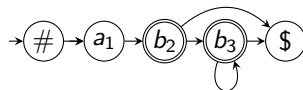


# Structure of regular expressions

$\#a_1b_2^*b_3\$$



$\#a_1b_2b_3^*\$$





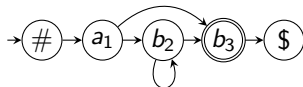
# Deterministic regular expressions (a.k.a. one-unambiguous)

Expression is *non deterministic* if:

$$b_i \xrightarrow{\quad} a_j \rightarrow a_k \quad (j \neq k)$$

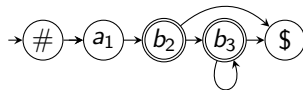
$\#a_1b_2^*b_3\$$

$\Rightarrow$  *non deterministic*



$\#a_1b_2b_3^*\$$

$\Rightarrow$  *deterministic*



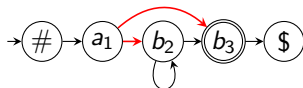
# Deterministic regular expressions (a.k.a. one-unambiguous)

Expression is *non deterministic* if:

$$b_i \xrightarrow{\quad} a_j \xrightarrow{\quad} a_k \quad (j \neq k)$$

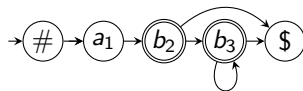
$\#a_1b_2^*b_3\$$

$\Rightarrow$  *non deterministic*



$\#a_1b_2b_3^*\$$

$\Rightarrow$  *deterministic*



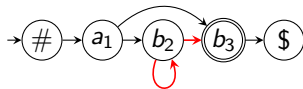
# Deterministic regular expressions (a.k.a. one-unambiguous)

Expression is *non deterministic* if:

$$b_i \xrightarrow{\quad} a_j \rightarrow a_k \quad (j \neq k)$$

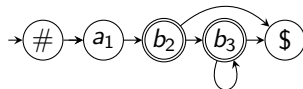
$\#a_1b_2^*b_3\$$

$\Rightarrow$  *non deterministic*



$\#a_1b_2b_3^*\$$

$\Rightarrow$  *deterministic*



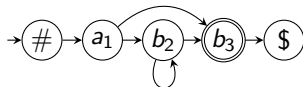
# Deterministic regular expressions (a.k.a. one-unambiguous)

Expression is *non deterministic* if:

$$b_i \xrightarrow{\quad} a_j \xrightarrow{\quad} a_k \quad (j \neq k)$$

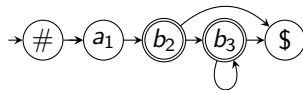
$\#a_1b_2^*b_3\$$

$\Rightarrow$  *non deterministic*



$\#a_1b_2b_3^*\$$

$\Rightarrow$  *deterministic*



Ambiguity parsing  $w = a**b**$

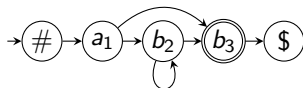
# Deterministic regular expressions (a.k.a. one-unambiguous)

Expression is *non deterministic* if:

$$b_i \xrightarrow{\quad} a_j \quad a_k \quad (j \neq k)$$

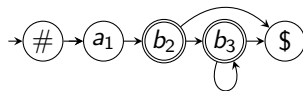
$\#a_1b_2^*b_3\$$

$\Rightarrow$  *non deterministic*



$\#a_1b_2b_3^*\$$

$\Rightarrow$  *deterministic*



$e = (a + b)b?(ab)^* \quad ?$

$e = (ab+ba?)^* \quad ?$

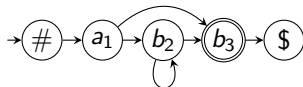
# Deterministic regular expressions (a.k.a. one-unambiguous)

