

# Language, Automata and Logic for Finite Trees

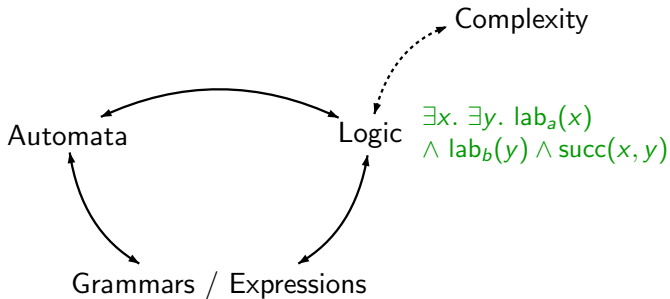
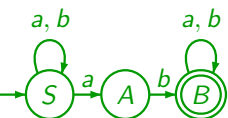
Olivier Gauwin

UMons

Feb/March 2010

# Languages, Automata, Logic

Example for regular word languages:  $\Sigma^*.a.b.\Sigma^*$



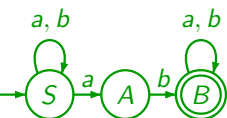
$\exists x. \exists y. lab_a(x)$   
 $\wedge lab_b(y) \wedge succ(x, y)$

$$\left\{ \begin{array}{llll} S \rightarrow aS & S \rightarrow aB & X \rightarrow aX & X \rightarrow \epsilon \\ S \rightarrow bS & B \rightarrow bX & X \rightarrow bX & \end{array} \right.$$

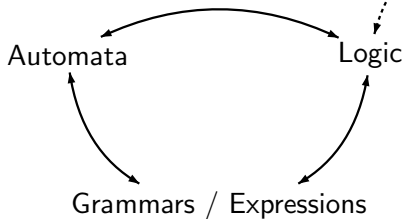
# Languages, Automata, Logic

Example for regular word languages:  $\Sigma^*.a.b.\Sigma^*$

Example for regular tree languages: all trees with an  $a$ -node having a  $b$ -child



$$\left\{ \begin{array}{l} \neg(S, X) \rightarrow S \\ \neg(X, S) \rightarrow S \\ a(B, X) \rightarrow S \\ a(X, B) \rightarrow S \\ \dots \\ b \rightarrow B \\ \dots \end{array} \right.$$



$$\begin{array}{l} \exists x. \exists y. lab_a(x) \\ \wedge lab_b(y) \wedge succ(x, y) \\ \\ \exists x. \exists y. lab_a(x) \\ \wedge lab_b(y) \wedge child(x, y) \end{array}$$

$$\left\{ \begin{array}{llll} S \rightarrow aS & S \rightarrow aB & X \rightarrow aX & X \rightarrow \epsilon \\ S \rightarrow bS & B \rightarrow bX & X \rightarrow bX & \end{array} \right.$$

$$\left\{ \begin{array}{llll} S \rightarrow \neg(S, X) & S \rightarrow a(B, X) & B \rightarrow b(X, X) & X \rightarrow \neg(X, X) \\ S \rightarrow \neg(X, S) & S \rightarrow a(X, B) & B \rightarrow b & X \rightarrow \_ \end{array} \right.$$

# References

The main reference for this talk is the *TATA* book [CDG<sup>+</sup>07]:

*Tree Automata, Techniques and Applications*

by

Hubert Comon, Max Dauchet, Rémi Gilleron, Christof Löding,  
Florent Jacquemard, Denis Lugier, Sophie Tison, Marc Tommasi.

Other references will be mentioned progressively.

## 1 Ranked Trees

- Trees on Ranked Alphabet
- Tree Automata
- Tree Grammars
- Logic

## 2 Unranked Trees

- Unranked Trees
- Automata
- Logic

## 1 Ranked Trees

- Trees on Ranked Alphabet
- Tree Automata
- Tree Grammars
- Logic

## 2 Unranked Trees

- Unranked Trees
- Automata
- Logic

# Trees on Ranked Alphabet

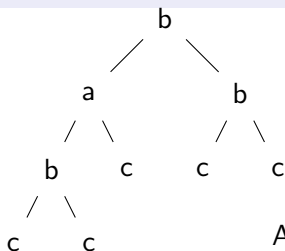
## Ranked alphabet

Ranked alphabet  $\Sigma_r$  = finite alphabet  $\Sigma = \{a, b, c\}$  + arity function  $\text{ar} : \Sigma \rightarrow \mathbb{N}$

$\text{ar}(a) = 2$   
 $\text{ar}(b) = 2$   
 $\text{ar}(c) = 0$

## Ranked trees over $\Sigma_r$

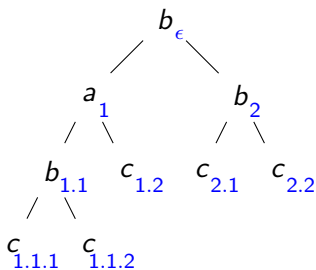
$\mathcal{T}_{\Sigma_r}$ , the set of **ranked trees**, is the smallest set of terms  $f(t_1, \dots, t_k)$  such that:  $f \in \Sigma_r$ ,  $k = \text{ar}(f)$ , and  $t_i \in \mathcal{T}_{\Sigma_r}$  for all  $1 \leq i \leq k$ .



A **tree language**  $T$  is a set of trees:  $T \subseteq \mathcal{T}_{\Sigma_r}$ .

## Trees as relational structures

We will sometimes consider a ranked tree  $t$  as a **relational structure** ( $\text{nodes}^t, \{\text{lab}_a^t, \text{lab}_b^t, \text{lab}_c^t, \text{ch}_1^t, \text{ch}_2^t\}$ ).



- $\text{lab}_a^t = \{1\}$
- $\text{lab}_b^t = \{\epsilon, 1.1, 2\}$
- $\text{lab}_c^t = \{1.1.1, 1.1.2, 1.2, 2.1, 2.2\}$
- $\text{ch}_1^t = \{(\epsilon, 1), (1, 1.1), (2, 2.1) \dots\}$
- $\text{ch}_2^t = \{(\epsilon, 2), (1, 1.2), (2, 2.2) \dots\}$

For convenience we write  $\text{lab}(\pi)$  for the label of node  $\pi$ .



## 1 Ranked Trees

- Trees on Ranked Alphabet
- **Tree Automata**
- Tree Grammars
- Logic

## 2 Unranked Trees

- Unranked Trees
- Automata
- Logic

# Tree Automata

## over Ranked Trees

### Definition

A tree automaton (TA) over  $\Sigma_r = (\Sigma, \text{ar})$  is a tuple  $A = (Q, F, \Delta, \Sigma_r)$  where:

- $Q$  is a finite set of **states**,
- $F \subseteq Q$  a set of **final** states,
- $\Delta$  are **rules** of type:  $a(q_1, \dots, q_k) \rightarrow q$   
with  $a \in \Sigma$ ,  $k = \text{ar}(a)$  and  $q, q_1, \dots, q_k \in Q$

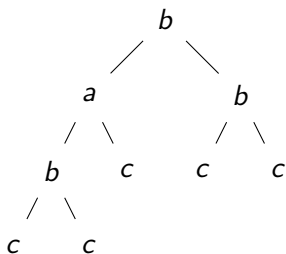
### Runs

A **run**  $\rho$  of  $A$  on  $t$  is a function  $\rho : \text{nodes}^t \rightarrow Q$  such that:  
if  $\pi \in \text{nodes}^t$  with children  $\pi_1, \dots, \pi_k$  and label  $a$  then  
 $a(\rho(\pi_1), \dots, \rho(\pi_k)) \rightarrow \rho(\pi) \in \Delta$

# Bottom-up vs top-down

## Bottom-up view

$a(q_1, \dots, q_k) \rightarrow q$  is a **bottom-up** point of view:

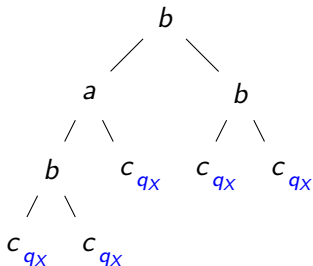


Rules:  $\left\{ \begin{array}{l} a(q_S, q_X) \rightarrow q_S \\ a(q_X, q_S) \rightarrow q_S \\ b(q_S, q_X) \rightarrow q_S \\ b(q_X, q_S) \rightarrow q_S \\ a(q_B, q_X) \rightarrow q_S \\ a(q_X, q_B) \rightarrow q_S \\ b(q_X, q_X) \rightarrow q_B \\ a(q_X, q_X) \rightarrow q_X \\ b(q_X, q_X) \rightarrow q_X \\ c \rightarrow q_X \end{array} \right.$

# Bottom-up vs top-down

## Bottom-up view

$a(q_1, \dots, q_k) \rightarrow q$  is a **bottom-up** point of view:



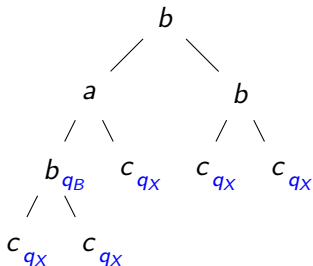
Rules:

$$\left\{ \begin{array}{l} a(q_S, q_X) \rightarrow q_S \\ a(q_X, q_S) \rightarrow q_S \\ b(q_S, q_X) \rightarrow q_S \\ b(q_X, q_S) \rightarrow q_S \\ a(q_B, q_X) \rightarrow q_S \\ a(q_X, q_B) \rightarrow q_S \\ b(q_X, q_X) \rightarrow q_B \\ a(q_X, q_X) \rightarrow q_X \\ b(q_X, q_X) \rightarrow q_X \\ c \rightarrow q_X \end{array} \right.$$

# Bottom-up vs top-down

## Bottom-up view

$a(q_1, \dots, q_k) \rightarrow q$  is a **bottom-up** point of view:



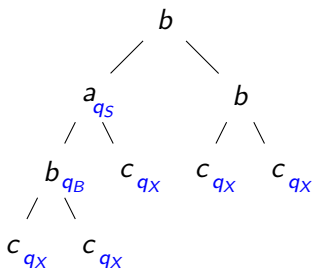
Rules:

$$\left\{ \begin{array}{l} a(q_S, q_X) \rightarrow q_S \\ a(q_X, q_S) \rightarrow q_S \\ b(q_S, q_X) \rightarrow q_S \\ b(q_X, q_S) \rightarrow q_S \\ a(q_B, q_X) \rightarrow q_S \\ a(q_X, q_B) \rightarrow q_S \\ b(q_X, q_X) \rightarrow q_B \\ a(q_X, q_X) \rightarrow q_X \\ b(q_X, q_X) \rightarrow q_X \\ c \rightarrow q_X \end{array} \right.$$

# Bottom-up vs top-down

## Bottom-up view

$a(q_1, \dots, q_k) \rightarrow q$  is a **bottom-up** point of view:



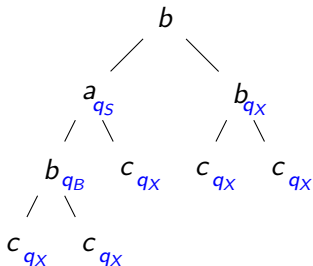
Rules:

$$\left\{ \begin{array}{l} a(q_S, q_X) \rightarrow q_S \\ a(q_X, q_S) \rightarrow q_S \\ b(q_S, q_X) \rightarrow q_S \\ b(q_X, q_S) \rightarrow q_S \\ a(q_B, q_X) \rightarrow q_S \\ a(q_X, q_B) \rightarrow q_S \\ b(q_X, q_X) \rightarrow q_B \\ a(q_X, q_X) \rightarrow q_X \\ b(q_X, q_X) \rightarrow q_X \\ c \rightarrow q_X \end{array} \right.$$

# Bottom-up vs top-down

## Bottom-up view

$a(q_1, \dots, q_k) \rightarrow q$  is a **bottom-up** point of view:



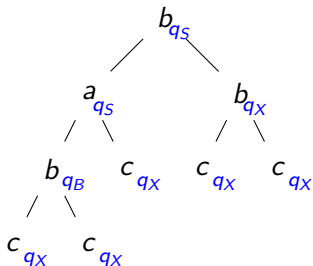
Rules:

$$\left\{ \begin{array}{l} a(q_S, q_X) \rightarrow q_S \\ a(q_X, q_S) \rightarrow q_S \\ b(q_S, q_X) \rightarrow q_S \\ b(q_X, q_S) \rightarrow q_S \\ a(q_B, q_X) \rightarrow q_S \\ a(q_X, q_B) \rightarrow q_S \\ b(q_X, q_X) \rightarrow q_B \\ a(q_X, q_X) \rightarrow q_X \\ b(q_X, q_X) \rightarrow q_X \\ c \rightarrow q_X \end{array} \right.$$

# Bottom-up vs top-down

## Bottom-up view

$a(q_1, \dots, q_k) \rightarrow q$  is a **bottom-up** point of view:



Rules:

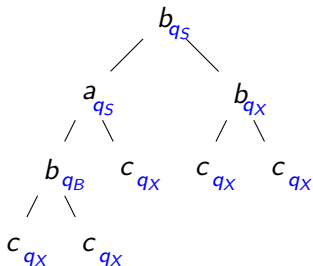
$$\left\{ \begin{array}{l} a(q_S, q_X) \rightarrow q_S \\ a(q_X, q_S) \rightarrow q_S \\ \mathbf{b(q_S, q_X) \rightarrow q_S} \\ b(q_X, q_S) \rightarrow q_S \\ a(q_B, q_X) \rightarrow q_S \\ a(q_X, q_B) \rightarrow q_S \\ b(q_X, q_X) \rightarrow q_B \\ a(q_X, q_X) \rightarrow q_X \\ b(q_X, q_X) \rightarrow q_X \\ c \rightarrow q_X \end{array} \right.$$



# Bottom-up vs top-down

## Bottom-up view

$a(q_1, \dots, q_k) \rightarrow q$  is a **bottom-up** point of view:



Rules:

$$\left\{ \begin{array}{l} a(q_S, q_X) \rightarrow q_S \\ a(q_X, q_S) \rightarrow q_S \\ b(q_S, q_X) \rightarrow q_S \\ b(q_X, q_S) \rightarrow q_S \\ a(q_B, q_X) \rightarrow q_S \\ a(q_X, q_B) \rightarrow q_S \\ b(q_X, q_X) \rightarrow q_B \\ a(q_X, q_X) \rightarrow q_X \\ b(q_X, q_X) \rightarrow q_X \\ c \rightarrow q_X \end{array} \right.$$

A run of  $A$  on  $t$  is **accepting** if  $\rho(\epsilon) \in F$ .

