

# XML and Databases

Lecture 5  
XML Validation using Automata

Sebastian Maneth  
NICTA and UNSW

CSE@UNSW -- Semester 1, 2010

## Outline

1. Recap: deterministic Reg Expr's / Glushkov Automaton
2. Complexity of DTD validation
3. Beyond DTDs: XML Schema and RELAX NG
4. Static Methods, based on Tree Automata

## Previous Lecture

### XML type definition languages

want to specify a certain subset of XML doc's = a "type" of XML documents

#### Remember

The specification/type definition should be **simple**, so that

- a **validator** can be built automatically (and efficiently)
- the **validator** runs efficient on any XML input

(similar demands as for a **parser**)

→ Type def. language must be SIMPLE!

(similarly: parser generators use EBNF or smaller subclasses: LL / LR)

↖  $O(n^3)$  parsing

## XML Type Definition Languages

**DTD** (Document Type Definition, W3C)  
Originated from SGML. Now part of XML

→ DTD may appear at the beginning of an XML document

**XML Schema** (W3C)  
Now at version 1.1  
HUGE language, many built-in simple types

→ Schemas themselves: written in XML

See the "Schema Primer" at <http://www.w3.org/TR/xml-schema-0/>

**RELAX NG** (Oasis)  
For tree structure definition, more powerful than Schemas&DTDs

Reg Exprs  
must be  
deterministic  
(=1-unambiguous)

same!!

"Unique  
Particle Attribution"

## XML Type Definition Languages

**DTD** (Document Type Definition)

<!DOCTYPE root-element [ doctype declaration ...]>

<!ELEMENT element-name content-model >

#### content-models

- EMPTY
- ANY
- (#PCDATA | element-name\_1 | ... | element-name\_n)\*
- deterministic Reg Expr over element names

<![ATTLIST element-name attr-name attr-type attr-default ...]>

Types: CDATA, (v1|...), ID, IDREFs  
Defaults: #REQUIRED, #IMPLIED, "value", #FIXED

## XML Type Definition Languages

**DTD** (Document Type Definition)

<!DOCTYPE root-element [ doctype declaration ...]>

<!ELEMENT element-name content-model >

#### content-models

- EMPTY
  - ANY
  - (#PCDATA | element-name\_1 | ... | element-name\_n)\*
  - **deterministic Reg Expr**
- Most interesting /  
challenging aspect  
of DTDs

<![ATTLIST element-name attr-name attr-type attr-default ...]>

Types: CDATA, (v1|...), ID, IDREFs  
Defaults: #REQUIRED, #IMPLIED, "value", #FIXED

7

### Summary

In order to check whether a (large) **document** is **valid** wrt to a given **DTD** ("it validates") you need to

→ check if children lists match the given **Reg Expr's**

This can be done *efficiently*, using **finite-automata (FAs)**!

---

To check if a **Reg Expr e** is **allowed in a DTD** we have to construct a particular finite automaton: the **Glushkov automaton**.  
 $\text{Glu}(e)$  must be *deterministic*.

**Note** If  $\text{Glu}(e)$  is *deterministic*, then its size (# transitions) is *linear* in  $\text{size}(e)$ !

8

### Summary

In order to check whether a (large) **document** is **valid** wrt to a given **DTD** ("it validates") you need to

→ check if children lists match the given **Reg Expr's**

This can be done *efficiently*, using **finite-automata (FAs)**!

---

To check if a **Reg Expr e** is **allowed in a DTD** we have to construct a particular finite automaton: the **Glushkov automaton**.  
 $\text{Glu}(e)$  must be *deterministic*.

**Note** If  $\text{Glu}(e)$  is *deterministic*, then its size (# transitions) is *linear* in  $\text{size}(e)$ !

**Question** Can you explain why this is the case?

9

### Summary

In order to check whether a (large) **document** is **valid** wrt to a given **DTD** ("it validates") you need to

→ check if children lists match the given **Reg Expr's**

This can be done *efficiently*, using **finite-automata (FAs)**!

---

To check if a **Reg Expr e** is **allowed in a DTD** we have to construct a particular finite automaton: the **Glushkov automaton**.  
 $\text{Glu}(e)$  must be *deterministic*.

**Note** If  $\text{Glu}(e)$  is *deterministic*, then its size (# transitions) is ~~linear~~ in  $\text{size}(e)$ !

**Question** Can you explain why this is the case? **not correct: linear in  $\text{size}(e) * \#\text{letters}(e)$**

10

### More Notes

(1) From a *deterministic* FA you **cannot** obtain a deterministic (= 1-unambiguous) regular expression!!

Example:  $e = (a|b)^* a (a|b)$  ← NO 1-unambiguous reg exp exists for  $e$

**Question** Can you build a **Deterministic Automaton** for the expression  $e$ ?

11

### More Notes

(1) From a *deterministic* FA you **cannot** obtain a deterministic (= 1-unambiguous) regular expression!!

Example:  $e = (a|b)^* a (a|b)$  ← NO 1-unambiguous reg exp exists for  $e$

**Deterministic Automaton for  $e$ :**  
 E.g., first nondeterministic FA, then **determinize** (subset construction)  
 CAVE: can cause exponential size blow-up!

---

Other important constructions on Finite Automata:

- **Union** (easy)
- **Intersection** (product construction)
- **Complementation**

**Question** Size blow-up for these?

12

### More Notes

(1) From a *deterministic* FA you **cannot** obtain a deterministic (= 1-unambiguous) regular expression!!

