

# Efficient Inclusion Checking for Deterministic Tree Automata and DTDs

Jérôme Champavère

Joint work with R. Gilleron, A. Lemay and J. Niehren

Université de Lille, INRIA, Mostrare

LATA 2008, Tarragona, España

March 14, 2008

# Tree Automata Inclusion

- Tree automata  $A$  and  $B$
- Instance:  $L(A) \subseteq L(B)$ ?
- Non deterministic  $B$ : DEXPTIME-complete [Seidl 1994]
- Deterministic  $B$ 
  - $B$  is complete:  $O(|A| * |B|)$
  - Not complete: ?

# Deterministic Tree Automata Inclusion

## Usual algorithm

- Ranked automata, arity  $n$
- Decide  $L(A) \cap L(B)^c = \emptyset$

## Complexity

- Running time:  $O(|A| * |\Sigma| * |B|^n)$
- Involves computation of  $B$  complement
- Linear reduction to binary automata:  $O(|A| * |\Sigma| * |B|^2)$
- Our goal:  $O(|A| * |B|)$ , independently of  $|\Sigma|$

# Inclusion into DTDs

- $L(A)$ : set of XML documents
  - Automata for unranked trees [TATA 2007]
  - Possibly non deterministic
- Representations for DTDs [Martens & Niehren 2007]
  - Hedge automata
  - Stepwise automata
- Algorithm for testing  $L(A) \subseteq L(DTD)$

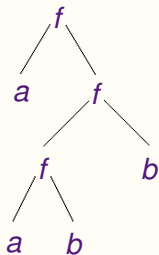
# Outline

- 1 Inclusion in Deterministic Binary Tree Automata
- 2 Inclusion in Deterministic DTDs

# Tree Automata

## Basic notions

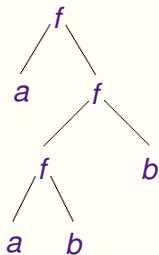
- Binary tree  $t$ :



# Tree Automata

## Basic notions

- Binary tree  $t$ :



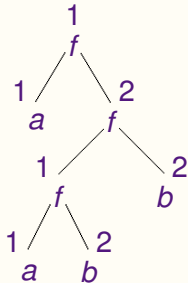
- Automaton  $A$ :

- States:  $\{1, 2\}$
- Final states:  $\{2\}$
- Rules: {
  - $a \rightarrow 1,$
  - $b \rightarrow 2,$
  - $f(1, 2) \rightarrow 1,$
  - $f(1, 2) \rightarrow 2 \quad \}$

# Tree Automata

## Basic notions

- Binary tree  $t$ :



- Automaton  $A$ :

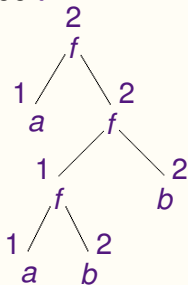
- States:  $\{1, 2\}$
- Final states:  $\{2\}$
- Rules: {
  - $a \rightarrow 1,$
  - $b \rightarrow 2,$
  - $f(1, 2) \rightarrow 1,$
  - $f(1, 2) \rightarrow 2 \quad \}$



# Tree Automata

## Basic notions

- Binary tree  $t$ :



- $\text{eval}_A(t) = \{1, 2\}$

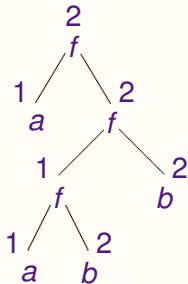
- Automaton  $A$ :

- States:  $\{1, 2\}$
- Final states:  $\{2\}$
- Rules: {
  - $a \rightarrow 1,$
  - $b \rightarrow 2,$
  - $f(1, 2) \rightarrow 1,$
  - $f(1, 2) \rightarrow 2$  }

# Tree Automata

## Basic notions

- Binary tree  $t$ :



- $\text{eval}_A(t) = \{1, 2\}$

- Automaton  $A$ :

- States:  $\{1, 2\}$
- Final states:  $\{2\}$
- Rules: {
  - $a \rightarrow 1,$
  - $b \rightarrow 2,$
  - $f(1, 2) \rightarrow 1,$
  - $f(1, 2) \rightarrow 2 \quad \}$

$$L(A) = \{t \mid \text{eval}_A(t) \cap \text{final}(A) \neq \emptyset\}$$

# Emptiness/Accessibility

## Emptiness

**Input:**  $A, \Sigma$

**Output:**  $L(A) = \emptyset$

Algorithm: does  $A$  have accessible final states?