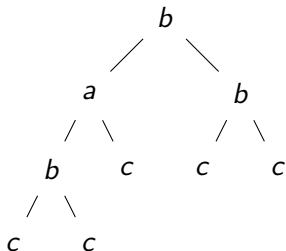
$\mathcal{L}(A) = \{t \mid \text{there exists an accepting run } \rho \text{ of } A \text{ on } t\}$

# Bottom-up vs top-down

## Top-down view

We could have written rules this way:  $a(q) \rightarrow (q_1, \dots, q_k)$ , and name  $F$  the **initial** states.

This corresponds to a **top-down** definition:



Initial:  $\{q_S\}$

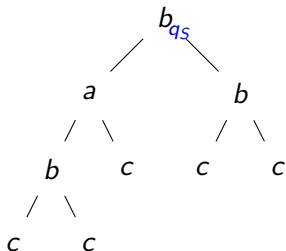
Rules:  $\left\{ \begin{array}{l} a(q_S) \rightarrow (q_S, q_X) \\ a(q_S) \rightarrow (q_X, q_S) \\ b(q_S) \rightarrow (q_S, q_X) \\ b(q_S) \rightarrow (q_X, q_S) \\ a(q_S) \rightarrow (q_B, q_X) \\ a(q_S) \rightarrow (q_X, q_B) \\ b(q_B) \rightarrow (q_X, q_X) \\ a(q_X) \rightarrow (q_X, q_X) \\ b(q_X) \rightarrow (q_X, q_X) \\ c(q_X) \end{array} \right.$

# Bottom-up vs top-down

## Top-down view

We could have written rules this way:  $a(q) \rightarrow (q_1, \dots, q_k)$ , and name  $F$  the **initial** states.

This corresponds to a **top-down** definition:



Initial:  $\{q_S\}$

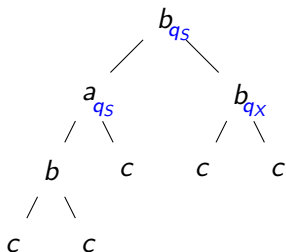
Rules:  $\left\{ \begin{array}{l} a(q_S) \rightarrow (q_S, q_X) \\ a(q_S) \rightarrow (q_X, q_S) \\ b(q_S) \rightarrow (q_S, q_X) \\ b(q_S) \rightarrow (q_X, q_S) \\ a(q_S) \rightarrow (q_B, q_X) \\ a(q_S) \rightarrow (q_X, q_B) \\ b(q_B) \rightarrow (q_X, q_X) \\ a(q_X) \rightarrow (q_X, q_X) \\ b(q_X) \rightarrow (q_X, q_X) \\ c(q_X) \end{array} \right.$

# Bottom-up vs top-down

## Top-down view

We could have written rules this way:  $a(q) \rightarrow (q_1, \dots, q_k)$ , and name  $F$  the **initial** states.

This corresponds to a **top-down** definition:



Initial:  $\{q_S\}$

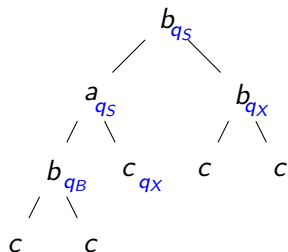
Rules:  $\left\{ \begin{array}{l} a(q_S) \rightarrow (q_S, q_X) \\ a(q_S) \rightarrow (q_X, q_S) \\ b(q_S) \rightarrow (q_S, q_X) \\ b(q_S) \rightarrow (q_X, q_S) \\ a(q_S) \rightarrow (q_B, q_X) \\ a(q_S) \rightarrow (q_X, q_B) \\ b(q_B) \rightarrow (q_X, q_X) \\ a(q_X) \rightarrow (q_X, q_X) \\ b(q_X) \rightarrow (q_X, q_X) \\ c(q_X) \end{array} \right.$

# Bottom-up vs top-down

## Top-down view

We could have written rules this way:  $a(q) \rightarrow (q_1, \dots, q_k)$ , and name  $F$  the **initial** states.

This corresponds to a **top-down** definition:



Initial:  $\{q_S\}$

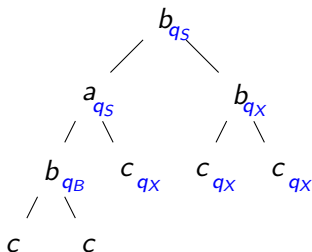
Rules:  $\left\{ \begin{array}{l} a(q_S) \rightarrow (q_S, q_X) \\ a(q_S) \rightarrow (q_X, q_S) \\ b(q_S) \rightarrow (q_S, q_X) \\ b(q_S) \rightarrow (q_X, q_S) \\ a(q_S) \rightarrow (q_B, q_X) \\ a(q_S) \rightarrow (q_X, q_B) \\ b(q_B) \rightarrow (q_X, q_X) \\ a(q_X) \rightarrow (q_X, q_X) \\ b(q_X) \rightarrow (q_X, q_X) \\ c(q_X) \end{array} \right.$

# Bottom-up vs top-down

## Top-down view

We could have written rules this way:  $a(q) \rightarrow (q_1, \dots, q_k)$ , and name  $F$  the **initial** states.

This corresponds to a **top-down** definition:



Initial:  $\{q_S\}$

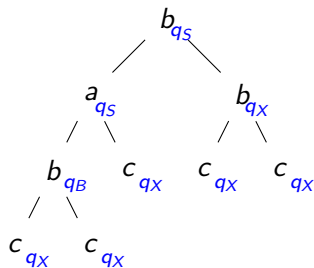
Rules:  $\left\{ \begin{array}{l} a(q_S) \rightarrow (q_S, q_X) \\ a(q_S) \rightarrow (q_X, q_S) \\ b(q_S) \rightarrow (q_S, q_X) \\ b(q_S) \rightarrow (q_X, q_S) \\ a(q_S) \rightarrow (q_B, q_X) \\ a(q_S) \rightarrow (q_X, q_B) \\ b(q_B) \rightarrow (q_X, q_X) \\ a(q_X) \rightarrow (q_X, q_X) \\ b(q_X) \rightarrow (q_X, q_X) \\ c(q_X) \end{array} \right.$

# Bottom-up vs top-down

## Top-down view

We could have written rules this way:  $a(q) \rightarrow (q_1, \dots, q_k)$ , and name  $F$  the **initial** states.

This corresponds to a **top-down** definition:



Initial:  $\{q_S\}$

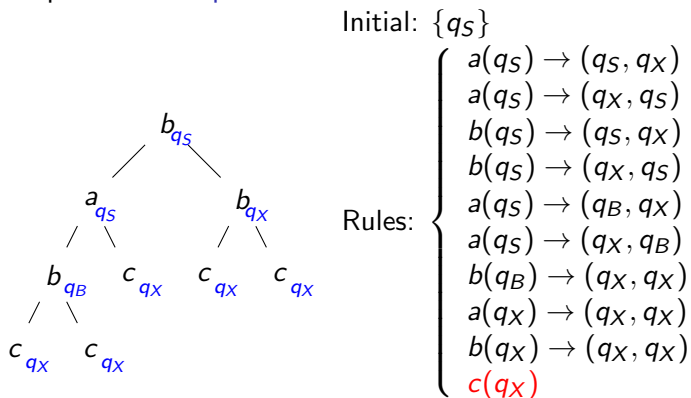
Rules:  $\left\{ \begin{array}{l} a(q_S) \rightarrow (q_S, q_X) \\ a(q_S) \rightarrow (q_X, q_S) \\ b(q_S) \rightarrow (q_S, q_X) \\ b(q_S) \rightarrow (q_X, q_S) \\ a(q_S) \rightarrow (q_B, q_X) \\ a(q_S) \rightarrow (q_X, q_B) \\ \mathbf{b(q_B) \rightarrow (q_X, q_X)} \\ a(q_X) \rightarrow (q_X, q_X) \\ b(q_X) \rightarrow (q_X, q_X) \\ c(q_X) \end{array} \right.$

# Bottom-up vs top-down

## Top-down view

We could have written rules this way:  $a(q) \rightarrow (q_1, \dots, q_k)$ , and name  $F$  the **initial** states.

This corresponds to a **top-down** definition:





# Bottom-up vs top-down

## Comparison

These definitions of runs **coincide**: a bottom-up run exists iff a top-down run exists, and they are strictly the same:

$$\text{bottom-up TA } (\uparrow\text{TA}) = \text{top-down TA } (\downarrow\text{TA}) = \text{"TA"}$$

# Bottom-up vs top-down

## Comparison

These definitions of runs **coincide**: a bottom-up run exists iff a top-down run exists, and they are strictly the same:

$$\text{bottom-up TA } (\uparrow\text{TA}) = \text{top-down TA } (\downarrow\text{TA}) = \text{"TA"}$$



However, notions of **determinism** differ!

$$\text{det. } \downarrow\text{TA } (d\downarrow\text{TA}) \neq \text{det. } \uparrow\text{TA } (d\uparrow\text{TA})$$

For instance:  $\{f(a, b), f(b, a)\}$  can be recognized by a det. bottom-up TA, but any det. top-down TA accepting  $f(a, b)$  and  $f(b, a)$  would recognize  $f(a, a)$ .

# Determinization

of  $\uparrow$ TA

## Proposition

For every  $\uparrow$ TA  $A$ , there exists a  $d\uparrow$ TA  $A_d$  such that  $\mathcal{L}(A) = \mathcal{L}(A_d)$ .

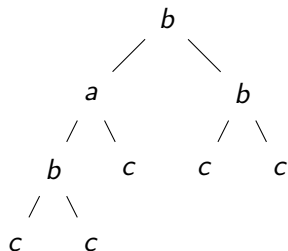
a subset construction, very similar to the determinization of NFAs:

- $Q_d = 2^Q$
- $\Delta_d = \{a(s_1, \dots, s_n) \rightarrow s \mid s = \{q \in Q \mid \exists q_1 \in s_1, \dots, \exists q_n \in s_n, f(q_1, \dots, q_n) \rightarrow q \in \Delta\}\}$
- $F_d = \{S \mid S \subseteq Q \text{ and } S \cap F \neq \emptyset\}$

This simulates all runs of  $A$ .

# Determinization

## Example run



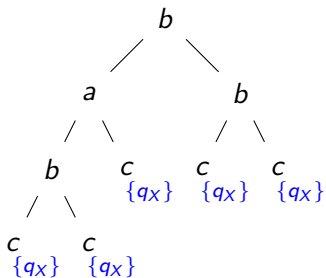
$$F = \{q_S\}$$

Rules of  $\uparrow$ TA:

$$\left\{ \begin{array}{l} a(q_S, q_X) \rightarrow q_S \\ a(q_X, q_S) \rightarrow q_S \\ b(q_S, q_X) \rightarrow q_S \\ b(q_X, q_S) \rightarrow q_S \\ a(q_B, q_X) \rightarrow q_S \\ a(q_X, q_B) \rightarrow q_S \\ b(q_X, q_X) \rightarrow q_B \\ a(q_X, q_X) \rightarrow q_X \\ b(q_X, q_X) \rightarrow q_X \\ c \rightarrow q_X \end{array} \right.$$

# Determinization

## Example run



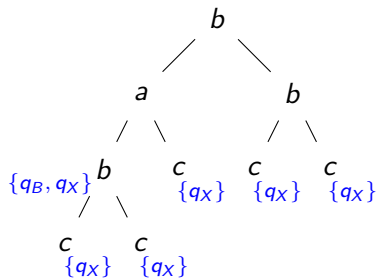
$$F = \{q_s\}$$

Rules of  $\uparrow$ TA:

$$\left\{ \begin{array}{l} a(q_s, q_x) \rightarrow q_s \\ a(q_x, q_s) \rightarrow q_s \\ b(q_s, q_x) \rightarrow q_s \\ b(q_x, q_s) \rightarrow q_s \\ a(q_b, q_x) \rightarrow q_s \\ a(q_x, q_b) \rightarrow q_s \\ b(q_x, q_x) \rightarrow q_b \\ a(q_x, q_x) \rightarrow q_x \\ b(q_x, q_x) \rightarrow q_x \\ c \rightarrow q_x \end{array} \right.$$

# Determinization

## Example run



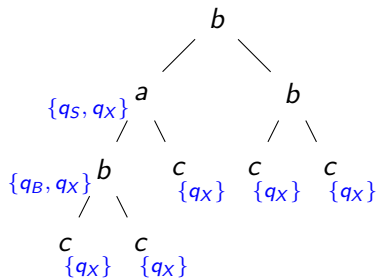
$$F = \{q_S\}$$

Rules of  $\uparrow$ TA:

$$\left\{ \begin{array}{l} a(q_S, q_X) \rightarrow q_S \\ a(q_X, q_S) \rightarrow q_S \\ b(q_S, q_X) \rightarrow q_S \\ b(q_X, q_S) \rightarrow q_S \\ a(q_B, q_X) \rightarrow q_S \\ a(q_X, q_B) \rightarrow q_S \\ b(q_X, q_X) \rightarrow q_B \\ a(q_X, q_X) \rightarrow q_X \\ b(q_X, q_X) \rightarrow q_X \\ c \rightarrow q_X \end{array} \right.$$