Example:  $e = (a|b)^* a (a|b)$  ← NO 1-unambiguous reg exp exists for  $e$

(2)  $\text{Glu}(e)$  is closely related to → **Thomson(e)** [remove  $\epsilon$ -transitions]  
 and to → **Berry/Sethi(e)** [same]  
 and → **Brzowski(e)**

---

To check if a **Reg Expr e** is **allowed in a DTD** we have to construct a particular finite automaton: the **Glushkov automaton**.  
 $\text{Glu}(e)$

**Note** If  $\text{Glu}(e)$  is *deterministic*, then its size is *linear* in  $\text{size}(e) * \#\text{letters}(e)$

For more details: See paper by Brüggemann-Klein. Linked from the course web-page.

Glushkov automaton  $\text{Glu}(e)$

Each letter-position in the Reg Expr  $e$  becomes one state of  $\text{Glu}$ ; plus,  $\text{Glu}$  has one extra begin state.

$\text{FIRST}(e) =$  all possible begin positions of words matching  $e$

e.g.  $\text{FIRST}(R(E|G)(EX)^*) = \{R_1\}$

13



## Glushkov's automaton

$R(E|G)(EX)^*$

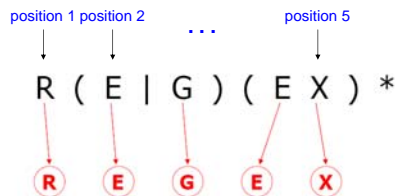
Following slides from: <http://www.cs.ut.ee/~varmo/day-rouge/tammeoja-slides.pdf>

10



## Glushkov's automaton

- Character in RE = **state** in automaton



11



## Glushkov's automaton

- Character in RE = **state** in automaton + one state for the beginning of the RE



12



## Glushkov's automaton

- Character in RE = **state** in automaton + one state for the beginning of the RE
- Transitions** show which characters/positions can precede each other

$R(E|G)(EX)^*$



13



## Glushkov's automaton

- Character in RE = **state** in automaton + one state for the beginning of the RE
- Transitions** show which characters/positions can precede each other



14

## Glushkov's automaton

- Character in RE = **state** in automaton + one state for the beginning of the RE
- Transitions** show which characters/positions can precede each other

$R ( E | G ) ( E X ) ^*$



15

Glushkov automaton  $G(e)$

Each position in the Reg Expr  $e$  becomes one state of  $G$ ; plus,  $G$  has one extra begin state.

$FIRST(e)$  = all possible begin positions of words matching  $e$

e.g.  $FIRST(R (E | G) (EX)^*) = \{ R_1 \}$

$FOLLOW(e, x)$  = all possible positions following position  $x$  in  $e$

e.g.  $FOLLOW(R (E | G) (EX)^*, R_1) = \{ E_2, G_3 \}$

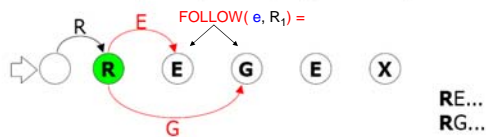
→ From state " $R_1$ ": add E-transition to  $E_2$   
G-transition to  $G_3$

20

## Glushkov's automaton

- Character in RE = **state** in automaton + one state for the beginning of the RE
- Transitions** show which characters/positions can precede each other

$R ( E | G ) ( E X ) ^*$

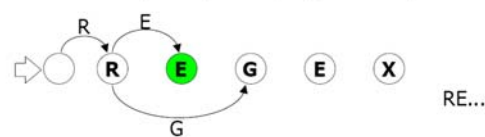


16

## Glushkov's automaton

- Character in RE = **state** in automaton + one state for the beginning of the RE
- Transitions** show which characters/positions can precede each other

$R ( E | G ) ( E X ) ^*$

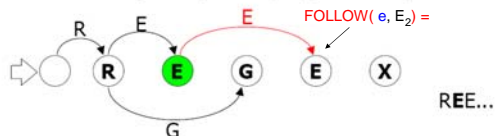


17

## Glushkov's automaton

- Character in RE = **state** in automaton + one state for the beginning of the RE
- Transitions** show which characters/positions can precede each other

$R ( E | G ) ( E X ) ^*$

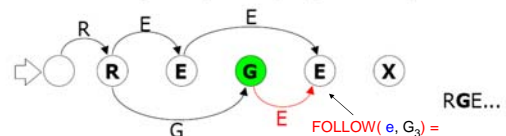


18

## Glushkov's automaton

- Character in RE = **state** in automaton + one state for the beginning of the RE
- Transitions** show which characters/positions can precede each other

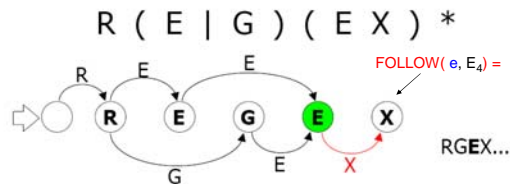
$R ( E | G ) ( E X ) ^*$



19

## Glushkov's automaton

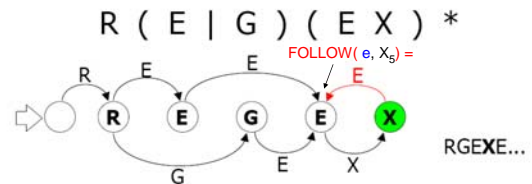
- Character in RE = **state** in automaton + one state for the beginning of the RE
- Transitions** show which characters/positions can precede each other



20

## Glushkov's automaton

- Character in RE = **state** in automaton + one state for the beginning of the RE
- Transitions** show which characters/positions can precede each other



21

Glushkov automaton  $G(e)$

Each position in the Reg Expr  $e$  becomes one state of  $G$ ; plus,  $G$  has one extra begin state.