## Accessibility

$\text{acc}(p)$  in  $A$  if  $\exists t$  such that  $p \in \text{eval}_A(t)$ .

# Efficiency of Ground Datalog

Inferring a ground Datalog program for accessibility

$$\frac{a \rightarrow p \in \text{rules}(A)}{\text{acc}(p).} \qquad \frac{f(p_1, p_2) \rightarrow p \in \text{rules}(A)}{\text{acc}(p) :- \text{acc}(p_1), \text{acc}(p_2).}$$

Number of clauses:  $O(|A|)$

**Theorem (Dantsin et al. 2001)**

*The least fixed point  $\text{lfp}(P)$  of a ground Datalog program  $P$  can be computed in time  $O(|P|)$ , even for unbounded predicate arities.*

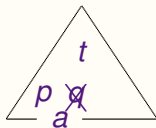
# Guiding Ideas of Characterization

- Detection of inclusion failure  $L(A) \not\subseteq L(B)$
- Find  $t \in L(A)$  such that  $t \notin L(B)$
- Accessibility on product  $A \times B$
- Ground Datalog program  $D_{A,B}$
- Conditions:  $B$  deterministic,  $A$  productive

# Inclusion Failure

## Constant rules

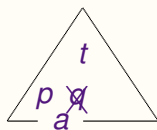
$$\frac{\exists q. \begin{array}{l} a \rightarrow p \text{ in } A \\ a \rightarrow q \text{ in } B \end{array}}{\text{fail.}}$$



# Inclusion Failure

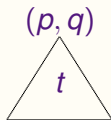
## Constant rules

$$\frac{\exists q. \begin{array}{l} a \rightarrow p \text{ in } A \\ a \rightarrow q \text{ in } B \end{array}}{\text{fail.}}$$



## Final states

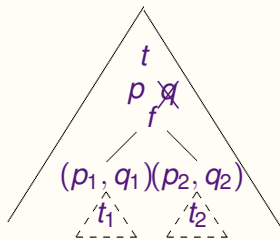
$$\frac{\begin{array}{l} p \text{ final in } A \\ q \text{ not final in } B \end{array}}{\text{fail}:- \text{acc}(p, q).}$$



## Inclusion Failure (cont.)

## Binary rules

$$\frac{\begin{array}{l} f(p_1, p_2) \rightarrow p \text{ in } A \\ \exists q. f(q_1, q_2) \rightarrow q \text{ in } B \end{array}}{\text{fail} :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2).}$$





# Datalog Characterization

## Proposition

*For deterministic  $B$  and productive  $A$ ,  $L(A) \not\subseteq L(B)$  if and only if  $\text{fail} \in \text{lfp}(D_{A,B})$ .*

# Datalog Characterization

## Proposition

For deterministic  $B$  and productive  $A$ ,  $L(A) \not\subseteq L(B)$  if and only if  $fail \in \text{lfp}(D_{A,B})$ .

## Computing failures

- For constant rules and final states:  $O(|A| * |B|)$
- For binary rules:  $O(|A| * |B|^2)$

$$\frac{\begin{array}{l} f(p_1, p_2) \rightarrow p \text{ in } A \\ \exists q. f(q_1, q_2) \rightarrow q \text{ in } B \end{array}}{fail :- acc(p_1, q_1), acc(p_2, q_2).}$$

# Forbidden States

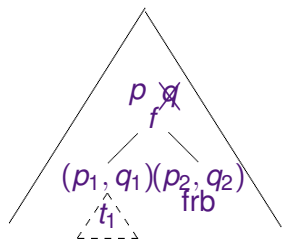
$$\text{frb}(p, q) \Leftrightarrow (\text{acc}(p, q) \Rightarrow L(A) \not\subseteq L(B))$$