# Determinization

## Example run



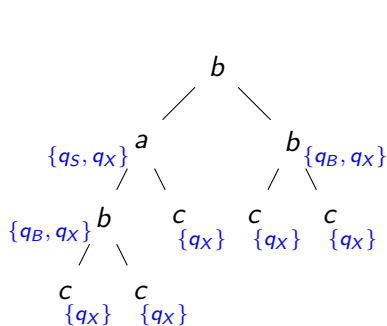
$$F = \{q_s\}$$

Rules of  $\uparrow$ TA:

$$\left\{ \begin{array}{l} a(q_s, q_x) \rightarrow q_s \\ a(q_x, q_s) \rightarrow q_s \\ b(q_s, q_x) \rightarrow q_s \\ b(q_x, q_s) \rightarrow q_s \\ a(q_B, q_x) \rightarrow q_s \\ a(q_x, q_B) \rightarrow q_s \\ b(q_x, q_x) \rightarrow q_B \\ a(q_x, q_x) \rightarrow q_x \\ b(q_x, q_x) \rightarrow q_x \\ c \rightarrow q_x \end{array} \right.$$

# Determinization

## Example run



$$F = \{q_S\}$$

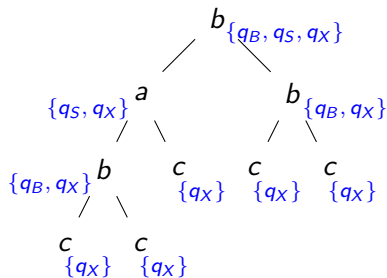
Rules of  $\uparrow$ TA:

$$\left\{ \begin{array}{l} a(q_S, q_X) \rightarrow q_S \\ a(q_X, q_S) \rightarrow q_S \\ b(q_S, q_X) \rightarrow q_S \\ b(q_X, q_S) \rightarrow q_S \\ a(q_B, q_X) \rightarrow q_S \\ a(q_X, q_B) \rightarrow q_S \\ b(q_X, q_X) \rightarrow q_B \\ a(q_X, q_X) \rightarrow q_X \\ b(q_X, q_X) \rightarrow q_X \\ c \rightarrow q_X \end{array} \right.$$



# Determinization

## Example run



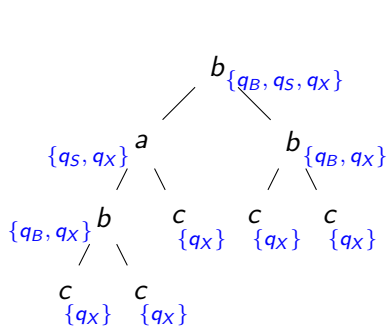
$$F = \{q_S\}$$

Rules of  $\uparrow$ TA:

$$\left\{ \begin{array}{l} a(q_S, q_X) \rightarrow q_S \\ a(q_X, q_S) \rightarrow q_S \\ b(q_S, q_X) \rightarrow q_S \\ b(q_X, q_S) \rightarrow q_S \\ a(q_B, q_X) \rightarrow q_S \\ a(q_X, q_B) \rightarrow q_S \\ b(q_X, q_X) \rightarrow q_B \\ a(q_X, q_X) \rightarrow q_X \\ b(q_X, q_X) \rightarrow q_X \\ c \rightarrow q_X \end{array} \right.$$

# Determinization

## Example run



$$F = \{q_S\}$$

Rules of  $\uparrow$ TA:

$$\left\{ \begin{array}{l} a(q_S, q_X) \rightarrow q_S \\ a(q_X, q_S) \rightarrow q_S \\ b(q_S, q_X) \rightarrow q_S \\ b(q_X, q_S) \rightarrow q_S \\ a(q_B, q_X) \rightarrow q_S \\ a(q_X, q_B) \rightarrow q_S \\ b(q_X, q_X) \rightarrow q_B \\ a(q_X, q_X) \rightarrow q_X \\ b(q_X, q_X) \rightarrow q_X \\ c \rightarrow q_X \end{array} \right.$$

$\{q_B, q_S, q_X\} \cap F \neq \emptyset$  so this tree is accepted by  $A_d$

# TA classes

$$d\uparrow\text{TA} = \uparrow\text{TA}$$

# TA classes

$$d\uparrow\text{TA} = \uparrow\text{TA} = \downarrow\text{TA}$$

# TA classes

$$d\downarrow\text{TA} \subsetneq d\uparrow\text{TA} = \uparrow\text{TA} = \downarrow\text{TA}$$

$$d\downarrow\text{TA} \subsetneq d\uparrow\text{TA} = \uparrow\text{TA} = \downarrow\text{TA}$$

## Definition

A ranked tree language  $L \subseteq \mathcal{T}_{\Sigma_r}$  is **recognizable** if there is a  $\uparrow\text{TA}$  recognizing  $L$ .

# Closure Properties

of recognizable tree languages

If  $L_1$  and  $L_2$  are recognizable tree languages, then:

$\overline{L_1}$  is recognizable

if  $A = (Q, F, \Delta, \Sigma_r)$  is a complete  $\uparrow$ TA, then  $A' = (Q, Q \setminus F, \Delta, \Sigma_r)$  recognizes  $\overline{\mathcal{L}(A)}$ .

Completing is easy, by adding a sink state.

$L_1 \cup L_2$  is recognizable

$L_1 \cap L_2$  is recognizable

# Closure Properties

of recognizable tree languages

If  $L_1$  and  $L_2$  are recognizable tree languages, then:

$\overline{L_1}$  is recognizable

$L_1 \cup L_2$  is recognizable

Let  $\begin{cases} A_1 = (Q_1, F_1, \Delta_1, \Sigma_r) \text{ be a complete } \uparrow\text{TA recognizing } L_1 \\ A_2 = (Q_2, F_2, \Delta_2, \Sigma_r) \text{ be a complete } \uparrow\text{TA recognizing } L_2 \end{cases}$

We build the product automaton

$A_1 \times A_2 = (Q_1 \times Q_2, F_1 \times Q_2 \cup Q_1 \times F_2, \Delta', \Sigma_r)$  with

$$\frac{a(q_1, \dots, q_n) \rightarrow q \in \Delta_1 \quad a(q'_1, \dots, q'_n) \rightarrow q' \in \Delta_2}{a((q_1, q'_1), \dots, (q_n, q'_n)) \rightarrow (q, q') \in \Delta'}$$

Then  $\mathcal{L}(A_1 \times A_2) = \mathcal{L}(A_1) \cup \mathcal{L}(A_2)$ .

$L_1 \cap L_2$  is recognizable



# Closure Properties

of recognizable tree languages

If  $L_1$  and  $L_2$  are recognizable tree languages, then:

$\overline{L_1}$  is recognizable

$L_1 \cup L_2$  is recognizable

$L_1 \cap L_2$  is recognizable

$L_1 \cap L_2 = \overline{\overline{L_1} \cup \overline{L_2}}$  (but more efficient using another product construction)

# Pumping Lemma

**Context** = tree over  $\Sigma_r \uplus \{\square\}$  where  $\begin{cases} \text{ar}(\square) = 0 \text{ and} \\ \square \text{ appears exactly once.} \end{cases}$

Intuitively, *context* = *tree with a hole*.

(context  $C$ ) + (tree  $t$ ) = a new tree  $C[t]$ , where  $t$  replaces  $\square$  in  $C$ .

# Pumping Lemma

**Context** = tree over  $\Sigma_r \uplus \{\square\}$  where  $\begin{cases} \text{ar}(\square) = 0 \text{ and} \\ \square \text{ appears exactly once.} \end{cases}$

Intuitively, *context* = *tree with a hole*.

(context  $C$ ) + (tree  $t$ ) = a new tree  $C[t]$ , where  $t$  replaces  $\square$  in  $C$ .

## Pumping lemma

Let  $L$  be a *recognizable* tree language. Then there exists  $k > 0$  such that for every tree  $t \in L$  of depth  $> k$ , there exist contexts  $C_1, C_2$  (with  $C_2 \neq \square$ ) and a tree  $t'$  such that  $\begin{cases} t = C_1[C_2[t']] \\ \forall n \geq 0. C_1[C_2^n[t']] \in L \end{cases}$

idea: for  $A$  s.t.  $\mathcal{L}(A) = L$ , take  $k = |Q|$  and consider a branch of length  $k$ .

# Tree homomorphisms

Tree homomorphism =  $\left\{ \begin{array}{l} \text{transformation } h : \mathcal{T}_\Sigma \rightarrow \mathcal{T}_{\Sigma'} \\ \text{based on a mapping } m : \Sigma \rightarrow \mathcal{T}_{\Sigma' \cup \mathcal{X}_{\text{ar}(m)}} \end{array} \right.$

## Theorem

$\left. \begin{array}{l} L: \text{recognizable tree language in } \mathcal{T}_\Sigma \\ h: \text{linear tree homom. } h : \mathcal{T}_\Sigma \rightarrow \mathcal{T}_{\Sigma'} \end{array} \right\} \Rightarrow h(L) \text{ is recognizable (over } \mathcal{T}_{\Sigma'})$

*Linear* means that each variable appears at most once in  $m(a)$ .

Simple counter-example with  $m(a) = a(x_1, x_1)$ .

# Tree homomorphisms

Tree homomorphism =  $\left\{ \begin{array}{l} \text{transformation } h : \mathcal{T}_\Sigma \rightarrow \mathcal{T}_{\Sigma'} \\ \text{based on a mapping } m : \Sigma \rightarrow \mathcal{T}_{\Sigma' \cup \mathcal{X}_{\text{ar}}(m)} \end{array} \right.$

## Theorem

$\left. \begin{array}{l} L: \text{recognizable tree language in } \mathcal{T}_\Sigma \\ h: \text{linear tree homom. } h : \mathcal{T}_\Sigma \rightarrow \mathcal{T}_{\Sigma'} \end{array} \right\} \Rightarrow h(L) \text{ is recognizable (over } \mathcal{T}_{\Sigma'})$

*Linear* means that each variable appears at most once in  $m(a)$ .

Simple counter-example with  $m(a) = a(x_1, x_1)$ .

## Theorem

$\left. \begin{array}{l} L: \text{recognizable tree language in } \mathcal{T}_{\Sigma'} \\ h: \text{any tree homom. } h : \mathcal{T}_\Sigma \rightarrow \mathcal{T}_{\Sigma'} \end{array} \right\} \Rightarrow h^{-1}(L) \text{ is recognizable (over } \mathcal{T}_\Sigma)$

# Minimization

## of Tree Automata

Let  $\equiv$  be an equivalence relation. It is:

- a **congruence** if  $\forall a \in \mathcal{T}_{\Sigma, r}$ ,

if  $t_i \equiv t'_i$  for all  $1 \leq i \leq n$  then  $a(t_1, \dots, t_n) \equiv a(t'_1, \dots, t'_n)$

- of **finite index** if there are only finitely many  $\equiv$ -classes

# Minimization

## of Tree Automata

Let  $\equiv$  be an equivalence relation. It is:

- a **congruence** if  $\forall a \in \mathcal{T}_{\Sigma_r}$ ,

if  $t_i \equiv t'_i$  for all  $1 \leq i \leq n$  then  $a(t_1, \dots, t_n) \equiv a(t'_1, \dots, t'_n)$

- of **finite index** if there are only finitely many  $\equiv$ -classes



Given a tree language  $L$ , we define the congruence  $\equiv_L$ :

$t \equiv_L t'$  if for all contexts  $C$  over  $\Sigma_r$ ,  $C[t] \in L \Leftrightarrow C[t'] \in L$

# Minimization

## of Tree Automata

Def:  $t \equiv_L t'$  if  $\forall C. C[t] \in L \Leftrightarrow C[t'] \in L$

## Myhill-Nerode Theorem

$L$  is a **recognizable** tree language iff  $\equiv_L$  is of **finite index**.

- ( $\Rightarrow$ ) Let  $A$  be a complete  $d\uparrow$ TA recognizing  $L$ . Let  $\equiv_A$  defined by:  
 $t \equiv_A t'$  iff  $\Delta(t) = \Delta(t')$  (state of  $A$ ). It is of finite index ( $\leq |Q_A|$ ),  
and  $L = \bigcup_{t \mid \Delta(t) \in F_A} \text{class}_{\equiv_A}(t)$ .



# Minimization

## of Tree Automata

Def:  $t \equiv_L t'$  if  $\forall C. C[t] \in L \Leftrightarrow C[t'] \in L$

## Myhill-Nerode Theorem

$L$  is a **recognizable** tree language iff  $\equiv_L$  is of **finite index**.

- ( $\Rightarrow$ ) Let  $A$  be a complete  $d\uparrow$ TA recognizing  $L$ . Let  $\equiv_A$  defined by:  $t \equiv_A t'$  iff  $\Delta(t) = \Delta(t')$  (state of  $A$ ). It is of finite index ( $\leq |Q_A|$ ), and  $L = \bigcup_{t \mid \Delta(t) \in F_A} \text{class}_{\equiv_A}(t)$ .  
As  $\equiv_A$  is of finite index, it suffices to prove that  $t \equiv_A t' \Rightarrow t \equiv_L t'$ .  
Assume  $t \equiv_A t'$ . By an easy induction,  $C[t] \equiv_A C[t']$  for all  $C$ . As  $L$  is a union of classes of  $\equiv_A$ ,  $C[t] \in L \Leftrightarrow C[t'] \in L$ , and thus  $t \equiv_L t'$ .

# Minimization

## of Tree Automata

Def:  $t \equiv_L t'$  if  $\forall C. C[t] \in L \Leftrightarrow C[t'] \in L$

## Myhill-Nerode Theorem

$L$  is a **recognizable** tree language iff  $\equiv_L$  is of **finite index**.