$FIRST(e)$  = all possible begin positions of words matching  $e$

e.g.  $FIRST(R(E|G)(EX)^*) = \{R_1\}$

$FOLLOW(e, x)$  = all possible positions following position  $x$  in  $e$

e.g.  $FOLLOW(R(E|G)(EX)^*, R_1) = \{E_2, G_3\}$

→ From state " $R_1$ ": add E-transition to  $E_2$   
G-transition to  $G_3$

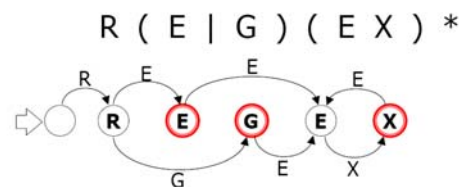
$LAST(e)$  = all possible end positions of words matching  $e$

e.g.  $LAST(R(E|G)(EX)^*) = \{E_2, G_3, X_5\}$

27

## Glushkov's automaton

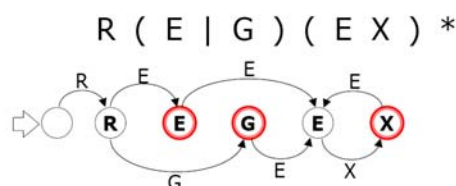
- Character in RE = **state** in automaton + one state for the beginning of the RE
- Transitions** show which characters/positions can precede each other



22

## Glushkov's automaton

- Character in RE = **state** in automaton + one state for the beginning of the RE
- Transitions** show which characters/positions can precede each other



Is this automaton deterministic ??

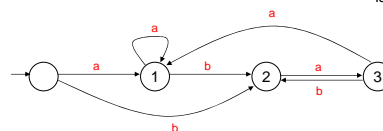
22

Glushkov automaton  $G(e)$

Another example

$(a^* | ba)^*$

This FA is deterministic.



Which of these is deterministic?

- $(ab) | (ac)$
- $a(b | c)$
- $a(a | b)^*ac$

30

31

Glushkov automaton  $G(e)$

Each position in the Reg Expr  $e$  becomes one state of  $G$ ; plus,  $G$  has one extra begin state.

$FIRST(e)$  = all possible *begin* positions of words matching  $e$

e.g.  $FIRST(R(E|G)(EX)^*) = \{R_1\}$

$FOLLOW(e, x)$  = all possible positions *following* position  $x$  in  $e$

$LAST(e)$  = all possible *end* positions of words matching  $e$

---

Naïve implementation:  $O(m^3)$  time, where  $m = size(e)$

(for each position: computing **FOLLOW** goes through every position at each step, needs to compute *union*  $\rightarrow O(m^3)$ )

32

Glushkov automaton  $G(e)$

Each position in the Reg Expr  $e$  becomes one state of  $G$ ; plus,  $G$  has one extra begin state.

$FIRST(e)$  = all possible *begin* positions of words matching  $e$

e.g.  $FIRST(R(E|G)(EX)^*) = \{R_1\}$

$FOLLOW(e, x)$  = all possible positions *following* position  $x$  in  $e$

$LAST(e)$  = all possible *end* positions of words matching  $e$

---

Naïve implementation:  $O(m^3)$  time, where  $m = size(e)$

(for each position: computing **FOLLOW** goes through every position at each step, needs to compute *union*  $\rightarrow O(m^3)$ )

Not really needed. Can be improved to  $O(m^2)$

33

Glushkov automaton  $G(e)$

Each position in the Reg Expr  $e$  becomes one state of  $G$ ; plus,  $G$  has one extra begin state.

$FIRST(e)$  = all possible *begin* positions of words matching  $e$

e.g.  $FIRST(R(E|G)(EX)^*) = \{R_1\}$

$FOLLOW(e, x)$  = all possible positions *following* position  $x$  in  $e$

$LAST(e)$  = all possible *end* positions of words matching  $e$

---

Naïve implementation:  $O(m^3)$  time, where  $m = size(e)$

(for each position: computing **FOLLOW** goes through every position at each step, needs to compute *union*  $\rightarrow O(m^3)$ )

Can be improved to  $O(size(e) + size(G(e)))$   $\leftarrow$  Not really needed. Can be improved to  $O(m^2)$

34

Glushkov automaton  $G(e)$

**Note** If  $G(e)$  is *deterministic*, then its size (# transitions) is *quadratic* in  $size(e)$ !

$\rightarrow$  **Linear** in  $size(e) * \#letters(e)$ , if  $G(e)$  is deterministic!  
 $\rightarrow O(size(e) * \#letters(e))$

Can this be improved?  $\rightarrow$  E.g., to  $O(size(e) * \log(\#letters(e)))$ ?  
 $\rightarrow$  Are there known *lower bounds*?

---

Naïve implementation:  $O(m^3)$  time, where  $m = size(e)$

(for each position: computing **FOLLOW** goes through every position at each step, needs to compute *union*  $\rightarrow O(m^3)$ )

$O(size(e) + size(G(e)))$   $\leftarrow$  Not really needed. Can be improved to  $O(m^2)$

35

To avoid these expensive running times

**DTD** requires that  $FA=G(e)$  must be *deterministic*!

$m = size(e)$   
 $n = length(w)$

**Total Running time**  $O(m + n)$

If  $s = \#letters(e)$  is assumed fixed (not part of the input)  
 Otherwise:  $O(ms + n)$

---

How can you **implement** a regular expression?

Input: Reg Expr  $e$ , string  $w$   
 Question: Does  $w$  match  $e$ ?