Expression is *non deterministic* if:

$$b_i \xrightarrow{\quad} a_j \quad a_k \quad (j \neq k)$$

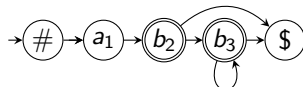
$\#a_1b_2^*b_3\$$

$\Rightarrow$  *non deterministic*



$\#a_1b_2b_3^*\$$

$\Rightarrow$  *deterministic*



$e = (a + b)b?(ab)^* \Rightarrow$  *deterministic*

$e = (ab + ba)^* \Rightarrow$  *non deterministic*:  $w = b\mathbf{a}$

# Outline

- 1 Glushkov relations: First, Last, Follow ... and determinism
- 2 Problem statement
- 3 Structure of the expression
- 4 Algorithms to test membership

# Problems of interest

## Testing determinism:

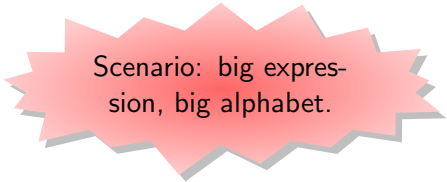
Input: expression  $e$ ,

Question: is  $e$  deterministic?

## Membership:

Input: word  $w$ , deterministic expression  $e$ ,

Question:  $w \in L(e)$ ?



Scenario: big expression,  
big alphabet.

## Remark:

size of  $e$  = number of nodes in the parse tree  
 $\simeq$  number of positions.



# Problems of interest

## Testing determinism:

Input: expression  $e$ ,

Question: is  $e$  deterministic?

## Membership:

Input: word  $w$ , deterministic expression  $e$ ,

Question:  $w \in L(e)$ ?

Straightforward solution through Glushkov automaton.

Build Glushkov in  $O(|\Sigma| \times |e|)$  [Brüggeman-Klein TCS'93].

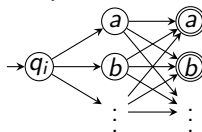
$\implies$  (quadratic in  $|e|$ )

Number of transitions of Glushkov can be quadratic:

$$e = (a + b + c \dots)(a + b + c \dots),$$

$$e' = (a + b + c \dots)^*,$$

$$e'' = (a?b?c?\dots)$$



# Problems of interest

## Testing determinism:

Input: expression  $e$ ,

Question: is  $e$  deterministic?

## Membership:

Input: word  $w$ , deterministic expression  $e$ ,

Question:  $w \in L(e)$ ?

Straightforward solution through Glushkov automaton.

Build Glushkov in  $O(|\Sigma| \times |e|)$  [Brüggeman-Klein TCS'93].

$\implies$  (quadratic in  $|e|$ )

Number of transitions of Glushkov can be quadratic:

$$e = (a + b + c \dots)(a + b + c \dots),$$

$$e' = (a + b + c \dots)^*,$$

$$e'' = (a?b?c? \dots)$$

# Problems of interest

## Testing determinism:

Input: expression  $e$ ,

Question: is  $e$  deterministic?

## Membership:

Input: word  $w$ , deterministic expression  $e$ ,

Question:  $w \in L(e)$ ?

Straightforward solution through Glushkov automaton.

Build Glushkov in  $O(|\Sigma| \times |e|)$  [Brüggeman-Klein TCS'93].

$\implies$  (quadratic in  $|e|$ )

*Can we do better?*

# Summary

	Glushkov	Our results
testing determinism	$O( \Sigma  \times  e )$	$O( e )$
membership	$O( \Sigma  \times  e  +  w )$	$O( e  +  w  \log \log(e))$
★ $k$ -occurrence		$O(k \times  w )$
★ restrictions on +		$O( e  +  w )$
★ star free		$O( e  +  w )$

This  
talk

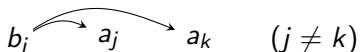
# Roadmap

## Preprocessing

We do not construct automaton. Instead we work on parse tree which we preprocess in linear time to build some pointers+datastructures.

- **Testing determinism.**

We search a witness for non-determinism in  $e$ : pair of two positions with same label that follow a common position.



*We limit the number of pairs examined to  $O(|e|)$  using skeleta from [Bojańczyk and Parys JACM'11].*

- **Testing membership.**

We simulate transitions on-the-fly using the pointers from preprocessing.