- ( $\Rightarrow$ ) Let  $A$  be a complete  $d\uparrow$ TA recognizing  $L$ . Let  $\equiv_A$  defined by:  $t \equiv_A t'$  iff  $\Delta(t) = \Delta(t')$  (state of  $A$ ). It is of finite index ( $\leq |Q_A|$ ), and  $L = \bigcup_{t \mid \Delta(t) \in F_A} \text{class}_{\equiv_A}(t)$ .  
As  $\equiv_A$  is of finite index, it suffices to prove that  $t \equiv_A t' \Rightarrow t \equiv_L t'$ .  
Assume  $t \equiv_A t'$ . By an easy induction,  $C[t] \equiv_A C[t']$  for all  $C$ . As  $L$  is a union of classes of  $\equiv_A$ ,  $C[t] \in L \Leftrightarrow C[t'] \in L$ , and thus  $t \equiv_L t'$ .
- ( $\Leftarrow$ )  $A_{min} = (Q_{min}, F_{min}, \Delta_{min}, \Sigma_r)$  with:
  - ▶  $Q_{min} = \text{class}_{\equiv_L}$
  - ▶  $F_{min} = \{\text{class}_{\equiv_L}(t) \mid t \in L\}$
  - ▶  $\Delta_{min}$  contains rules:  
 $f(\text{class}_{\equiv_L}(t_1), \dots, \text{class}_{\equiv_L}(t_n)) \rightarrow \text{class}_{\equiv_L}(f(t_1, \dots, t_n))$

# Minimization

## of Tree Automata

### Consequence of Myhill-Nerode

The minimum  $d\uparrow$ TA recognizing  $L$  is unique, up to a renaming of states.

if  $A$  recognizes  $L$ ,  $\equiv_A$  is a refinement of  $\equiv_L$ , so  $|Q_A| \geq |Q_{A_{min}}|$ .

### Minimization algorithm (sketch)

- input: complete and reduced  $d\uparrow$ TA  $A$
- start from  $\equiv_A$  and build  $\equiv_L$ ,
- by merging  $q$  and  $q'$  if  $\forall a \in \Sigma_r, \forall i, \forall q_1, \dots, q_{i-1}, q_{i+1}, \dots, q_n \in Q_A$ ,  
 $\Delta(a(q_1, \dots, q_{i-1}, q, q_{i+1}, \dots, q_n)) \equiv \Delta(a(q_1, \dots, q_{i-1}, q', q_{i+1}, \dots, q_n))$
- until fixed point

## Complexity of some decision problems

Name	Input	Output	Complexity
Membership	$\uparrow$ TA $A$ , tree $t$	$t \in \mathcal{L}(A)$ ?	PTIME
Emptiness	$\uparrow$ TA $A$	$\mathcal{L}(A) = \emptyset$ ?	PTIME
Intersection non-emptiness	set $S$ of $\uparrow$ TAs	$\bigcap_{A \in S} \mathcal{L}(A) = \emptyset$ ?	EXPTIME-compl.
	set $S$ of $d\uparrow$ TAs		EXPTIME-compl.
Universality	$\uparrow$ TA $A$	$\mathcal{L}(A) = T_{\Sigma_r}$ ?	EXPTIME-compl.
	$d\uparrow$ TA $A$		PTIME
Equivalence	$\uparrow$ TAs $A_1, A_2$	$\mathcal{L}(A_1) = \mathcal{L}(A_2)$ ?	EXPTIME-compl.
	$d\uparrow$ TAs $A_1, A_2$		PTIME

## 1 Ranked Trees

- Trees on Ranked Alphabet
- Tree Automata
- **Tree Grammars**
- Logic

## 2 Unranked Trees

- Unranked Trees
- Automata
- Logic

# Tree Grammars

Let  $\mathcal{X}$  be a set of variables.

## Tree grammar

A **tree grammar** is a tuple  $G = (S, N, F, R)$  where:

- $N$  is a set of **non-terminal** symbols,
- $S \in N$  an **axiom**,
- $F$  a set of **terminal** symbols ( $F \cap N = \emptyset$ ),
- $R$  a set of **production rules**  $\alpha \rightarrow \beta$  with  $\alpha, \beta \in \mathcal{T}_{N \cup F \cup \mathcal{X}}$  and  $\alpha$  contains at least one non-terminal.

Each element of  $N \cup F$  has a fixed arity,  $\text{ar}(S) = 0$ , and  $\text{ar}(x) = 0$ , for  $x \in \mathcal{X}$ .

$$a(b(x, d(N_1(c)))) \rightarrow a(d(x), b(N_2, c))$$

# Regular Tree Grammars

## Definition

### Regular tree grammars

A tree grammar  $G$  is **regular** if  $\text{ar}(N_i) = 0$  for all  $N_i \in N$ , and production rules are of the form  $N_i \rightarrow \beta$  with  $N_i \in N$  and  $\beta \in \mathcal{T}_{NUF}$ .

$\mathcal{L}(G)$  is the set of trees obtained by applying a series of rules, starting from  $S$ .

$L \subseteq \mathcal{T}_{\Sigma_r}$  is said **regular** if  $L = \mathcal{L}(G)$  for some regular tree grammar  $G$ .

### Example

$$R = \begin{cases} S \rightarrow a(S, X) & S \rightarrow b(S, X) & S \rightarrow a(B, X) & B \rightarrow b(X, X) & X \rightarrow b(X, X) \\ S \rightarrow a(X, S) & S \rightarrow b(X, S) & S \rightarrow a(X, B) & X \rightarrow a(X, X) & X \rightarrow c \end{cases}$$

$a(b(c, c), c) \in \mathcal{L}(G)$  because:

$$S \rightarrow_G a(B, X) \rightarrow_G a(b(X, X), X) \rightarrow_G a(b(c, X), X) \rightarrow_G a(b(c, c), X) \rightarrow_G a(b(c, c), c)$$

Regular tree grammars can be **normalized**, so that rules have the form  $N_0 \rightarrow a(N_1, \dots, N_n)$ .

# Regular Tree Grammars

## Expressiveness

### Proposition

regular tree languages = recognizable tree languages

Normalized regular tree grammars are  $\downarrow$ TAs...



# Regular Tree Grammars

## Expressiveness

### Proposition

regular tree languages = recognizable tree languages

Normalized regular tree grammars are  $\downarrow$ TAs...

### Regular Expressions

There also exists a notion of **regular expressions** for trees, with the same expressiveness as regular tree grammar (Kleene's theorem).

skipped in this talk...

## 1 Ranked Trees

- Trees on Ranked Alphabet
- Tree Automata
- Tree Grammars
- **Logic**

## 2 Unranked Trees

- Unranked Trees
- Automata
- Logic

# Monadic Second-Order (MSO) Logic

## Syntax

- $\Sigma_r$  a ranked signature
- $\mathcal{X}$  contains first-order  $(x, y\dots)$  and second-order  $(X, Y\dots)$  variables.
- $\Omega = \{\text{lab}_a \mid a \in \Sigma_r\} \cup \{\text{ch}_i \mid \exists a \in \Sigma_r. \text{ar}(a) \geq i\}$
- predicates  $\text{lab}_a$  are unary, while  $\text{ch}_i$  are binary

## MSO[ $\Omega$ ]: Syntax

$\phi ::= \text{lab}_a(x) \mid \text{ch}_i(x, y) \mid \phi \wedge \phi \mid \neg\phi \mid \exists x. \phi \mid \exists X. \phi \mid x \in X$

where  $a \in \Sigma$ ,  $\exists a \in \Sigma_r. \text{ar}(a) \geq i$ , and  $x, x_1, \dots, x_k, X \in \mathcal{X}$ .

# Monadic Second-Order (MSO) Logic

## Syntax

- $\Sigma_r$  a ranked signature
- $\mathcal{X}$  contains first-order  $(x, y\dots)$  and second-order  $(X, Y\dots)$  variables.
- $\Omega = \{\text{lab}_a \mid a \in \Sigma_r\} \cup \{\text{ch}_i \mid \exists a \in \Sigma_r. \text{ar}(a) \geq i\}$
- predicates  $\text{lab}_a$  are unary, while  $\text{ch}_i$  are binary

## MSO[ $\Omega$ ]: Syntax

$\phi ::= \text{lab}_a(x) \mid \text{ch}_i(x, y) \mid \phi \wedge \phi \mid \neg\phi \mid \exists x. \phi \mid \exists X. \phi \mid x \in X$

where  $a \in \Sigma$ ,  $\exists a \in \Sigma_r. \text{ar}(a) \geq i$ , and  $x, x_1, \dots, x_k, X \in \mathcal{X}$ .

MSO[ $\Omega$ ] is also known as **WSkS**, the *Weak Second-order logic with k Successors*.  
“Weak” means “interpreted over finite structures” (here, terms).

# Monadic Second-Order (MSO) Logic

## Example

For instance: all trees having an even number of  $a$ -labeled nodes.

*hint*: define  $X$  as the set of nodes having an even number of  $a$ -descendants

# Monadic Second-Order (MSO) Logic

## Example

For instance: all trees having an even number of  $a$ -labeled nodes.

*hint:* define  $X$  as the set of nodes having an even number of  $a$ -descendants

$$\text{even}_a = \exists X \forall x \text{lab}_a(x) \Rightarrow \begin{cases} \exists y_1 \text{ch}_1(x, y_1) \wedge \\ \exists y_2 \text{ch}_2(x, y_2) \wedge \\ x \in X \Leftrightarrow (y_1 \in X \oplus y_2 \in X) \end{cases}$$

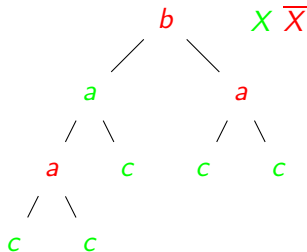
$$\bigwedge_{\alpha \neq a} \text{lab}_\alpha(x) \Rightarrow \begin{cases} \bigwedge_{1 \leq i \leq \text{ar}(\alpha)} \exists y_i \text{ch}_i(x, y_i) \wedge \\ x \notin X \Leftrightarrow (y_1 \notin X \oplus \dots \oplus y_{\text{ar}(\alpha)} \notin X) \end{cases}$$

$$\wedge \exists x \text{root}(x) \wedge x \in X$$

where:

$$\phi \oplus \phi' = (\phi \wedge \neg \phi') \vee (\neg \phi \wedge \phi')$$

$$\text{root}(x) = \neg \exists y. \text{ch}_1(y, x)$$



# Monadic Second-Order (MSO) Logic

## Semantics

*convention:*  $\phi(\bar{x}, \bar{X})$  means that  $\phi$  has free FO variables  $\bar{x}$  and free SO variables  $\bar{X}$ .

A formula  $\phi(\bar{x}, \bar{X}) \in MSO[\Omega]$  is interpreted over a tree  $t$  under an assignment  $\mu : \bar{x} \cup \bar{X} \rightarrow \text{nodes}(t)$ .

Satisfiability  $t, \mu \models \phi$  is defined inductively by:

$$t, \mu \models \text{lab}_a(x) \quad \text{iff} \quad \text{lab}_a^t(\mu(x))$$

$$t, \mu \models \text{ch}_i(x, y) \quad \text{iff} \quad \text{ch}_i(\mu(x), \mu(y))$$

$$t, \mu \models \phi \wedge \phi' \quad \text{iff} \quad t, \mu \models \phi \text{ and } t, \mu \models \phi'$$

$$t, \mu \models \neg\phi \quad \text{iff} \quad t, \mu \not\models \phi$$

$$t, \mu \models \exists x\phi \quad \text{iff} \quad \text{there exists } \pi \in \text{nodes}(t) \text{ s.t. } t, \mu[x \leftarrow \pi] \models \phi$$

$$t, \mu \models \exists X\phi \quad \text{iff} \quad \text{there exists } S \subseteq \text{nodes}(t) \text{ s.t. } t, \mu[X \leftarrow S] \models \phi$$

$$t, \mu \models x \in X \quad \text{iff} \quad \mu(x) \in \mu(X)$$

# MSO and Tree Languages

Expressiveness of MSO vs  $\uparrow$ TA?

For closed formulas  $\phi$ , define  $\mathcal{L}_\phi = \{t \mid t \models \phi\} \subseteq \mathcal{T}_{\Sigma, r}$ .

What about  $\phi(\bar{x}, \bar{X})$ ?

$\phi(\bar{x}, \bar{X}) \xrightarrow{\text{semantics}} \{(t, \mu) \mid t, \mu \models \phi\} \xrightarrow{\text{tree language}} \mathcal{L}_\phi = \{t * \mu \mid t, \mu \models \phi\}$



# MSO and Tree Languages

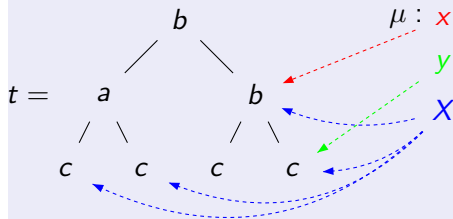
Expressiveness of MSO vs  $\uparrow$ TA?

For closed formulas  $\phi$ , define  $\mathcal{L}_\phi = \{t \mid t \models \phi\} \subseteq \mathcal{T}_{\Sigma, r}$ .

What about  $\phi(\overline{x}, \overline{X})$ ?

$\phi(\overline{x}, \overline{X}) \xrightarrow{\text{semantics}} \{(t, \mu) \mid t, \mu \models \phi\} \xrightarrow{\text{tree language}} \mathcal{L}_\phi = \{t * \mu \mid t, \mu \models \phi\}$

$t * \mu \in \mathcal{T}_{\Sigma, r \times \mathbb{B}^n}$



where  $n = |\overline{x}| + |\overline{X}|$  and  $\mathbb{B} = \{0, 1\}$ .