*deterministic FA*: run on  $w$  takes time **linear** in  $length(w)$

Unrestricted Reg Expr  $e \rightarrow$

Algorithm

$FA = BuildFA(e)$ ;  
 $DFA = BuildDFA(FA)$ ;

Size of FA is linear in  $size(e)=m$   
 Size of DFA is exponential in  $m$

**Total Running time**  $O(2^m m + n)$   
 $\rightarrow$  Other alternative:  $O(mn)$

36

**Summary**

Deterministic (1-unambiguous) content models give rise to *efficient matching algorithms*.

(they avoid  $O(mn)$  or  $O(2^m m + n)$  complexities)

---

**Disadvantages**

$\rightarrow$  Hard to know whether given reg expr is OK (deterministic)

$\rightarrow$  Det. reg exprs are NOT closed under union. (not so nice.)

**Question** Can you see why?

Maybe, can find det reg exprs, such that their union equals  $(a|b)^* a (a|b)^*$ ?

37

Now that we know how to check all the different **content-models** (in particular det. **Reg Expr's**) how to build a full validator for a given DTD?

elem-name\_1 → RegExpr\_1

elem-name\_2 → RegExpr\_2

...

elem-name\_k → RegExpr\_k

Automata A\_1, A\_2, ..., A\_k

The Validation Problem

Given a DTD T and a document D, is D valid wrt T?

Top-Down Implementation

→ at element node w. label elem-name\_i, run automaton A\_i

→ check attribute constraints

→ check ID/IDREF constraints

Total Running time

(Given A\_1, A\_2, ..., A\_k)

linear in the sum of size(DTD) \* #letters(DTD) and size of document: O( size(D) \* #letters(D) + size(T) )

38

DTDs have the

**"label-guarded subtree exchange"** property:

t1, t2 trees in a DTD language T

v1 node in t1, labeled "lab"

v2 node in t2, labeled "lab"

aka "local"

→ content model only depends on label of parent

trees obtained by exchanging the subtrees rooted at v1 and v2 are also in T

39

**Beyond DTDs**

Often, the expressive power of DTDs is *not sufficient*.

**Problem** each element name has precisely one content-model in a DTD. Would like to distinguish, depending on the context (parent).

dealer

used

new

car

model

year

model

(content model depends on the parent node)

car has different structure, in different contexts.

40

**Beyond DTDs**

Often, the expressive power of DTDs is *not sufficient*.

**Problem** each element name has precisely one content-model in a DTD. Would like to distinguish, depending on the context (parent).

dealer

used

new

car<sub>used</sub>

model

year

model

car<sub>new</sub>

"specialization"

car has different structure, in different contexts.

41

**Specialized DTDs**

dealer → used, new

used → (car<sub>used</sub>)\*

new → (car<sub>new</sub>)\*

car<sub>used</sub> → model, year

car<sub>new</sub> → model

42

**Specialized DTDs**

dealer → used, new

used → (car<sub>used</sub>)\*

new → (car<sub>new</sub>)\*

car<sub>used</sub> → model, year

car<sub>new</sub> → model

New notation. Use *capitalized TYPE Names*

Dealer → dealer [Used, New]

Used → used [(Car<sub>used</sub>)\*]

New → new [(Car<sub>new</sub>)\*]

Car<sub>used</sub> → car [Model, Year]

Car<sub>new</sub> → car [Model]

7

43

New notation. Use *capitalized* TYPE Names

```

Dealer → dealer [Used, New]
Used → used [(Carused)*]
New → new [(Carnew)*]
Carused → car [Model, Year]
Carnew → car [Model]

```

Let us call this new concept a "grammar". the "local" restriction

A grammar *G* is **local**, if for any label[RegExpr<sub>1</sub>], label[RegExpr<sub>2</sub>] present in *G* it holds that RegExpr<sub>1</sub> = RegExpr<sub>2</sub>.

By definition: Every DTD is a **local grammar**, and vice versa.

44

New notation. Use *capitalized* TYPE Names

```

Dealer → dealer [Used, New]
Used → used [(Carused)*]
New → new [(Carnew)*]
Carused → car [Model, Year]
Carnew → car [Model]

```

Alternatively:

Call two TYPE Names *T1* and *T2* "competing" if they have the same element name (but not identical rules)

**Classes of Grammars**

**local** no competing TYPE names! (DTDs)

**single-type** TYPE names in the *same content model* do not compete! (XML Schema's)

**regular** no restriction. (RELAX NG)

45

New notation. Use *capitalized* TYPE Names

```

Dealer → dealer [Used, New]
Used → used [(Carused)*]
New → new [(Carnew)*]
Carused → car [Model, Year]
Carnew → car [Model]

```

**Question** Are there single-type grammars (XML Schemas) which cannot be expressed by local grammars (DTDs).

**Classes of Grammars**

**local** no competing TYPE names! (DTDs)

**single-type** TYPE names in the *same content model* do not compete! (XML Schema's)

**regular** no restriction. (RELAX NG)

46

New notation. Use *capitalized* TYPE Names

```

Person → person [PersonName, Gender, Spouse?, Pet*]
PersonName → name [First, Last]
Pet → pet [Kind, PetName]
PetName → name [#PCDATA]
...

```

but are not in same content model!

