

# Recapitulation: Hedge Automata for XML Languages

$$\Sigma = \{a, b, c\}$$

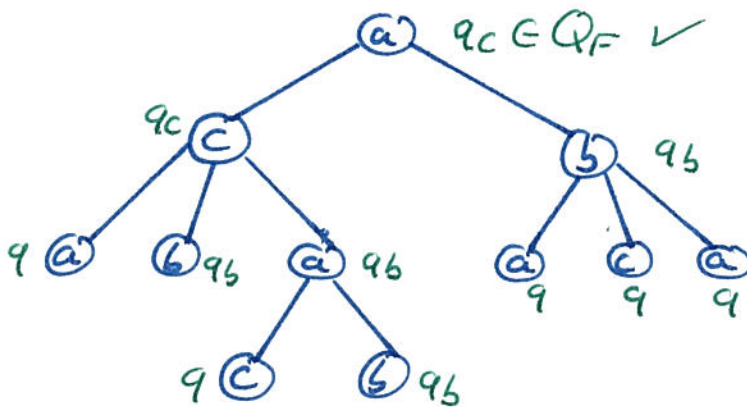
$$\mathcal{H} = (Q, \rightarrow, Q_F), \quad Q = \{q, q_b, q_c\}, \quad Q_F = \{q_c\}$$

Transitions:

$$Q^* \xrightarrow{a|c} q \quad Q^* q_b Q^* \xrightarrow{a|c} q_b \quad Q^* q_c Q^* \xrightarrow{a|b|c} q_c$$

$$Q^* \xrightarrow{b} q_b \quad Q^* q_b Q^* q_b Q^* \xrightarrow{c} q_c$$

Run on



## 9.2 Membership

Goal: Given a DTD  $D$  and an XML document  $t$ .  
Decide  $t \in L(D)$  ( $= L(\mathcal{H}_D)$ ).

Theorem:

Membership for NHTFs with horizontal languages represented by NFRs can be decided in polynomial time (in the size of the tree and the size of the automaton).

Algorithm:

Given: NHTF  $\mathcal{H} = (Q, \rightarrow, Q_F)$

• tree  $t: T \rightarrow \Sigma$ .

Construct a mapping

$$S: T \rightarrow \mathbb{P}(Q). \quad (S \text{ like run})$$

↳ It labels each node  $w \in T$   
with the set of states  $S(w) \subseteq Q$   
that are reachable at  $w$   
in a run of  $\mathcal{N}$  on  $\epsilon$ .

↳ To this end, pick a node  $w \in T$

- whose children  $w.0, \dots, w.(n-1)$  we already labelled  
by  $S(w.0) = Q_0, \dots, S(w.(n-1)) = Q_{n-1}$ .
- Check for each transition  $R \xrightarrow{a} q$  with  $a = t(w)$   
whether it can be applied to a string  $q_0 \dots q_{n-1}$   
obtained from selecting  $q_i \in Q_i$ .

The algorithm:

input: NFA  $\mathcal{N} = (Q, \rightarrow, Q_{\text{in}})$ , tree  $t: T \rightarrow \Sigma$

begin:

set  $S(w) := \perp$  f.a.  $w \in T$ .

while ( $\exists w \in T$  with  $S(w) = \perp$  and

$S(w.i) \neq \perp$  f.a. children  $w.0, \dots, w.(n-1)$  of  $w$ ) {

$M := \emptyset$

for each transition  $R \xrightarrow{a} q$  with  $a = t(w)$  {

if ( $\exists q_0 \in S(w.0), \dots, q_{n-1} \in S(w.(n-1))$   
with  $q_0 \dots q_{n-1} \in R$ ) {