# MSO and Tree Languages

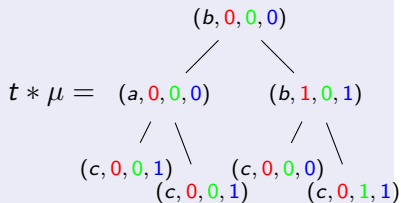
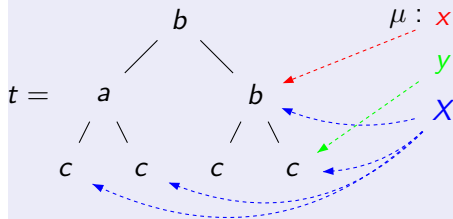
Expressiveness of MSO vs  $\uparrow$ TA?

For closed formulas  $\phi$ , define  $\mathcal{L}_\phi = \{t \mid t \models \phi\} \subseteq \mathcal{T}_{\Sigma, r}$ .

What about  $\phi(\overline{x}, \overline{X})$ ?

$\phi(\overline{x}, \overline{X}) \xrightarrow{\text{semantics}} \{(t, \mu) \mid t, \mu \models \phi\} \xrightarrow{\text{tree language}} \mathcal{L}_\phi = \{t * \mu \mid t, \mu \models \phi\}$

$t * \mu \in \mathcal{T}_{\Sigma, r \times \mathbb{B}^n}$



where  $n = |\overline{x}| + |\overline{X}|$  and  $\mathbb{B} = \{0, 1\}$ .

# Recognizable languages = MSO-definable languages

## Theorem [TW68, Don70]

$L \subseteq \mathcal{T}_{\Sigma_r}$  is recognizable iff there exists  $\phi \in \text{MSO}[\Omega]$  such that  $\mathcal{L}_\phi = L$ .

consequence:  $\text{MSO}[\Omega]$  is decidable (convert to  $\uparrow$ TA, test emptiness)

( $\Rightarrow$ ) *idea: encode a run of A in MSO.*

$A = (Q, F, \Delta, \Sigma_r)$  a complete  $d\uparrow$ TA recognizing  $L$ . Let  $\{q_1, \dots, q_n\} = Q$ .

$\phi_A = \exists X_{q_1} \dots \exists X_{q_n} \text{partition}(X_{q_1}, \dots, X_{q_n}) \wedge$

$\bigwedge_{a \rightarrow q \in \Delta} \text{leaf}(x) \wedge \text{lab}_a(x) \Rightarrow x \in X_q$

$\bigwedge_{a(q_{i_1}, \dots, q_{i_k}) \rightarrow q_i \in \Delta} \forall x \forall y_1 \dots \forall y_k \left\{ \begin{array}{l} \text{lab}_a(x) \wedge \\ \text{ch}_1(x, y_1) \wedge \dots \wedge \text{ch}(x, y_k) \wedge \\ y_1 \in X_{q_{i_1}} \wedge \dots \wedge y_k \in X_{q_{i_k}} \end{array} \right\} \Rightarrow x \in X_{q_i}$

$\bigwedge \exists x. \text{root}(x) \wedge \bigvee_{q \in F} x \in X_q$

with:  $\left\{ \begin{array}{l} \text{leaf}(x) = \nexists y. \text{ch}_1(x, y) \\ \text{partition}(X_{q_1}, \dots, X_{q_n}) = \dots \end{array} \right.$

( $\Leftarrow$ ) *idea: equivalence betw. logical connectives and automata operators.*

$$\bullet \mathcal{L}_{\text{lab}_a} = \left\{ \dots, \begin{array}{c} (a,1) \\ / \quad \backslash \\ (a,0) \quad (b,0) \end{array}, \begin{array}{c} (a,0) \\ / \quad \backslash \\ (a,1) \quad (b,0) \end{array}, \dots \right\} \text{recognizable}$$

( $\Leftarrow$ ) *idea: equivalence betw. logical connectives and automata operators.*

$$\bullet \mathcal{L}_{lab_a} = \left\{ \dots, \begin{array}{c} (a,1) \\ / \quad \backslash \\ (a,0) \quad (b,0) \end{array}, \begin{array}{c} (a,0) \\ / \quad \backslash \\ (a,1) \quad (b,0) \end{array}, \dots \right\} \text{recognizable}$$

$$\bullet \mathcal{L}_{ch_1} = \left\{ \dots, \begin{array}{c} (a,1,0) \\ / \quad \backslash \\ (a,0,1) \quad (b,0,0) \end{array}, \dots \right\} \text{recognizable } (ch_i \text{ also})$$

( $\Leftarrow$ ) *idea: equivalence betw. logical connectives and automata operators.*

$$\bullet \mathcal{L}_{lab_a} = \left\{ \dots, \begin{array}{c} (a,1) \\ / \quad \backslash \\ (a,0) \quad (b,0) \end{array}, \begin{array}{c} (a,0) \\ / \quad \backslash \\ (a,1) \quad (b,0) \end{array}, \dots \right\} \text{recognizable}$$

$$\bullet \mathcal{L}_{ch_1} = \left\{ \dots, \begin{array}{c} (a,1,0) \\ / \quad \backslash \\ (a,0,1) \quad (b,0,0) \end{array}, \dots \right\} \text{recognizable } (ch_i \text{ also})$$

$$\bullet \mathcal{L}_{\phi \wedge \phi'} = \mathcal{L}_{\phi} \cap \mathcal{L}_{\phi'} \text{ recognizable (product construction)}$$

( $\Leftarrow$ ) *idea: equivalence betw. logical connectives and automata operators.*

$$\bullet \mathcal{L}_{lab_a} = \left\{ \dots, \begin{array}{c} (a,1) \\ / \quad \backslash \\ (a,0) \quad (b,0) \end{array}, \begin{array}{c} (a,0) \\ / \quad \backslash \\ (a,1) \quad (b,0) \end{array}, \dots \right\} \text{recognizable}$$

$$\bullet \mathcal{L}_{ch_1} = \left\{ \dots, \begin{array}{c} (a,1,0) \\ / \quad \backslash \\ (a,0,1) \quad (b,0,0) \end{array}, \dots \right\} \text{recognizable } (ch_i \text{ also})$$

$$\bullet \mathcal{L}_{\phi \wedge \phi'} = \mathcal{L}_{\phi} \cap \mathcal{L}_{\phi'} \text{ recognizable (product construction)}$$

$$\bullet \mathcal{L}_{\neg \phi} = \mathcal{L}_{valid} \setminus \mathcal{L}_{\phi} \text{ where } \mathcal{L}_{valid} = \{t * \mu \mid t \in \mathcal{T}_{\Sigma_r}, \mu \text{ an assignment of free vars of } \phi\}$$

$\mathcal{L}_{valid}$  is recognizable so  $\mathcal{L}_{\neg \phi}$  is recognizable

( $\Leftarrow$ ) *idea: equivalence betw. logical connectives and automata operators.*

- $$\mathcal{L}_{lab_a} = \left\{ \dots, \begin{array}{c} (a,1) \\ / \quad \backslash \\ (a,0) \quad (b,0) \end{array}, \begin{array}{c} (a,0) \\ / \quad \backslash \\ (a,1) \quad (b,0) \end{array}, \dots \right\} \text{ recognizable}$$

- $$\mathcal{L}_{ch_1} = \left\{ \dots, \begin{array}{c} (a,1,0) \\ / \quad \backslash \\ (a,0,1) \quad (b,0,0) \end{array}, \dots \right\} \text{ recognizable } (ch_i \text{ also})$$

- $$\mathcal{L}_{\phi \wedge \phi'} = \mathcal{L}_{\phi} \cap \mathcal{L}_{\phi'} \text{ recognizable (product construction)}$$

- $$\mathcal{L}_{\neg \phi} = \mathcal{L}_{valid} \setminus \mathcal{L}_{\phi} \text{ where } \mathcal{L}_{valid} = \{t * \mu \mid t \in \mathcal{T}_{\Sigma_r}, \mu \text{ an assignment of free vars of } \phi\}$$

$$\mathcal{L}_{valid} \text{ is recognizable so } \mathcal{L}_{\neg \phi} \text{ is recognizable}$$

- $$\mathcal{L}_{\exists x \phi} \text{ is obtained from } \mathcal{L}_{\phi} \text{ by removing the } x\text{-component:}$$

$$\begin{array}{c} (a,1,0) \\ / \quad \backslash \\ (a,0,1) \quad (b,0,0) \end{array} \in \mathcal{L}_{ch_1(x,y)} \quad \rightarrow \quad \begin{array}{c} (a,0) \\ / \quad \backslash \\ (a,1) \quad (b,0) \end{array} \in \mathcal{L}_{\exists x \ ch_1(x,y)}$$

recognizable (remove the component in rules)



( $\Leftarrow$ ) *idea: equivalence betw. logical connectives and automata operators.*

$$\bullet \mathcal{L}_{lab_a} = \left\{ \dots, \begin{array}{c} (a,1) \\ / \quad \backslash \\ (a,0) \quad (b,0) \end{array}, \begin{array}{c} (a,0) \\ / \quad \backslash \\ (a,1) \quad (b,0) \end{array}, \dots \right\} \text{ recognizable}$$

$$\bullet \mathcal{L}_{ch_1} = \left\{ \dots, \begin{array}{c} (a,1,0) \\ / \quad \backslash \\ (a,0,1) \quad (b,0,0) \end{array}, \dots \right\} \text{ recognizable } (ch_i \text{ also})$$

$$\bullet \mathcal{L}_{\phi \wedge \phi'} = \mathcal{L}_{\phi} \cap \mathcal{L}_{\phi'} \text{ recognizable (product construction)}$$

$$\bullet \mathcal{L}_{\neg \phi} = \mathcal{L}_{valid} \setminus \mathcal{L}_{\phi} \text{ where } \mathcal{L}_{valid} = \{t * \mu \mid t \in \mathcal{T}_{\Sigma_r}, \mu \text{ an assignment of free vars of } \phi\}$$

$\mathcal{L}_{valid}$  is recognizable so  $\mathcal{L}_{\neg \phi}$  is recognizable

$$\bullet \mathcal{L}_{\exists x \phi} \text{ is obtained from } \mathcal{L}_{\phi} \text{ by removing the } x\text{-component:}$$

$$\begin{array}{c} (a,1,0) \\ / \quad \backslash \\ (a,0,1) \quad (b,0,0) \end{array} \in \mathcal{L}_{ch_1(x,y)} \quad \rightarrow \quad \begin{array}{c} (a,0) \\ / \quad \backslash \\ (a,1) \quad (b,0) \end{array} \in \mathcal{L}_{\exists x \ ch_1(x,y)}$$

recognizable (remove the component in rules)

$$\bullet \mathcal{L}_{\exists X \phi}: \text{ same idea}$$

( $\Leftarrow$ ) *idea: equivalence betw. logical connectives and automata operators.*

- $\mathcal{L}_{lab_a} = \left\{ \dots, \begin{array}{c} (a,1) \\ / \quad \backslash \\ (a,0) \quad (b,0) \end{array}, \begin{array}{c} (a,0) \\ / \quad \backslash \\ (a,1) \quad (b,0) \end{array}, \dots \right\}$  recognizable

- $\mathcal{L}_{ch_1} = \left\{ \dots, \begin{array}{c} (a,1,0) \\ / \quad \backslash \\ (a,0,1) \quad (b,0,0) \end{array}, \dots \right\}$  recognizable ( $ch_i$  also)

- $\mathcal{L}_{\phi \wedge \phi'} = \mathcal{L}_{\phi} \cap \mathcal{L}_{\phi'}$  recognizable (product construction)

- $\mathcal{L}_{\neg\phi} = \mathcal{L}_{valid} \setminus \mathcal{L}_{\phi}$  where  $\mathcal{L}_{valid} = \{t * \mu \mid t \in \mathcal{T}_{\Sigma_r}, \mu \text{ an assignment of free vars of } \phi\}$   
 $\mathcal{L}_{valid}$  is recognizable so  $\mathcal{L}_{\neg\phi}$  is recognizable

- $\mathcal{L}_{\exists x\phi}$  is obtained from  $\mathcal{L}_{\phi}$  by removing the  $x$ -component:

$$\begin{array}{c} (a,1,0) \\ / \quad \backslash \\ (a,0,1) \quad (b,0,0) \end{array} \in \mathcal{L}_{ch_1(x,y)} \quad \rightarrow \quad \begin{array}{c} (a,0) \\ / \quad \backslash \\ (a,1) \quad (b,0) \end{array} \in \mathcal{L}_{\exists x \ ch_1(x,y)}$$

recognizable (remove the component in rules)

- $\mathcal{L}_{\exists X\phi}$ : same idea

- $\mathcal{L}_{x \in X}$ : easily recognizable (a 1 on  $x$ -component implies a 1 on  $X$ -component)

$$\begin{array}{c} (a,1,1) \\ / \quad \backslash \\ (a,0,1) \quad (b,0,0) \end{array} \in \mathcal{L}_{x \in X} \quad \text{but} \quad \begin{array}{c} (a,0,1) \\ / \quad \backslash \\ (a,1,0) \quad (b,0,1) \end{array} \notin \mathcal{L}_{x \in X}$$

## Bonus: Recognizable tree relations

see [CDG<sup>+</sup>07] Chapter 3, and also [BLN07, Gau09]

## 1 Ranked Trees

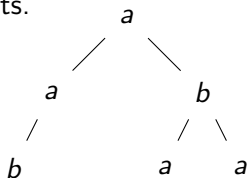
- Trees on Ranked Alphabet
- Tree Automata
- Tree Grammars
- Logic

## 2 Unranked Trees

- Unranked Trees
- Automata
- Logic

# Unranked Trees

An **unranked tree** over  $\Sigma$  is a tree, labeled by elements of  $\Sigma$ , without arity constraints.



We write  $\mathcal{T}_\Sigma$  for the set of unranked trees over  $\Sigma$ .