**Question** Are there single-type grammars (XML Schemas) which cannot be expressed by local grammars (DTDs). **YES!**

**Classes of Grammars**

**local** no competing TYPE names! (DTDs)

**single-type** TYPE names in the *same content model* do not compete! (XML Schema's)

**regular** no restriction. (RELAX NG)

47

New notation. Use *capitalized* TYPE Names

```

Dealer → dealer [Used, New]
Used → used [(Carused)*]
New → new [(Carnew)*]
Carused → car [Model, Year]
Carnew → car [Model]

```

Through the use of TYPE Names (nonterminals / states) you can distinguish **deep context!**

```

graph TD
    dealer --> used
    dealer --> new
    used --> car1[car]
    new --> car2[car]
    car1 --> model1[model]
    car1 --> year[year]
    car2 --> model2[model]

```

"specialization" through parent

48

New notation. Use *capitalized* TYPE Names

```

Dealer → dealer [Used, New]
Used → used [(Carused)*]
New → new [(Carnew)*]
Carused → car [Model, Year]
Carnew → car [Model]

```

Through the use of TYPE Names (nonterminals / states) you can distinguish **deep context!**

Can we model context that is far away from the specialized node?

```

graph TD
    dealer --> used
    dealer --> new
    used --> car1[car]
    new --> car2[car]
    car1 --> model1[model]
    car1 --> year[year]
    car2 --> model2[model]

```

"specialization" through parent



49

New notation. Use *capitalized* TYPE Names

```

Dealer → dealer [Used, New]
Used → used [(Carused)*]
New → new [(Carnew)*]
Carused → car [Model, Year]
Carnew → car [Model]

```

Through the use of **TYPE Names** (nonterminals / states) you can distinguish **deep context!**

Can we model context that is far away from the specialized node?

**Sure!**

"specialization" through a following node...

50

```

DB → db [Dealer, User]
Dealer → dealer [Used, New]
Used → used [(Carused)*]
New → new [(Carnew)*]
Carused → car [Model, Year]
Carnew → car [Model]

```

Through the use of **TYPE Names** (nonterminals / states) you can distinguish **deep context!**

Can we model context that is far away from the specialized node?

**Sure!**

51

```

DB → db [Dealer, User]
Dealer → dealer [Used, New]
Used → used [(Carused)*]
New → new [(Carnew)*]
Carused → car [Model, Year]
Carnew → car [Model]

```

Through the use of **TYPE Names** (nonterminals / states) you can distinguish **deep context!**

Can we model context that is far away from the specialized node?

**Sure!**

**Question**

Sure this grammar is **not local** (DTD). But, is it **single-type**?

52

**Question** Is this grammar **single-type**?

53

**Classes XML Type Formalisms**

**local** no competing TYPE names! (DTDs)

**single-type** TYPE names in the *same content model* do not compete! (XML Schema's)

**regular** no restriction... (RELAX NG)

**Increasing Expressiveness** of defining sets of trees ("tree languages")

**Questions**

Given two DTDs **D1** and **D2** can we check if

- all documents valid for **D1** are also valid for **D2**? (DTD inclusion problem)
- **D1** and **D2** describe the same set of documents? (DTD equality problem)

Given a **Relax NG** grammar **G**, can we check if

- there exists any document that is valid for **G**? (emptiness problem)
- there is a document valid for **G** and valid for **G2**? (intersection & emptiness)

54

**Classes XML Type Formalisms**

**local** no competing TYPE names! (DTDs)

**single-type** TYPE names in the *same content model* do not compete! (XML Schema's)

**regular** no restriction... (RELAX NG)

**Increasing Expressiveness** of defining sets of trees ("tree languages")

**Questions**

Given two DTDs **D1** and **D2** can we check if

- all documents valid for **D1** are also valid for **D2**? (DTD inclusion problem)
- **D1** and **D2** describe the same set of documents? (DTD equality problem)

Given a **Relax NG** grammar **G**, can we check if

- there exists any document that is valid for **G**? (emptiness problem)
- there is a document valid for **G** and valid for **G2**? (intersection & emptiness)

**If we can do it for regular tree grammars, then also works for single-type/local!!**

55

All of the checks can be done automatically, for **regular tree** grammars!

↑  
equivalent to tree automata

**Tree Automata:** very powerful framework,  
→ Have all the good properties of string automata!  
→ Yet, they are more expressive!

**Questions**

Given two DTDs D1 and D2 can we check if  
→ all documents valid for D1 are also valid for D2? (DTD inclusion problem)  
→ D1 and D2 describe the same set of documents? (DTD equality problem)

Given a Relax NG grammar G, can we check if  
→ there exists any document that is valid for G? (emptiness problem)  
→ there is a document valid for G and valid for G2? (intersection & emptiness)

→ If we can do it for **regular tree** grammars, then also works for single-type/local!!

56

All of the checks can be done automatically, for **regular tree** grammars!

↑  
equivalent to tree automata

**Tree Automata:** very powerful framework,  
→ Have all the good properties of string automata!  
→ Yet, they are more expressive!

**Note**

String automata are **not** sufficient to check DTDs / Schemas!  
Even if we only consider well-bracketed strings!

**Example 1**

c → c[ a, c, b ]  
a → empty  
b → empty  
c → empty

**Example 2**

a → a[ c, a ]  
a → a[ a, b ]  
a / b / c → empty

57

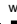

All of the checks can be done automatically, for **regular tree** grammars!

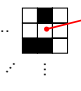
constant memory computation

↑  
equivalent to tree automata

**Finite-state** automata are important:

→ Think you are in a maze, with only fixed memory and you can only read the maze (cannot mark anything).

Model by **finite automaton**. In state q1, (to [N|S|E|W], ) → ( q2, [N|S|E|W] )  
q2, (to [N|S|E|W], ) → ( q3, [N|S|E|W] )  
empty



Can an automaton search the maze?

58

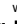

All of the checks can be done automatically, for **regular tree** grammars!

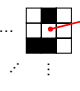
constant memory computation

↑  
equivalent to tree automata

**Finite-state** automata are important:

→ Think you are in a maze, with only fixed memory and you can only read the maze (cannot mark anything).