- Inferred Datalog clauses:

$$\frac{\begin{array}{l} f(p_1, p_2) \rightarrow p \in \text{rules}(A) \\ \nexists q. f(q_1, q_2) \rightarrow q \in \text{rules}(B) \end{array}}{\begin{array}{l} \text{frb}(p_2, q_2) :- \text{acc}(p_1, q_1). \\ \text{frb}(p_1, q_1) :- \text{acc}(p_2, q_2). \end{array}}$$

$$\frac{p \in \text{states}(A) \quad q \in \text{states}(B)}{\text{fail} :- \text{acc}(p, q), \text{frb}(p, q).}$$

- Number of clauses:  $O(|A| * |B|^2)$



# Grouping Forbidden States

For  $Q \subseteq \text{states}(B)$ :

$$\text{frb}^c(p, Q) \Leftrightarrow \forall q \notin Q, \text{frb}(p, q)$$

- Inferred Datalog clauses:

$$\begin{array}{l}
 f(p_1, p_2) \rightarrow p \in \text{rules}(A) \\
 f(q_1, q_2^1) \rightarrow q^1 \in \text{rules}(B) \\
 \vdots \\
 f(q_1, q_2^n) \rightarrow q^n \in \text{rules}(B)
 \end{array}
 \left. \vphantom{\begin{array}{l} f(p_1, p_2) \rightarrow p \in \text{rules}(A) \\ f(q_1, q_2^1) \rightarrow q^1 \in \text{rules}(B) \\ \vdots \\ f(q_1, q_2^n) \rightarrow q^n \in \text{rules}(B) \end{array}} \right\} \begin{array}{l} \text{all rules} \\ \text{for } q_1 \end{array}$$


---


$$\text{frb}^c(p_2, \{q_2^1, \dots, q_2^n\}) :- \text{acc}(p_1, q_1).$$

# Grouping Forbidden States

For  $Q \subseteq \text{states}(B)$ :

$$\text{frb}^c(p, Q) \Leftrightarrow \forall q \notin Q, \text{frb}(p, q)$$

- Inferred Datalog clauses:

$$\left. \begin{array}{l} f(q_1, q_2^1) \rightarrow q^1 \in \text{rules}(B) \\ f(p_1, p_2) \rightarrow p \in \text{rules}(A) \\ \vdots \\ f(q_1, q_2^n) \rightarrow q^n \in \text{rules}(B) \end{array} \right\} \begin{array}{l} \text{all rules} \\ \text{for } q_1 \end{array}$$

---


$$\begin{array}{l} \text{frb}^c(p_2, \{q_2^1, \dots, q_2^n\}) :- \text{acc}(p_1, q_1). \\ \text{acc}(p_1, q^1) :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2^1). \\ \vdots \\ \text{acc}(p_1, q^n) :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2^n). \end{array}$$

# Grouping Forbidden States

For  $Q \subseteq \text{states}(B)$ :

$$\text{frb}^c(p, Q) \Leftrightarrow \forall q \notin Q, \text{frb}(p, q)$$

- Inferred Datalog clauses:

$$\begin{array}{l}
 f(q_1, q_2^1) \rightarrow q^1 \in \text{rules}(B) \\
 f(p_1, p_2) \rightarrow p \in \text{rules}(A) \\
 \vdots \\
 f(q_1, q_2^n) \rightarrow q^n \in \text{rules}(B)
 \end{array}
 \left. \vphantom{\begin{array}{l} f(q_1, q_2^1) \rightarrow q^1 \in \text{rules}(B) \\ f(p_1, p_2) \rightarrow p \in \text{rules}(A) \\ \vdots \\ f(q_1, q_2^n) \rightarrow q^n \in \text{rules}(B) \end{array}} \right\} \begin{array}{l} \text{all rules} \\ \text{for } q_1 \end{array}$$