## 1 Ranked Trees

- Trees on Ranked Alphabet
- Tree Automata
- Tree Grammars
- Logic

## 2 Unranked Trees

- Unranked Trees
- **Automata**
- Logic

# Which notion of automata can we use?

For ranked trees, rules were  $a(q_1, \dots, q_n) \rightarrow q$  but here  $n$  is not bounded.  
Several approaches:

- **horizontal languages**: use a string language on the states of children
  - ▶ hedge automata
  - ▶ DTDs
- **binary encodings**: encode unranked trees into binary trees
  - ▶ stepwise tree automata...
- **linearization**: serialize trees and use a pushdown system
  - ▶ nested word automata
  - ▶ visibly pushdown automata

# Horizontal Languages: Hedge Automata

*idea*: use a regular language on the states of children

“hedge” = finite series of trees

## Hedge automata [BKWM01]

A **hedge automaton** over  $\Sigma$  is a tuple  $A = (Q, F, \Delta, \Sigma)$  where:

- $Q$  is a finite set of **states**
- $F \subseteq Q$  is the set of **final** states
- $\Delta$  is a set of **rules**  $a(L) \rightarrow q$  where  $a \in \Sigma$ ,  $q \in Q$  and  $L$  is a regular string language over  $Q$ .

A run of  $A$  on  $t$  is a function  $\rho : \text{nodes}(t) \rightarrow Q$  such that for all nodes  $\pi$  of  $t$  with children  $\pi_1, \dots, \pi_n$  and label  $a$ , there is a rule  $a(L) \rightarrow \rho(\pi) \in \Delta$  with  $\rho(\pi_1) \dots \rho(\pi_n) \in L$ .

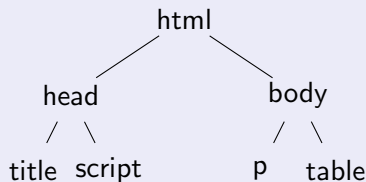


# Horizontal Languages: DTDs

- Document Type Definitions [BPSM<sup>+</sup>08]
- W3C standard for specifying valid XML documents
- hedge automata, where horizontal languages are specified by regexp.

## Example: HTML DTD

*html* → *head.body*  
*head* → *title.meta?.style?.script?...*  
*body* → *(p|div|table|h1|...)\**  
...



# Binary Encodings

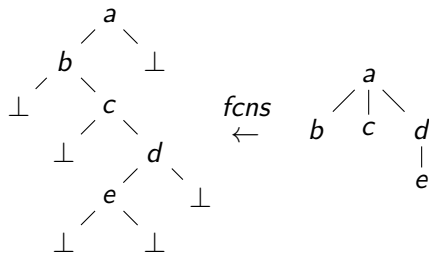
- define a bijection between unranked trees and binary trees
- 2 common encodings: *first-child next-sibling* and *Curryfication*

# Binary Encodings

- define a bijection between unranked trees and binary trees
- 2 common encodings: **first-child next-sibling** and **Curryfication**

first-child next-sibling [Rab69, Koc03]

$fcns : \mathcal{T}_\Sigma \rightarrow \mathcal{T}_{\Sigma_f}$  where  $\Sigma_f = (\Sigma \uplus \{\perp\}, ar)$  with  $ar(a)=2$  for  $a \in \Sigma$  and  $ar(\perp)=0$



# Binary Encodings

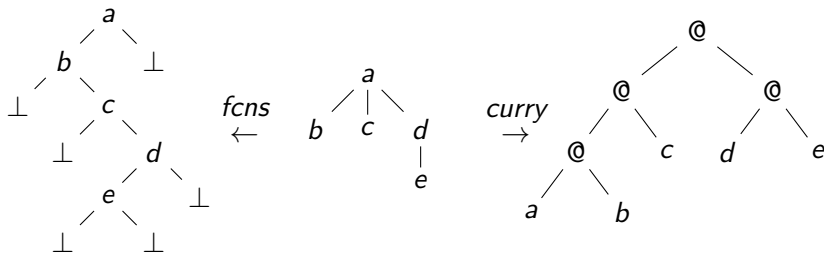
- define a bijection between unranked trees and binary trees
- 2 common encodings: **first-child next-sibling** and **Curryfication**

## first-child next-sibling [Rab69, Koc03]

$fcns: \mathcal{T}_\Sigma \rightarrow \mathcal{T}_{\Sigma_f}$  where  $\Sigma_f = (\Sigma \uplus \{\perp\}, \text{ar})$  with  $\text{ar}(a)=2$  for  $a \in \Sigma$  and  $\text{ar}(\perp)=0$

## Curryfication

$curry: \mathcal{T}_\Sigma \rightarrow \mathcal{T}_{\Sigma_c}$  where  $\Sigma_c = (\Sigma \uplus \{\textcircled{\@}\}, \text{ar})$  with  $\text{ar}(a)=0$  for  $a \in \Sigma$  and  $\text{ar}(\textcircled{\@})=2$   
 $t_1 \textcircled{\@} t_2 =$  “add  $t_2$  as last child of the root of  $t_1$ ”



# Binary Encodings

## Automata

We can use tree automata over ranked languages on encoded trees:

	<i>fcns</i>	<i>curry</i>
$\uparrow$ TA	$C_1$	$C_2$
$\downarrow$ TA	$C_3$	$C_4$

These 4 classes are equally expressive.

$C_2$  corresponds to [stepwise tree automata](#) [CNT04].

# Binary Encodings

## Automata

We can use tree automata over ranked languages on encoded trees:

	<i>fcns</i>	<i>curry</i>
$\uparrow$ TA	$C_1$	$C_2$
$\downarrow$ TA	$C_3$	$C_4$

These 4 classes are equally expressive.

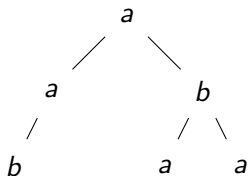
$C_2$  corresponds to [stepwise tree automata](#) [CNT04].

$d\downarrow$ TAs  $\circ \{fcns, curry\}$  define two other classes, that have different expressiveness.  $d\downarrow$ TAs  $\circ fcns$  is the determinism of DTDs.

# Linearization

Let  $\bar{\Sigma} = \{\bar{a} \mid a \in \Sigma\}$ .

$lin(t)$ , the linearization of  $t$ , is the word over  $\Sigma \cup \bar{\Sigma}$  produced by the pre-order traversal of  $t$ :



$$lin(t) = a.a.b.\bar{b}.\bar{a}.b.a.\bar{a}.a.\bar{a}.\bar{b}.\bar{a}$$

This corresponds to the **XML** serialization: `<a><a><b></b>...</b></a>`

# Visibly Pushdown Automata [AM04]

- a pushdown automaton over  $\Sigma \uplus \bar{\Sigma}$
- **visible** means that 1 action (push/pop) is performed by each letter:
  - ▶ rules using  $a \in \Sigma$  only push
  - ▶ rules using  $\bar{a} \in \bar{\Sigma}$  only pop

## Visibly Pushdown Automaton (VPA)

A VPA is a tuple  $(Q, I, F, \Gamma, \Delta, \Sigma \uplus \bar{\Sigma})$  with  $I, F \subseteq Q$  and  $\Delta$  is a set of rules with the form:

- $q_1, a \rightarrow \gamma, q_2$
- $q_1, \bar{a}, \gamma \rightarrow q_2$

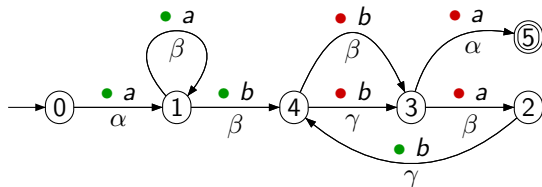
with  $q_1, q_2 \in Q$ ,  $a \in \Sigma$ ,  $\bar{a} \in \bar{\Sigma}$ ,  $\gamma \in \Gamma$ .

The semantics is defined as for usual pushdown automata (ends on final states, not on empty stack).



# VPA as Tree Automata

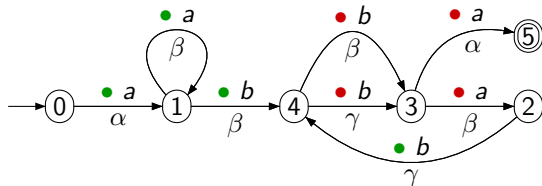
- a run of a VPA on a tree consists in:
  - ▶ when opening node  $\pi$ : update the current state, and assign  $\gamma \in \Gamma$  to  $\pi$
  - ▶ when closing node  $\pi$ : update the current state according to  $\gamma$



A: VPA on  $\mathcal{T}_\Sigma$  with  $\Sigma = \{a, b\}$   
and  $Q = \{0, 1, 2, 3, 4, 5\}$   
and  $\Gamma = \{\alpha, \beta, \gamma\}$

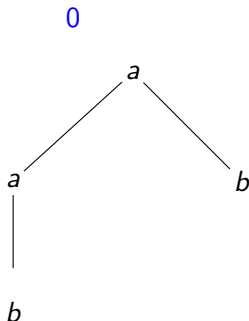
# VPA as Tree Automata

- a run of a VPA on a tree consists in:
  - ▶ when opening node  $\pi$ : update the current state, and assign  $\gamma \in \Gamma$  to  $\pi$
  - ▶ when closing node  $\pi$ : update the current state according to  $\gamma$



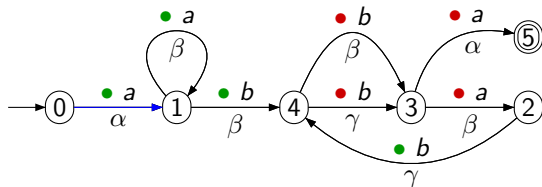
A: VPA on  $\mathcal{T}_\Sigma$  with  $\Sigma = \{a, b\}$   
 and  $Q = \{0, 1, 2, 3, 4, 5\}$   
 and  $\Gamma = \{\alpha, \beta, \gamma\}$

$(0, \emptyset)$



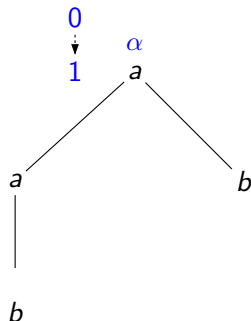
# VPA as Tree Automata

- a run of a VPA on a tree consists in:
  - ▶ when opening node  $\pi$ : update the current state, and assign  $\gamma \in \Gamma$  to  $\pi$
  - ▶ when closing node  $\pi$ : update the current state according to  $\gamma$



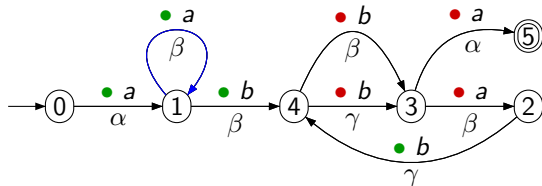
A: VPA on  $\mathcal{T}_\Sigma$  with  $\Sigma = \{a, b\}$   
 and  $Q = \{0, 1, 2, 3, 4, 5\}$   
 and  $\Gamma = \{\alpha, \beta, \gamma\}$

$$(0, \emptyset) \xrightarrow{\bullet a} (1, \alpha)$$



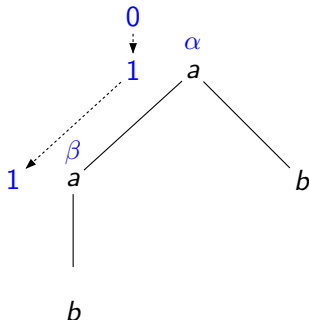
# VPA as Tree Automata

- a run of a VPA on a tree consists in:
  - ▶ when opening node  $\pi$ : update the current state, and assign  $\gamma \in \Gamma$  to  $\pi$
  - ▶ when closing node  $\pi$ : update the current state according to  $\gamma$



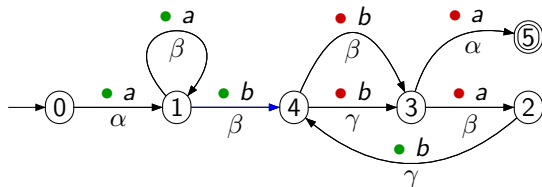
A: VPA on  $\mathcal{T}_\Sigma$  with  $\Sigma = \{a, b\}$   
 and  $Q = \{0, 1, 2, 3, 4, 5\}$   
 and  $\Gamma = \{\alpha, \beta, \gamma\}$

$(0, \emptyset) \xrightarrow{\bullet a} (1, \alpha) \xrightarrow{\bullet a} (1, \alpha.\beta)$



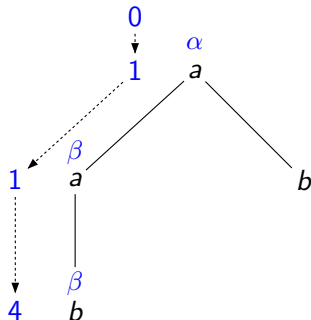
# VPA as Tree Automata

- a run of a VPA on a tree consists in:
  - ▶ when opening node  $\pi$ : update the current state, and assign  $\gamma \in \Gamma$  to  $\pi$
  - ▶ when closing node  $\pi$ : update the current state according to  $\gamma$



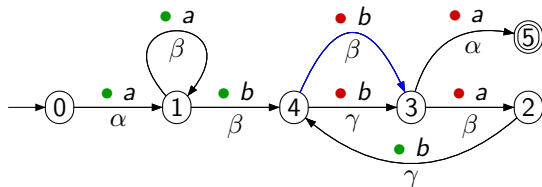
A: VPA on  $\mathcal{T}_\Sigma$  with  $\Sigma = \{a, b\}$   
 and  $Q = \{0, 1, 2, 3, 4, 5\}$   
 and  $\Gamma = \{\alpha, \beta, \gamma\}$

$(0, \emptyset) \xrightarrow{a} (1, \alpha) \xrightarrow{a} (1, \alpha.\beta) \xrightarrow{b} (4, \alpha.\beta.\beta)$