Model by **finite automaton**. In state q1, (to [N|S|E|W], ) → ( q2, [N|S|E|W] )  
q2, (to [N|S|E|W], ) → ( q3, [N|S|E|W] )  
empty



Can an automaton search the maze?

No!! → need markers ("pebbles").  
How many? 5? 2?

59

All of the checks can be done automatically, for **regular tree** grammars!

constant memory computation

↑  
equivalent to tree automata

**Finite-state** automata are important:

In our context, e.g., for

→ KMP (efficient string matching) [Knuth/Morris/Pratt]  
generalization using automata. Used, e.g., in **grep**

→ **Compression**

→ **Static analysis** of schemas & queries  
(= "everything you can do *before* running over the actual data")

60

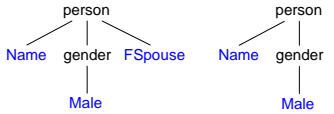
#### 4. Static Methods, based on Tree Automata

Person → MPerson | FPerson  
MPerson → person [Name, gender[Male], FSpouse?]  
FPerson → person [Name, gender[Female], MSpouse?]

**Regular Tree Grammar**

Rules of the form **TypeName** → Tree

Leaves may be labeled by **TypeNames**



Alternatively, regular tree languages are defined by **Tree Automata**.

state, element-name → state1, state2

conventionally, defined for **binary/ranked** trees.

#### 4. Static Methods, based on Tree Automata

61

Given grammars D1 and D2 can we check if  
 → all documents valid for D1 are also valid for D2? (inclusion problem)  
 → D1 and D2 describe the same set of documents? (equality problem)  
 → does there exists any document that is valid for D1? (emptiness problem)  
 → there is a document valid for D1 \*and\* valid for D2? (intersection & emptiness)

ALL these checks are possible for **regular tree grammars!!**

→ hence, they are also solvable for DTDs / XML Schemas / RELAX NG's

- (1) use binary tree encodings
- (2) translate XML Type Definition to a Tree Grammar (easy)

Alternatively, regular tree languages are defined by **Tree Automata**.

state, element-name → state1, state2 ← conventionally, defined for binary/ranked trees.

#### 4. Static Methods, based on Tree Automata

62

Given grammars D1 and D2 can we check if  
 → all documents valid for D1 are also valid for D2? (*inclusion problem*)  
 → D1 and D2 describe the same set of documents? (equality problem)  
 → does there exists any document that is valid for D1? (emptiness problem)  
 → there is a document valid for D1 \*and\* valid for D2? (intersection & emptiness)

ALL these checks are possible for **regular tree grammars!!**

→ The checks above give rise to very powerful optimization procedures for XML Databases!

For example:  
 documents d\_1, d\_2, ..., d\_n are valid for your schema "Small\_xhtml".

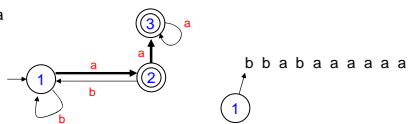
Are they also valid for schema XHTML?

→ Check *inclusion problem* for Small\_html and XHTML!

Automata

on words

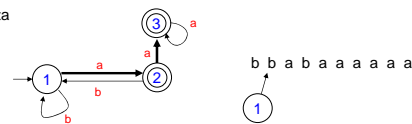
63



Automata

on words

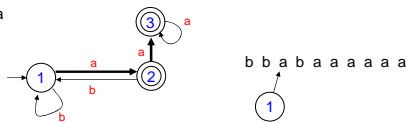
64



Automata

on words

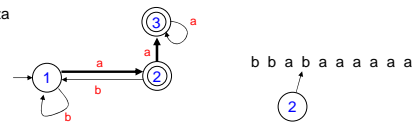
65



Automata

on words

66



Automata 67 on words

Expressiveness deterministic = nondeterministic word is **accepted** by the automaton  
left-to-right = right-to-left

Automata 68 on words

Expressiveness deterministic = nondeterministic word is **accepted** by the automaton  
left-to-right = right-to-left

Automata on trees 1. bottom-up LABEL( state1, state2 ) → state

$0() \rightarrow F$   
 $1() \rightarrow T$   
 $OR(F, F) \rightarrow F$   
 $OR(F, T) \rightarrow T$   
...

Automata 69 on words

Expressiveness deterministic = nondeterministic word is **accepted** by the automaton  
left-to-right = right-to-left

Automata on trees 1. bottom-up LABEL( state1, state2 ) → state

$0() \rightarrow F$   
 $1() \rightarrow T$   
 $OR(F, F) \rightarrow F$   
 $OR(F, T) \rightarrow T$   
...

Automata 70 on words

Expressiveness deterministic = nondeterministic word is **accepted** by the automaton  
left-to-right = right-to-left

Automata on trees 1. bottom-up LABEL( state1, state2 ) → state

$0() \rightarrow F$   
 $1() \rightarrow T$   
 $OR(F, F) \rightarrow F$   
 $OR(F, T) \rightarrow T$   
...  
 $AND(T, F) \rightarrow F$

Automata 71 on words

Expressiveness deterministic = nondeterministic word is **accepted** by the automaton  
left-to-right = right-to-left

Automata on trees 1. bottom-up LABEL( state1, state2 ) → state

$0() \rightarrow F$   
 $1() \rightarrow T$   
 $OR(F, F) \rightarrow F$   
 $OR(F, T) \rightarrow T$   
...  
 $AND(T, F) \rightarrow F$

Accepting States = { T }

tree is **accepted** by the automaton

Automata 72 on words

Expressiveness deterministic = nondeterministic word is **accepted** by the automaton  
left-to-right = right-to-left

Automata on trees 1. bottom-up LABEL( state1, state2 ) → state

$0() \rightarrow F$   
 $1() \rightarrow T$   
 $OR(F, F) \rightarrow F$   
 $OR(F, T) \rightarrow T$   
...  
 $AND(T, F) \rightarrow F$

Accepting States = { T }

This automaton is **deterministic**.

nondeterminism  
LABEL( st1, st2 ) → { st3, ... }

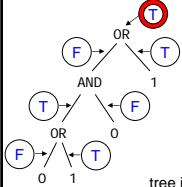
tree is **accepted** by the automaton

73

### Question

How much memory do you need exactly, to run such a bottom-up tree automaton?