# Roadmap: transition simulation

We want a procedure for transition simulation:

**Input:**

- a position  $a_i$  in the expression,
- a letter  $b$  (the next letter of the word)

**Output:**

- the unique  $b$ -labeled position that follows  $a_i$

# Outline

- 1 Glushkov relations: First, Last, Follow ... and determinism
- 2 Problem statement
- 3 Structure of the expression**
- 4 Algorithms to test membership

# Some tools



## LCA preprocessing [Harel&Tarjan, SICOMP'84]

One can preprocess a tree in linear time, to answer lowest common ancestor (LCA) queries and ancestor queries in constant time.

(Input for preprocessing:  $t$ )

### **LCA query:**

Input: two nodes  $n_1, n_2$  in  $t$

Output: the lowest common ancestor (LCA) of  $n_1$  and  $n_2$

### **Ancestor query:**

Input: two nodes  $n_1, n_2$  in  $t$

Question: is  $n_1$  an ancestor of  $n_2$ ?

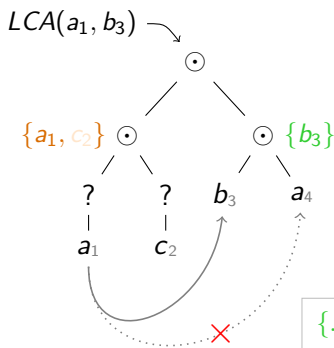


## Work on parse tree

Goal: Given two positions  $a_i$  and  $b_j$ ,  
test if  $b_j$  follows  $a_i$  in constant time.

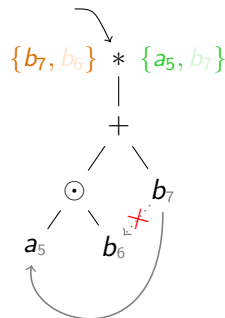
# Work on parse tree: Follow

When does  $b_j$  follow  $a_i$ ?



Case ❶ :  $\odot$

lowest star above  
 $LCA(b_7, a_5)$



Case ❷ :  $*$

$\{\dots\}$  : First set

$\{\dots\}$  : Last set

## Work on parse tree: Follow (2)

Assuming that after some preprocessing we can test *First* and *Last* in constant time,

### Theorem

We can test if  $b_j$  follows  $a_i$  in constant time.

## Work on parse tree: Follow (2)

Assuming that after some preprocessing we can test *First* and *Last* in constant time,

### Theorem

We can test if  $b_j$  follows  $a_i$  in constant time.

Preprocessing:

- pointer to lowest \* ancestor of each node.
- build *LCA*, *First* and *Last* structures.

Algo. to test if  $b_j$  follows  $a_i$ :

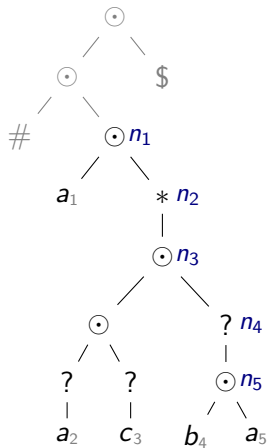
compute  $LCA(b_j, a_i)$

follow the \* pointer (for case 2: \*)

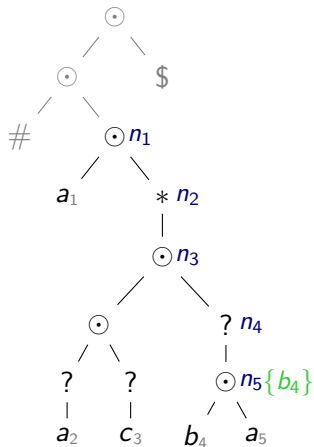
test that  $a_i, b_j$  in *First* and *Last* of appropriate nodes