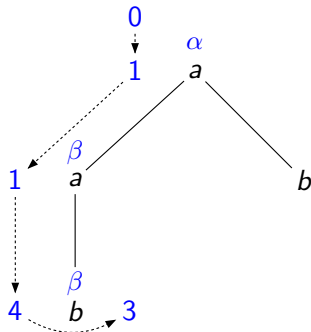
# VPA as Tree Automata

- a run of a VPA on a tree consists in:
  - ▶ when opening node  $\pi$ : update the current state, and assign  $\gamma \in \Gamma$  to  $\pi$
  - ▶ when closing node  $\pi$ : update the current state according to  $\gamma$



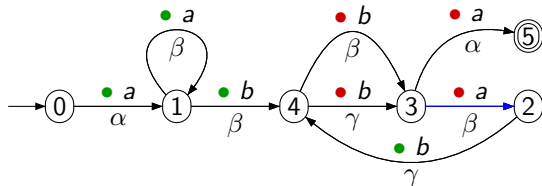
A: VPA on  $\mathcal{T}_\Sigma$  with  $\Sigma = \{a, b\}$   
 and  $Q = \{0, 1, 2, 3, 4, 5\}$   
 and  $\Gamma = \{\alpha, \beta, \gamma\}$

$(0, \emptyset) \xrightarrow{\overset{\bullet}{a}} (1, \alpha) \xrightarrow{\overset{\bullet}{a}} (1, \alpha.\beta) \xrightarrow{\overset{\bullet}{b}} (4, \alpha.\beta.\beta) \xrightarrow{\overset{\bullet}{b}} (3, \alpha.\beta)$



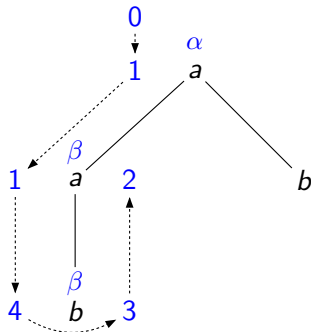
# VPA as Tree Automata

- a run of a VPA on a tree consists in:
  - ▶ when opening node  $\pi$ : update the current state, and assign  $\gamma \in \Gamma$  to  $\pi$
  - ▶ when closing node  $\pi$ : update the current state according to  $\gamma$



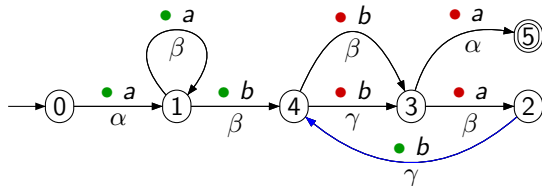
A: VPA on  $\mathcal{T}_\Sigma$  with  $\Sigma = \{a, b\}$   
 and  $Q = \{0, 1, 2, 3, 4, 5\}$   
 and  $\Gamma = \{\alpha, \beta, \gamma\}$

$(0, \emptyset) \xrightarrow{\bullet a} (1, \alpha) \xrightarrow{\bullet a} (1, \alpha.\beta) \xrightarrow{\bullet b} (4, \alpha.\beta.\beta)$   
 $\xrightarrow{\bullet b} (3, \alpha.\beta) \xrightarrow{\bullet a} (2, \alpha)$



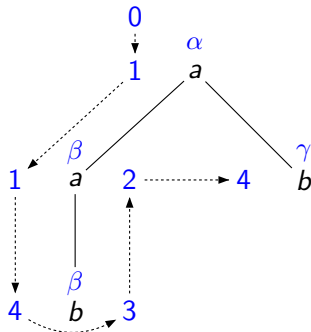
# VPA as Tree Automata

- a run of a VPA on a tree consists in:
  - ▶ when opening node  $\pi$ : update the current state, and assign  $\gamma \in \Gamma$  to  $\pi$
  - ▶ when closing node  $\pi$ : update the current state according to  $\gamma$



A: VPA on  $\mathcal{T}_\Sigma$  with  $\Sigma = \{a, b\}$   
 and  $Q = \{0, 1, 2, 3, 4, 5\}$   
 and  $\Gamma = \{\alpha, \beta, \gamma\}$

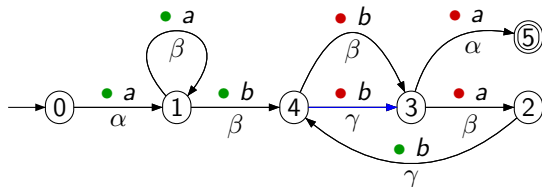
$(0, \emptyset) \xrightarrow{\bullet a} (1, \alpha) \xrightarrow{\bullet a} (1, \alpha.\beta) \xrightarrow{\bullet b} (4, \alpha.\beta.\beta)$   
 $\xrightarrow{\bullet b} (3, \alpha.\beta) \xrightarrow{\bullet a} (2, \alpha) \xrightarrow{\bullet b} (4, \alpha.\gamma)$





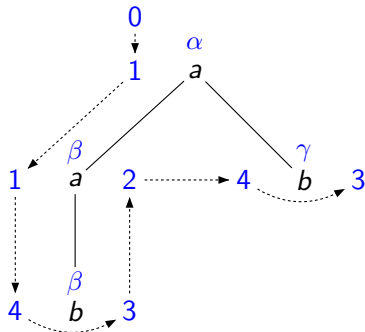
# VPA as Tree Automata

- a run of a VPA on a tree consists in:
  - ▶ when opening node  $\pi$ : update the current state, and assign  $\gamma \in \Gamma$  to  $\pi$
  - ▶ when closing node  $\pi$ : update the current state according to  $\gamma$



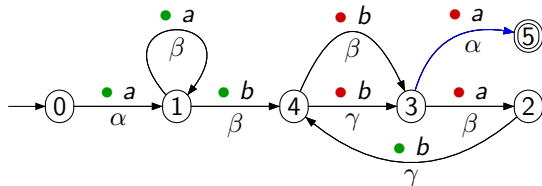
A: VPA on  $\mathcal{T}_\Sigma$  with  $\Sigma = \{a, b\}$   
 and  $Q = \{0, 1, 2, 3, 4, 5\}$   
 and  $\Gamma = \{\alpha, \beta, \gamma\}$

$(0, \emptyset) \xrightarrow{\bullet a} (1, \alpha) \xrightarrow{\bullet a} (1, \alpha.\beta) \xrightarrow{\bullet b} (4, \alpha.\beta.\beta)$   
 $\xrightarrow{\bullet b} (3, \alpha.\beta) \xrightarrow{\bullet a} (2, \alpha) \xrightarrow{\bullet b} (4, \alpha.\gamma)$   
 $\xrightarrow{\bullet b} (3, \alpha)$



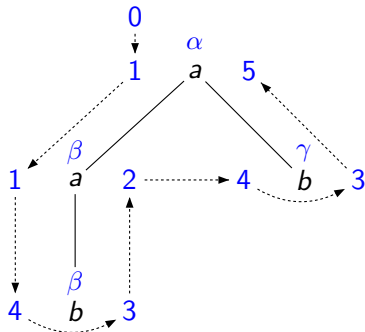
# VPA as Tree Automata

- a run of a VPA on a tree consists in:
  - when opening node  $\pi$ : update the current state, and assign  $\gamma \in \Gamma$  to  $\pi$
  - when closing node  $\pi$ : update the current state according to  $\gamma$



A: VPA on  $\mathcal{T}_\Sigma$  with  $\Sigma = \{a, b\}$   
 and  $Q = \{0, 1, 2, 3, 4, 5\}$   
 and  $\Gamma = \{\alpha, \beta, \gamma\}$

$(0, \emptyset) \xrightarrow{\bullet a} (1, \alpha) \xrightarrow{\bullet a} (1, \alpha.\beta) \xrightarrow{\bullet b} (4, \alpha.\beta.\beta)$   
 $\xrightarrow{\bullet b} (3, \alpha.\beta) \xrightarrow{\bullet a} (2, \alpha) \xrightarrow{\bullet b} (4, \alpha.\gamma)$   
 $\xrightarrow{\bullet b} (3, \alpha) \xrightarrow{\bullet a} (5, \emptyset)$



# VPAs: properties & related models

From now on, we will consider VPAs as [tree automata](#).

Translations between [VPAs](#) and  $\uparrow\text{TA} \circ \text{curry}$  exist [Gau09], so:

- VPAs are as expressive as the 4 classes using  $\{\uparrow\text{TA}, \downarrow\text{TA}\} + \{\text{fcns}, \text{curry}\}$
- [hedge automata](#) also have this expressiveness
- we will call an unranked tree language [recognizable](#) if it belongs to this class
- other equivalent models (in expressiveness):
  - ▶ [nested word automata](#) [Alu07]: reformulation of VPAs
  - ▶ [pushdown forest automata](#) [NS98]: VPAs on forests (i.e. hedges)
  - ▶ [streaming tree automata](#) [GNR08, Gau09]: VPAs on trees

# VPAs: determinism

the class of deterministic VPAs (dVPAs) has numerous interesting properties:

- as expressive as VPAs
  - ▶ so, as **string** acceptors:  $NFAs \subsetneq VPAs = dVPAs \subsetneq dPAs \subsetneq PAs$
- determinization procedure in  $O(2^{|Q|^2})$
- corresponds to **streaming XML** deterministically
  - ▶ yardstick class for streamability

## VPAs: determinization

VPA  $A = (Q, I, F, \Gamma, \Delta, \Sigma \uplus \bar{\Sigma})$  recognizing **encodings of trees**  
(see [AM04] otherwise)

For an hedge  $h$ , let

$$acc_A(h) = \{(q, q') \in Q^2 \mid \text{there is a run of } A \text{ on } h \text{ from } q \text{ to } q'\}$$

The determinization procedure computes  $acc_A(h)$  for all hedges  $h$  of the tree  $t$ . More precisely, the current state at node  $\pi$  is  $acc(h)$  where  $h$  is the hedge of left siblings of  $\pi$ :

- when opening a node,
  - ▶ the previous state is pushed on the stack
  - ▶ the current state is set to the identity of  $Q^2$  ( $h$ =empty hedge)
- when closing a node, the current state is updated from:
  - ▶ the top of the stack, i.e. hedge accessibility before traversing  $\pi$
  - ▶ the previous state, i.e. hedge accessibility through  $t.\pi$

# Determinization

$$Q^{A'} = 2^{Q^A \times Q^A}$$

$$I^{A'} = id_{I^A}$$

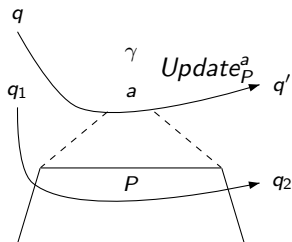
$$F^{A'} = \{P \mid \pi_2(P) \cap F^A \neq \emptyset\}$$

$$\frac{a \in \Sigma \quad P \in Q^{A'}}{P \xrightarrow{\bullet a:P} id_{Q^A} \in \Delta^{A'}}$$

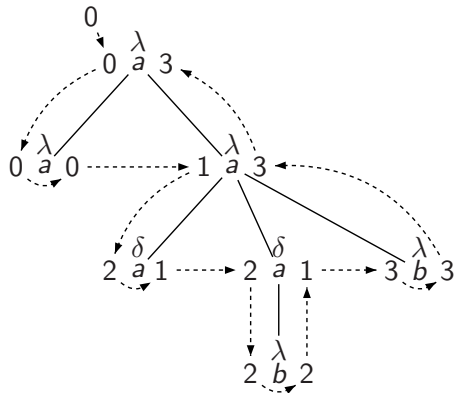
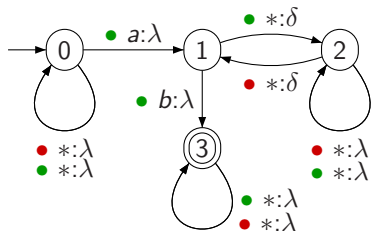
$$\frac{a \in \Sigma \quad P, P' \subseteq Q^A}{P \xrightarrow{\bullet a:P'} P' \circ Update_P^a \in \Delta^{A'}}$$

with

$$Update_P^a = \{(q, q') \mid \exists (q_1, q_2) \in P. \exists \gamma. q \xrightarrow{\bullet a:\gamma} q_1 \in \Delta^A \ \& \ q_2 \xrightarrow{\bullet a:\gamma} q' \in \Delta^A\}$$



# Determinization: example



exercise: run of  $\det(A)$  on  $t$

# Linearization

## Place of recognizable languages

$$NFAs \subsetneq VPAs = dVPAs \subsetneq dPAs \subsetneq PAs = CFLs$$

	closed by			det.	decidable		
	$\cap$	$\cup$	compl.		$\mathcal{L}(A)=\emptyset$	$\mathcal{L}(A)=\sum_{wf}^*$ $\mathcal{L}(A)=\mathcal{L}(B)$	$\mathcal{L}(A)\subseteq\mathcal{L}(B)$
NFAs	✓	✓	✓	✓	✓	✓	✓
VPAs	✓	✓	✓	✓	✓	✓	✓
dPAs	✗	✗	✓	✓	✓	✓	✗
PAs	✗	✓	✗	✗	✓	✗	✗



# Minimization

of unranked tree automata

	unique minimal automaton?	procedure cost	ref.
det. hedge automata using DFAs for horiz. lang.	×	not PTIME unless PTIME=NP	[MN07]
$d \uparrow \text{TA} \circ \text{curry}$ = stepwise tree automata	✓	PTIME	[MN07]
$d \uparrow \text{TA} \circ \text{fcns}$	✓	PTIME	[MN07]
dVPAs	×	open?	[AKMV05] [CW07]

Congruence of a language  $L \subseteq \hat{\Sigma}^*$  with  $\hat{\Sigma} = \Sigma_{\text{push}} \uplus \Sigma_{\text{pop}}$

For well-matched words  $w$  and  $w'$ ,

$$w \equiv_L w' \Leftrightarrow \forall u, v \in \hat{\Sigma}^*, uwv \in L \text{ iff } uw'v \in L$$

A well-matched language  $L \subseteq \hat{\Sigma}^*$  is VPA-recognizable iff  $\equiv_A$  is of finite index. This permits to define canonical VPAs, but not minimal.