#### Automata on trees



1. *bottom-up* LABEL( state1, state2 ) → state

0() → F      Accepting States = { T }

1() → T

OR( F, F ) → F

OR( F, T ) → T

...      This automaton is deterministic.

AND( T, F ) → F

nondeterminism LABEL( st1, st2 ) → { st3, ... }

tree is **accepted** by the automaton

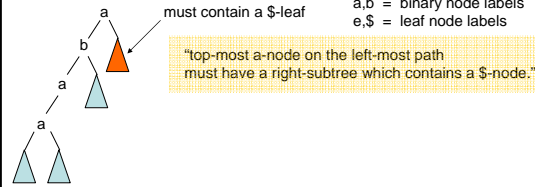
74

Similarly as for word automata:

For every **nondeterministic** bottom-up tree automaton there is an equivalent **deterministic** bottom-up tree automaton.

Again, the construction can cause exponential size blow-up.

2. *top-down* state, LABEL → (state1, state2)



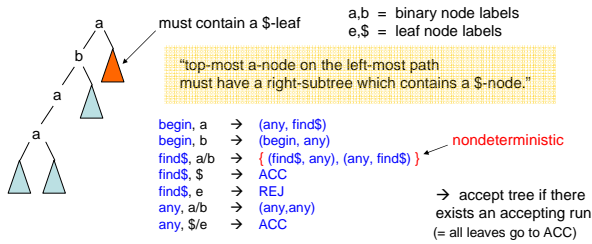
75

Similarly as for word automata:

For every **nondeterministic** bottom-up tree automaton there is an equivalent **deterministic** bottom-up tree automaton.

Again, the construction can cause exponential size blow-up.

2. *top-down* state, LABEL → (state1, state2)



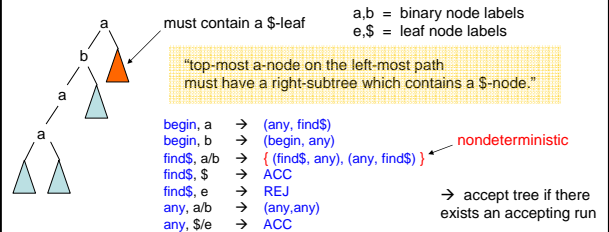
76

For every **nondeterministic** bottom-up tree automaton there is an equivalent **deterministic** bottom-up tree automaton.

### Question

Can you find an equivalent *bottom-up* automaton for this example?

2. *top-down* state, LABEL → (state1, state2)

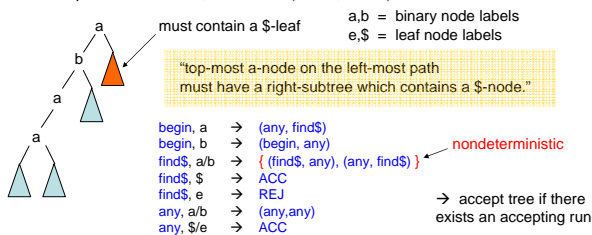


77

For every **nondeterministic** bottom-up tree automaton  
→ there is an equivalent **deterministic** bottom-up tree automaton, and  
→ there is an equivalent **nondeterministic top-down** tree automaton.

→ Yes! you can... ☺

2. *top-down* state, LABEL → (state1, state2)



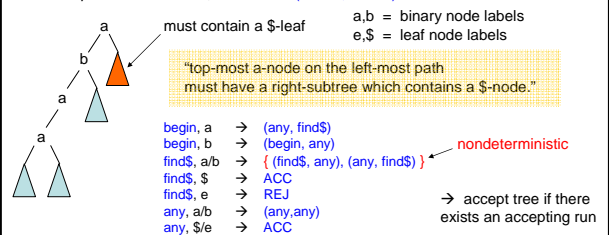
78

For every **nondeterministic** bottom-up tree automaton  
→ there is an equivalent **deterministic** bottom-up tree automaton, and  
→ there is an equivalent **nondeterministic top-down** tree automaton.

### Question

Is there an equivalent **deterministic top-down** automaton??

2. *top-down* state, LABEL → (state1, state2)

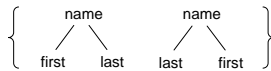


For every **nondeterministic** bottom-up tree automaton  
 → there is an equivalent **deterministic** bottom-up tree automaton, and  
 → there is an equivalent **nondeterministic top-down** tree automaton.

#### Question

Is there an equivalent **deterministic top-down** automaton??

→ NO! ☹ -- not a good model..



This set of two trees canNOT be recognized  
 by any **deterministic top-down** tree automaton!!

**Why?**

For every **nondeterministic** bottom-up tree automaton  
 → there is an equivalent **deterministic** bottom-up tree automaton, and  
 → there is an equivalent **nondeterministic top-down** tree automaton.

#### Question

Is there an equivalent **deterministic top-down** automaton??

→ NO! ☹ -- not a good model..

#### Questions

What about **local** tree languages (defined by DTDs).

→ Can they be accepted by **deterministic top-down** automata?

What about **single-type** tree languages (defined by XML Schema's)

→ Can they be accepted by **deterministic top-down** automata?