*New objective: test if  $a_i \in \text{First}(n)$  in constant time*

## Work on parse tree: First and Last

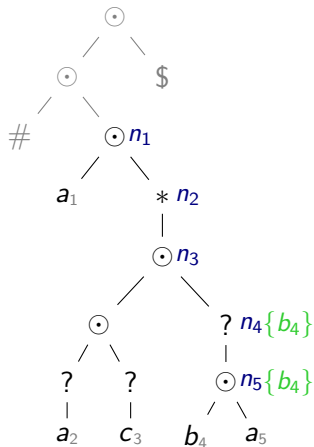


## Work on parse tree: First and Last



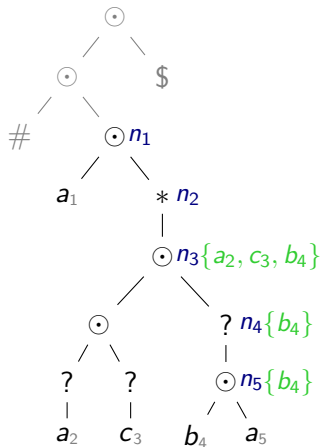
$\{a_2, \dots\}$  : *First set*

## Work on parse tree: First and Last



$\{a_2, \dots\}$  : *First set*

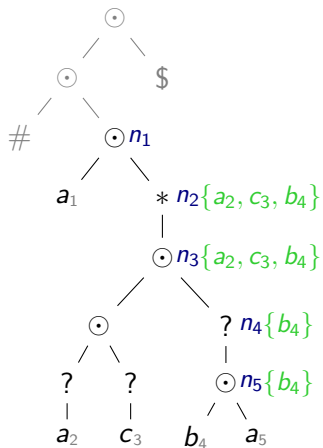
## Work on parse tree: First and Last



$\{a_2, \dots\}$  : *First set*

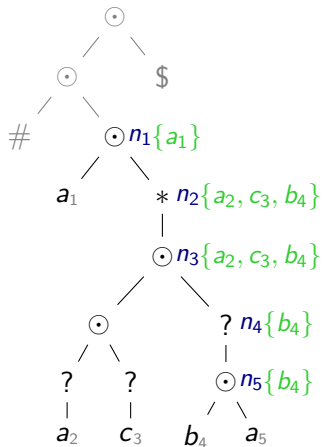


# Work on parse tree: First and Last



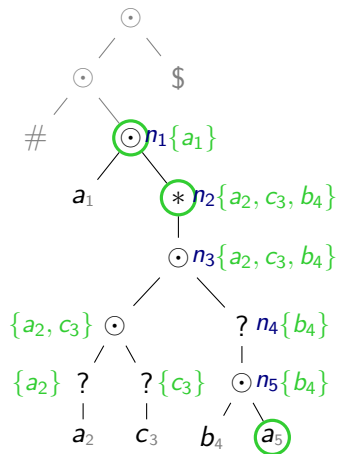
$\{a_2, \dots\}$  : *First set*

# Work on parse tree: First and Last



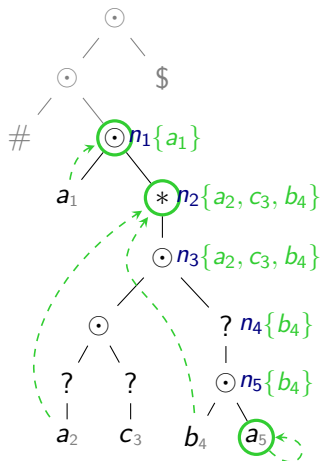
$\{a_2, \dots\}$  : First set

## Work on parse tree: First and Last



$\{a_2, \dots\}$  : *First set*

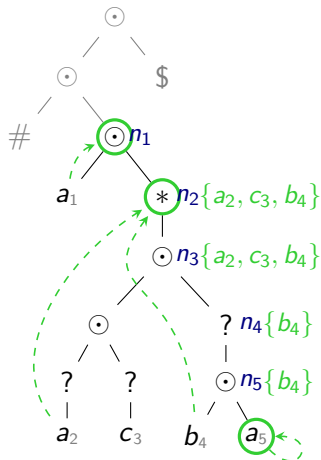
# Work on parse tree: First and Last



$\{a_2, \dots\}$  : First set

$\dashrightarrow$  : a SupFirst pointer

## Work on parse tree: First and Last

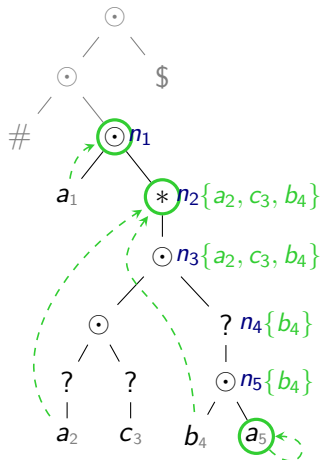


We compute *SupFirst* pointers in simple traversal.

$b_4 \in \text{First}(n_3)$  because  
 $n_2$  ancestor of  $n_3$   
and  $n_3$  ancestor of  $b_4$ .

$a_5 \notin \text{First}(n_3)$

## Work on parse tree: First and Last



We compute *SupFirst* pointers in simple traversal.

$b_4 \in \text{First}(n_3)$  because  
 $n_2$  ancestor of  $n_3$   
and  $n_3$  ancestor of  $b_4$ .

$a_5 \notin \text{First}(n_3)$

✓ We can test if  $a_i \in \text{First}(n)$  in constant time.

Symmetrically, we can test if  $a_i \in \text{Last}(n)$  in constant time.

# Reminder

## Theorem

We can test if  $b_j$  follows  $a_i$  in constant time.

# Outline

- 1 Glushkov relations: First, Last, Follow ... and determinism
- 2 Problem statement
- 3 Structure of the expression