---


$$\begin{array}{l}
 \text{frb}^c(p_2, \{q_2^1, \dots, q_2^n\}) :- \text{acc}(p_1, q_1). \\
 \text{acc}(p_1, q^1) :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2^1). \\
 \vdots \\
 \text{acc}(p_1, q^n) :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2^n).
 \end{array}
 \left. \vphantom{\begin{array}{l} \text{frb}^c(p_2, \{q_2^1, \dots, q_2^n\}) :- \text{acc}(p_1, q_1). \\ \text{acc}(p_1, q^1) :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2^1). \\ \vdots \\ \text{acc}(p_1, q^n) :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2^n). \end{array}} \right\} \begin{array}{l} \text{1 clause of size } O(n) \\ n \text{ clauses of size } O(1) \end{array}$$

# Grouping Forbidden States

For  $Q \subseteq \text{states}(B)$ :

$$\text{frb}^c(p, Q) \Leftrightarrow \forall q \notin Q, \text{frb}(p, q)$$

- Inferred Datalog clauses:

$$\begin{array}{l}
 f(q_1, q_2^1) \rightarrow q^1 \in \text{rules}(B) \\
 f(p_1, p_2) \rightarrow p \in \text{rules}(A) \\
 \vdots \\
 f(q_1, q_2^n) \rightarrow q^n \in \text{rules}(B)
 \end{array}
 \left. \vphantom{\begin{array}{l} f(q_1, q_2^1) \rightarrow q^1 \in \text{rules}(B) \\ f(p_1, p_2) \rightarrow p \in \text{rules}(A) \\ \vdots \\ f(q_1, q_2^n) \rightarrow q^n \in \text{rules}(B) \end{array}} \right\} \begin{array}{l} \text{all rules} \\ \text{for } q_1 \end{array}$$


---


$$\begin{array}{l}
 \text{frb}^c(p_2, \{q_2^1, \dots, q_2^n\}) :- \text{acc}(p_1, q_1). \\
 \text{acc}(p_1, q^1) :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2^1). \\
 \vdots \\
 \text{acc}(p_1, q^n) :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2^n).
 \end{array}
 \left. \vphantom{\begin{array}{l} \text{frb}^c(p_2, \{q_2^1, \dots, q_2^n\}) :- \text{acc}(p_1, q_1). \\ \text{acc}(p_1, q^1) :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2^1). \\ \vdots \\ \text{acc}(p_1, q^n) :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2^n). \end{array}} \right\} \begin{array}{l} \text{1 clause of size } O(n) \\ n \text{ clauses of size } O(1) \end{array}$$

- Size of the inferred Datalog program:  $O(|A| * |B|)$

# Testing Failure

## Failure for binary rules

$$\frac{p \in \text{states}(A) \quad q \in \text{states}(B)}{\text{fail} :- \text{acc}(p, q), \text{frb}(p, q).}$$

## Computing forbidden states

$$O(|A| * |B|) \left\{ \begin{array}{l} \text{frb}^c(p, Q_1) \Rightarrow \text{frb}(p, q_1), \forall q_1 \notin Q_1 \\ \vdots \\ \text{frb}^c(p, Q_m) \Rightarrow \text{frb}(p, q_m), \forall q_m \notin Q_m \end{array} \right\} O(|A| * |B|^2)$$

- Cannot enumerate **frb** predicates: need algorithmic solution



# Efficient data structure

- Example:  $\text{states}(B) = \{1, 2, 3, 4\}$

$$\left\{ \begin{array}{l} \text{frb}^c(p, \{1, 2\}) \Rightarrow \text{frb}(p, 3) \text{ and } \text{frb}(p, 4) \\ \text{frb}^c(p, \{2, 4\}) \Rightarrow \text{frb}(p, 1) \text{ and } \text{frb}(p, 3) \end{array} \right\}$$

- Array for state  $p$ :

$q$	1	2	3	4	<i>total</i>
<i>occurrence</i> ( $q$ )	1	2	0	1	2

- $\text{frb}(p, q)$  if  $\text{occurrence}(q) < \text{total}$
- Overall time for testing failure:  $O(|A| * |B|)$

# Testing Inclusion

## Theorem

$L(A) \subseteq L(B)$  can be tested in time  $O(|A| * |B|)$  for binary tree automata  $A$ , deterministic  $B$  over  $\Sigma$ .