## 1 Ranked Trees

- Trees on Ranked Alphabet
- Tree Automata
- Tree Grammars
- Logic

## 2 Unranked Trees

- Unranked Trees
- Automata
- Logic

# MSO

## Equivalence with recognizable tree languages

for unranked trees, use first-child/next-sibling predicates:

$$\Omega_u = \{\text{lab}_a \mid a \in \Sigma\} \cup \{\text{fc}, \text{ns}\}$$

## Equivalence with automata

A tree language  $L \subseteq \mathcal{T}_\Sigma$  is recognizable iff  $\exists \phi \in \text{MSO}[\Omega_u]$  s.t.  $L = \mathcal{L}_\phi$ .

( $\Rightarrow$ ) Similar to the ranked case: define a formula recognizing runs.

( $\Leftarrow$ ) From  $\phi$ , define  $\phi'$  recognizing  $\text{fcns}(\mathcal{L}_\phi)$ . Then use equivalence for ranked trees.

# XPath

XPath  $\rightarrow$  det. automata: in  $O(2^{2^{|\text{el}|}})$



For the fragment  $k$ -Downward XPath, it is in PTIME.

# $k$ -Downward XPath

## Syntax

axis	$d ::= self \mid ch \mid ch^*$
steps	$S ::= d::a \mid d::* \quad (\text{where } a \in \Sigma)$
paths	$P ::= S \mid P[F] \mid P_1/P_2$
filters	$F ::= P \mid \neg F \mid F_1 \wedge F_2$
rooted paths	$R ::= /P$

# $k$ -Downward XPath

## Syntax

axis	$d ::= self \mid ch \mid ch^*$
steps	$S ::= d::a \mid d::* \quad (\text{where } a \in \Sigma)$
paths	$P ::= S \mid P[F] \mid P_1/P_2$
filters	$F ::= P \mid \neg F \mid F_1 \wedge F_2$
rooted paths	$R ::= /P$

## Semantic

$$\llbracket d::* \rrbracket_{path}(t) = d^t$$

$$\llbracket d::a \rrbracket_{path}(t) = \{(\pi, \pi') \in d^t \mid \text{lab}_a^t(\pi')\}$$

$$\llbracket P_1/P_2 \rrbracket_{path}(t) = \llbracket P_1 \rrbracket_{path}(t) \circ \llbracket P_2 \rrbracket_{path}(t)$$

$$\llbracket P[F] \rrbracket_{path}(t) = \{(\pi, \pi') \in \llbracket P \rrbracket_{path}(t) \mid \pi' \in \llbracket F \rrbracket_{filter}(t)\}$$

# k-Downward XPath

## Syntax

axis	$d ::= self \mid ch \mid ch^*$
steps	$S ::= d::a \mid d::* \quad (\text{where } a \in \Sigma)$
paths	$P ::= S \mid P[F] \mid P_1/P_2$
filters	$F ::= P \mid \neg F \mid F_1 \wedge F_2$
rooted paths	$R ::= /P$

## Semantic

$$\llbracket d::* \rrbracket_{path}(t) = d^t$$

$$\llbracket d::a \rrbracket_{path}(t) = \{(\pi, \pi') \in d^t \mid \text{lab}_a^t(\pi')\}$$

$$\llbracket P_1/P_2 \rrbracket_{path}(t) = \llbracket P_1 \rrbracket_{path}(t) \circ \llbracket P_2 \rrbracket_{path}(t)$$

$$\llbracket P[F] \rrbracket_{path}(t) = \{(\pi, \pi') \in \llbracket P \rrbracket_{path}(t) \mid \pi' \in \llbracket F \rrbracket_{filter}(t)\}$$

$$\llbracket P \rrbracket_{filter}(t) = \{\pi \mid \exists \pi'. (\pi, \pi') \in \llbracket P \rrbracket_{path}(t)\}$$

$$\llbracket \neg F \rrbracket_{filter}(t) = \text{nodes} \setminus \llbracket F \rrbracket_{filter}(t)$$

$$\llbracket F_1 \wedge F_2 \rrbracket_{filter}(t) = \llbracket F_1 \rrbracket_{filter}(t) \cap \llbracket F_2 \rrbracket_{filter}(t)$$

$$\llbracket /P \rrbracket_{filter}(t) = \{\pi \mid (\text{root}, \pi) \in \llbracket P \rrbracket_{path}(t)\}$$

# $k$ -Downward XPath

Restrictions:

- $|\text{conjunctions} + \text{filters}| \leq k$
- if  $ch^*::a$  appears, then there are no 2  $a$ -nodes on the same branch

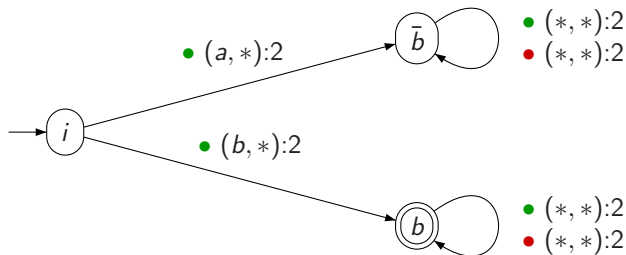


## $k$ -Downward XPath $\rightarrow$ dVPA

- inductive construction
- at each step, automata are **deterministic** and **pseudo-complete**:  
 $\rightarrow$  for every tree  $t$ , there is **exactly** one run on  $t$ .

## Example: $[ch::b]$

First step:  $A_b$  checks whether the root is labeled by  $b$



## Example: $[ch::b]$

Second step:  $A_{ch::b}$  runs  $A_b$  on every child of the root

Procedure:

- add 3 states: start, 0 and 1
- add the rules inferred from what follows
  - ▶ this builds rules for  $F = [ch[F']]$

## Example: $[ch::b]$

Second step:  $A_{ch::b}$  runs  $A_b$  on every child of the root

Procedure:

- add 3 states: start, 0 and 1
- add the rules inferred from what follows
  - ▶ this builds rules for  $F = [ch[F']]$

$$\frac{a \in \Sigma \quad V \in \{0, 1\}}{\text{start} \xrightarrow{\bullet (a, V): 0} 0}$$

opening the root:  
move to 0

$$\frac{q_1 \xrightarrow{\bullet (a, V): \gamma} q_2 \in \Delta^{A'} \quad q_1 \in I^{A'} \quad b \in \{0, 1\}}{q_1 \xrightarrow{\bullet (a, V): b} q_2}$$

opening a child:  
start testing  $F'$

$$\frac{q_1 \xrightarrow{\alpha (a, V): \gamma} q_2 \in \Delta^{A'}}{q_1 \xrightarrow{\alpha (a, V): \gamma} q_2}$$

run test of  $F'$

# Example: $[ch::b]$

$$\begin{array}{c}
 q_1 \xrightarrow{\bullet (a,V):\gamma} q_2 \in \Delta^{A'} \quad q_2 \notin F^{A'} \\
 q'_1 \xrightarrow{\bullet (a,V):\gamma} q'_2 \in \Delta^{A'} \quad q'_1 \in I^{A'} \quad b \in \{0, 1\} \\
 \hline
 q_1 \xrightarrow{\bullet (a,V): b} b
 \end{array}$$

failure of  $F'$ :  
no new match

$$\begin{array}{c}
 q_1 \xrightarrow{\bullet (a,V):\gamma} q_2 \in \Delta^{A'} \quad q_2 \in F^{A'} \\
 q'_1 \xrightarrow{\bullet (a,V):\gamma} q'_2 \in \Delta^{A'} \quad q'_1 \in I^{A'} \quad b \in \{0, 1\} \\
 \hline
 q_1 \xrightarrow{\bullet (a,V): b} 1
 \end{array}$$

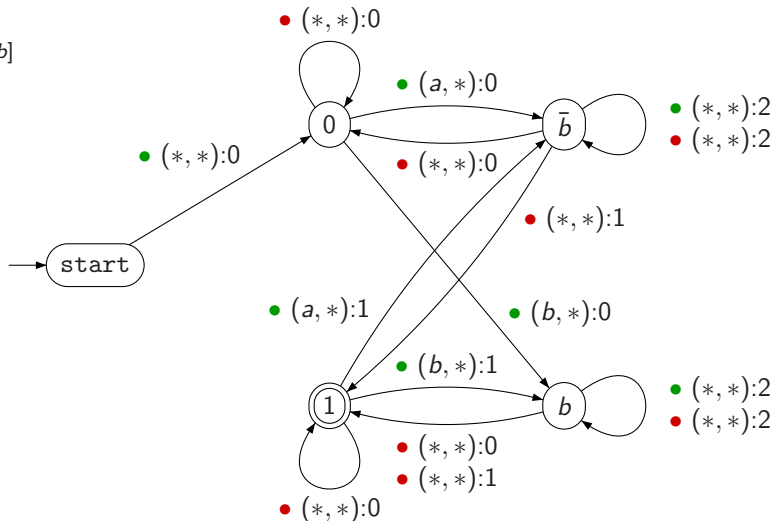
success of  $F'$ :  
move to 1

$$\begin{array}{c}
 a \in \Sigma \quad V \in \{0, 1\} \quad b \in \{0, 1\} \\
 \hline
 b \xrightarrow{\bullet (a,V):0} b
 \end{array}$$

closing the root

# Example: $[ch::b]$

$A_{[ch::b]}$



# $k$ -Downward XPath $\rightarrow$ dVPAs

other steps

- $ch^*$  is similar to  $ch$
- $P_1/P_2$  is transformed into filter (w/ variables):  
 $ch^*::a/ch::b$  becomes  $[ch^*::a[ch::b[x]]]$
- $A_{F_1 \wedge F_2} = A_{F_1} \wedge A_{F_2}$  (product construction)
- $A_{\neg F} = \neg A_F$

# Other logics

- Conditional XPath, Regular XPath
- other modal logics: Tree TL, CTL\*, PDL<sub>tree</sub>
- $\mu$ -calculus



# Language, Automata and Logic for Finite Trees

Olivier Gauwin

UMons

Feb/March 2010

# References

- [AKMV05] Rajeev Alur, Viraj Kumar, P. Madhusudan, and Mahesh Viswanathan.  
Congruences for visibly pushdown languages.  
*In Automata, Languages and Programming, 32nd International Colloquium*, volume 3580 of *Lecture Notes in Computer Science*, pages 1102–1114. Springer Verlag, 2005.
- [Alu07] Rajeev Alur.  
Marrying words and trees.  
*In 26th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*, pages 233–242. ACM-Press, 2007.
- [AM04] Rajeev Alur and P. Madhusudan.  
Visibly pushdown languages.  
*In 36th ACM Symposium on Theory of Computing*, pages 202–211. ACM-Press, 2004.
- [BKWM01] Anne Brüggemann-Klein, Derick Wood, and Makoto Murata.

Regular tree and regular hedge languages over unranked alphabets: Version 1, April 07 2001.

- [BLN07] Michael Benedikt, Leonid Libkin, and Frank Neven.  
Logical definability and query languages over ranked and unranked trees.  
*ACM Transactions on Computational Logics*, 8(2), April 2007.
- [BPSM<sup>+</sup>08] Tim Bray, Jean Paoli, C. M. Sperberg-McQueen, Eve Maler, and François Yergeau.  
Extensible Markup Language (XML) 1.0 (Fifth Edition), November 2008.  
<http://www.w3.org/TR/2008/REC-xml-20081126/>.
- [CDG<sup>+</sup>07] Hubert Comon, Max Dauchet, Rémi Gilleron, Christof Löding, Florent Jacquemard, Denis Lugiez, Sophie Tison, and Marc Tommasi.  
Tree automata techniques and applications.  
Available online since 1997:  
<http://tata.gforge.inria.fr>, October 2007.

Revised October, 12th 2007.

- [CNT04] Julien Carme, Joachim Niehren, and Marc Tommasi.  
Querying unranked trees with stepwise tree automata.  
*In 19th International Conference on Rewriting Techniques and Applications*, volume 3091 of *Lecture Notes in Computer Science*, pages 105–118. Springer Verlag, 2004.
- [CW07] Patrick Chervet and Igor Walukiewicz.  
Minimizing variants of visibly pushdown automata.  
*In Mathematical Foundations of Computer Science*, volume 4708 of *Lecture Notes in Computer Science*, pages 135–146. Springer Verlag, 2007.
- [Don70] John E. Doner.  
Tree acceptors and some of their applications.  
4:406–451, 1970.
- [Gau09] Olivier Gauwin.  
*Streaming Tree Automata and XPath*.  
PhD thesis, Université Lille 1, 2009.

- [GNR08] Olivier Gauwin, Joachim Niehren, and Yves Roos.  
Streaming tree automata.  
*Information Processing Letters*, 109(1):13–17, December 2008.
- [Koc03] Christoph Koch.  
Efficient processing of expressive node-selecting queries on XML data in secondary storage: A tree automata-based approach.  
*In Proc. VLDB 2003*, 2003.
- [MN07] Wim Martens and Joachim Niehren.  
On the minimization of XML schemas and tree automata for unranked trees.  
*Journal of Computer and System Science*, 73(4):550–583, 2007.
- [NS98] Andreas Neumann and Helmut Seidl.  
Locating matches of tree patterns in forests.  
*In 18th Conference on Foundations of Software Technology and Theoretical Computer Science*, volume 1530 of *Lecture*

*Notes in Computer Science*, pages 134–145. Springer Verlag, 1998.

[Rab69]

Michael O. Rabin.

Decidability of Second-Order Theories and Automata on Infinite Trees.

*Transactions of the American Mathematical Society*, 141:1–35, 1969.

[TW68]

J. W. Thatcher and J. B. Wright.

Generalized finite automata with an application to a decision problem of second-order logic.

*Mathematical System Theory*, 2:57–82, 1968.