- 4 Algorithms to test membership**



## Simple case: $k$ -occurrence expression

### Theorem

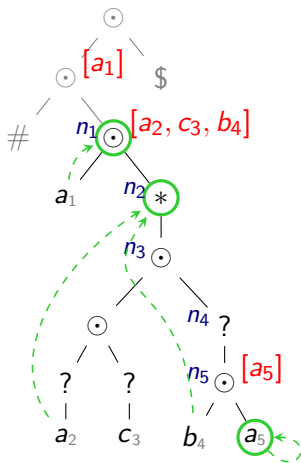
For deterministic  $k$ -occurrence expression, the membership problem can be solved in  $O(k|w|)$  after  $O(|e|)$  preprocessing.

each transition simulated in  $O(k)$ :

$$e = (ba)^* \boxed{c} aca\cancel{b}cc?b$$


## General case

We put color  $a$  in  $\text{parent}(\text{SupFirst}(a_i))$  and store  $a_i$  as the witness for color  $a$  in that node.



OBSERVATION: positions followed by  $a_5$  are below  $n_5$ , those followed by  $a_2$  or  $b_4$  are below  $n_1 \dots$

$\text{---} \rightarrow$  : a *SupFirst* pointer  
 $[a_5]$  : map  $a \mapsto a_5$   
(color  $a$ , witness  $a_5$ )

# Finding the right ancestor

We can use “nearest colored ancestor queries”.

**Nearest colored ancestor [Muthukrishnan et al. 96]**

We can preprocess a tree  $t$  in expected linear time  $O(|t|)$  to answer nearest colored ancestor queries in  $O(\log \log |t|)$ .

Expected time because of hashmaps, but becomes worst-case linear using lazy arrays.

**Evaluation algorithm**

Repeatedly jump to the nearest ancestor with color  $a$ , and test if its witness follows.

Why is it linear?

Use amortization argument

# Summary

	Glushkov	Our results
testing determinism	$O( \Sigma  \times  e )$	$O( e )$
membership	$O( \Sigma  \times  e  +  w )$	$O( e  +  w  \log \log(e))$
★ $k$ -occurrence ( $k$ -ORE)		$O(k \times  w )$
★ restrictions on +		$O( e  +  w )$
★ star free		$O( e  +  w )$

This  
talk

## Future work

- Close the log log gap for membership.
- Searching regular pattern instead of matching (KMP...)

Thanks for your attention!

For a few dollars more...

Questions are most welcome!

...but there is no guarantee for the answer

Questions ?