# Outline

1 Inclusion in Deterministic Binary Tree Automata

2 Inclusion in Deterministic DTDs

# XML Trees

- Unranked trees
- Binary encoding *via* Currying
- Stepwise tree automata [Carme et al. 2004]

# DTDs to Glushkov Automata

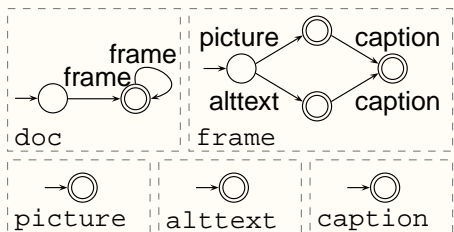
## DTD

```

<!ELEMENT doc      (frame+)>
<!ELEMENT frame    (picture|alttext),caption?>
<!ELEMENT picture  (#PCDATA)>
<!ELEMENT alttext  (#PCDATA)>
<!ELEMENT caption  (#PCDATA)>
  
```

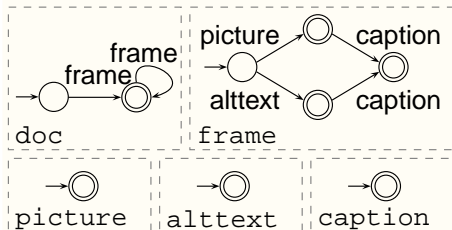
## Glushkov automata

- One-unambiguous regular expressions (W3C)
- To deterministic word automata in  $O(|\Sigma| * |DTD|)$  [Brüggemann-Klein & Wood 1998]

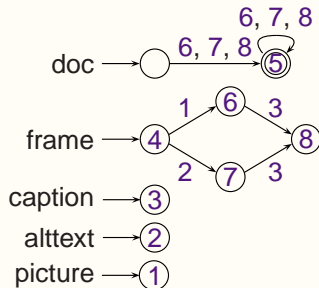


# DTDs to Stepwise Tree Automata

## Glushkov automata to stepwise tree automaton

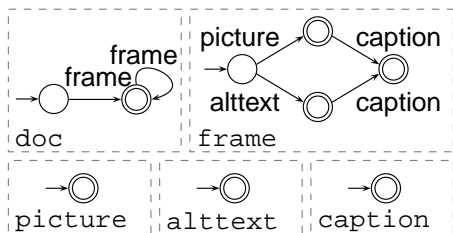


$$O(|\Sigma| * |DTD|)$$

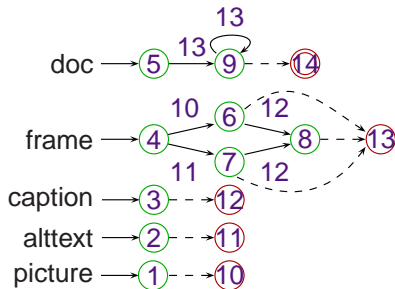


$$O(|\Sigma| * |DTD|^2)$$

## Factorized Tree Automata



$$O(|\Sigma| * |DTD|)$$



$$O(|\Sigma| * |DTD|)$$

# Inclusion into Factorized Tree Automata

- More inference rules to consider
  - Epsilon-rules of factorized automaton
  - Distinct cases for each sort of states
  - Other stuff
- Complexity:  $O(|A| * |\Sigma| * |DTD|)$



# Summary

- Inclusion in deterministic binary tree automata in time  $O(|A| * |B|)$ 
  - Characterization in terms of failure
  - Inferring a ground Datalog program
  - Efficient data structure
- Extension to inclusion in DTDs in time  $O(|A| * |\Sigma| * |DTD|)$ 
  - Unranked tree automata to stepwise tree automata *via* Curryng
  - DTDs to factorized tree automata
  - Adapt inclusion algorithm
  - Lower bound?