For every **nondeterministic** bottom-up tree automaton  
 → there is an equivalent **deterministic** bottom-up tree automaton, and  
 → there is an equivalent **nondeterministic top-down** tree automaton.

#### Question

Is there an equivalent **deterministic top-down** automaton??

→ NO! ☹ -- not a good model..

#### Questions

What about **local** tree languages (defined by DTDs).

→ Can they be accepted by **deterministic top-down** automata?

What about **single-type** tree languages (defined by XML Schema's)

→ Can they be accepted by **deterministic top-down** automata?

**Yes!**

Hence, there is **no DTD / Schema** for { name[first,last], name[last,first] }

For every **nondeterministic** bottom-up tree automaton  
 → there is an equivalent **deterministic** bottom-up tree automaton, and  
 → there is an equivalent **nondeterministic top-down** tree automaton.

#### Question

Is there an equivalent **deterministic top-down** automaton??

→ NO! ☹ -- not a good model..

Nevertheless, **XML Schemas** are a subclass of  
**deterministic top-down** automata.

#### Questions

What is the reasoning behind this?

Similarity to restriction to deterministic reg. expressions?

Recall: Deterministic (1-unambiguous) content models give rise to  
*efficient matching algorithms.*

(they avoid  $O(mn)$   
 $O(2^m + n)$  complexities)

For every **nondeterministic** bottom-up tree automaton  
 → there is an equivalent **deterministic** bottom-up tree automaton, and  
 → there is an equivalent **nondeterministic top-down** tree automaton.

#### Question

Is there an equivalent **deterministic top-down** automaton??

→ NO! ☹ -- not a good model..

Nevertheless, **XML Schemas** are a subclass of  
**deterministic top-down** automata.

#### Questions

What is the reasoning behind this?

Similarity to restriction to deterministic reg. expressions?

In total, given a **DTD D**, we can build one  
**deterministic top-down automaton** of size  $O(\text{size}(D) * \#\text{letters}(D))$

Thus, matching time is inside  $O(m^2 + n)$

For every deterministic bottom-up tree automaton  
 there exists a **minimal unique** equivalent one!

→ Equivalence is decidable

In fact, YOU have already produced  
 minimal bottom-up tree automata!

The **minimal DAG** of a tree  $t$  can be seen as the minimal unique  
 deterministic BU tree automaton that only accepts the tree  $t$ .

For every deterministic bottom-up tree automaton there exists a **minimal unique** equivalent one!

→ Equivalence is decidable

In fact, YOU have already produced minimal bottom-up tree automata!

The **minimal DAG** of a tree  $t$  can be seen as the minimal unique tree automaton that only accepts the tree  $t$ .

#### Question

Given deterministic BU tree automaton.  
How expensive (complexity) to find unique minimal one?

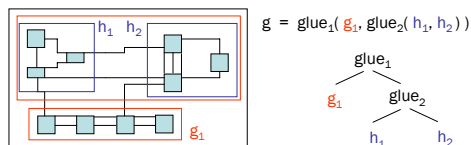
As for DFAs: merge equivalent states  
Typically: quadratic running time. → [Hopcroft's algorithm,  \$O\(n \log n\)\$](#)

More generally, once we have partitioned on a block and an input symbol, we need never partition on that block and input symbol again until the block is split and then we need only partition on one of the two subblocks. Since the time needed to **partition on a block is proportional to the transitions into the block and since we can always select the half with fewer transitions**, the total number of steps in the algorithm is bounded by  $n \log n$ .

See  
<ftp://db.stanford.edu/pub/cstr/reports/cs/tr/71/190/CS-TR-71-190.pdf>

Tree Automata are a very useful concept in CS!

→ Heavily used in **verification**  
"Derive a property of a complex object from the properties of its constituents..."



→ Do all graphs / chip-layouts produced in this way, have property P?

Use the hierarchical construction history of an object, in order to work on a "parse" tree instead of a complex graph.  
From there, use tree automata. ☺

Many NP-complete graph problems become tractable on "bounded-treewidth" graphs!

**XML** Tree Automata play crucial role for

→ Efficient validators against **XML Types**

→ Optimizations If doc1 is of TYPE1, then no need to validate against TYPE2, if we know TYPE2 included in TYPE1

- if only "slightly different" then only need to validate "there"
- incremental validation against updates
- etc, etc.

→ Efficient **query evaluators**, use richer automata which can select nodes and produce query answers

→ Optimizations If answer of QUERY1 is in cache, then no need to evaluate QUERY2, if "included" in QUERY1.

- if every possible answer set to QUERY1 (of TYPE X) is EMPTY, then no need to evaluate on the real data!

→ **XML Type Checking for Programming Languages**

#### The Future

In 5-10 years from now: ☺

You can write a function in Programming Language X

```
Function foo(XML document D: TYPE1): TYPE2
{
  traverse D
  & compute output;
  .
  .
  return output
}
```

Compiler (**XML Type Checker**) will complain, if your function does not compute documents of **TYPE2**.

→ If no complaint, then **guaranteed**:  
ALL outputs are ALWAYS of correct type!!

#### The Future

In 5-10 years from now: ☺

You can write a function in Programming Language X

```
Function foo(XML document D: TYPE1): TYPE2
{
  traverse D
  & compute output;
  .
  .
  return output
}
```

Compiler (**XML Type Checker**) will complain, if your function does not compute documents of **TYPE2**.

→ If no complaint, then correct type **guaranteed**.

Compilers will **have** to be able to give *static guarantees* about input/output behaviour of program!

Experimental PL's  
In this direction:  
→ CDuce  
→ XDuce

END  
Lecture 5