$M := M \cup \{q\}$

}

}

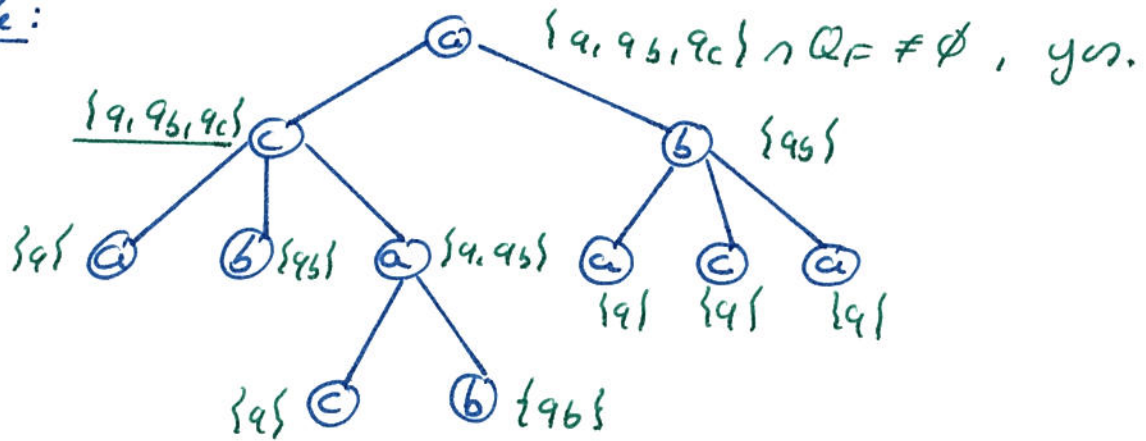
$S(w) := M$

}

output "yes" if  $S(\epsilon) \cap Q_{\text{in}} \neq \emptyset$ , "no" otherwise.

end

Example:



Why is  $\odot$  labelled by  $\{a, q_b, q_c\}$ ?

$Q^* \xrightarrow{a} q$  applies  $\rightsquigarrow M = \emptyset \cup \{q\}$

$Q^* q_b Q^* \xrightarrow{q_b} q_b$  applies  $\rightsquigarrow M = \{q\} \cup \{q_b\}$

$Q^* q_b Q^* q_b Q^* \xrightarrow{q_c} q_c$  applies  $\rightsquigarrow M = \{a, q_b\} \cup \{q_c\}$  ✓

The only nonpolynomial check:

↳ May be if  $(\exists q_0 \in S(w, 0), \dots$   
inside the loop.

↳ Works in polynomial time as follows:

- Start in initial state of the NFA for the horizontal language  $R$ .

- Collect all NFA states reachable (in a single step) with NFA states in  $S(w, 0) = Q_0$ .

transition labels

- Compute from these all NFA states reachable with elements in  $S(w, 1) = Q_1$ .

- etc.

- If last set contains a final state of the NFA, condition is satisfied.

↳ Algorithm can be understood as

"powerset construction along  $Q_0, \dots, Q_{n-1}$ "

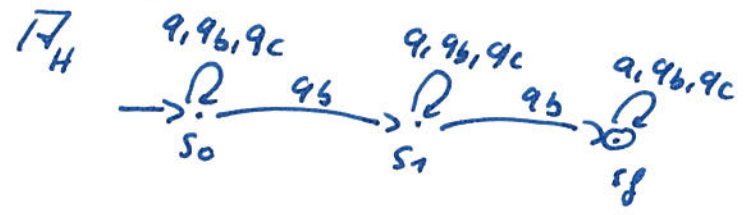


Example (continued):

$Q_0$     $Q_1$     $Q_2$   
 "   "   "  
 "   "   "

- Consider node  $\odot$  with child sequence  $\{a\}, \{ab\}, \{a, ab\}$ .
- Consider transition  $Q^* ab Q^* ab Q^* \rightarrow c qc$ .

↳ The regular expression is represented by

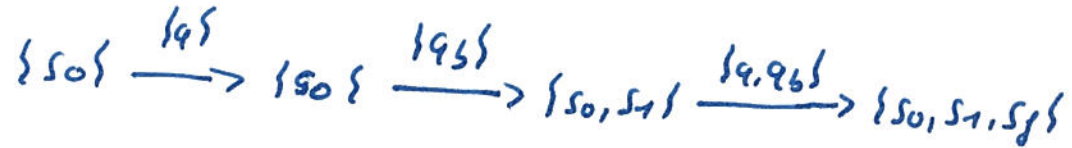


↳ To check whether

$$\exists q_0 \in \{a\}, q_1 \in \{ab\}, q_2 \in \{a, ab\}$$

$$\text{with } q_0 q_1 q_2 \in L(\mathcal{R}) = Q^* ab Q^* ab Q^*$$

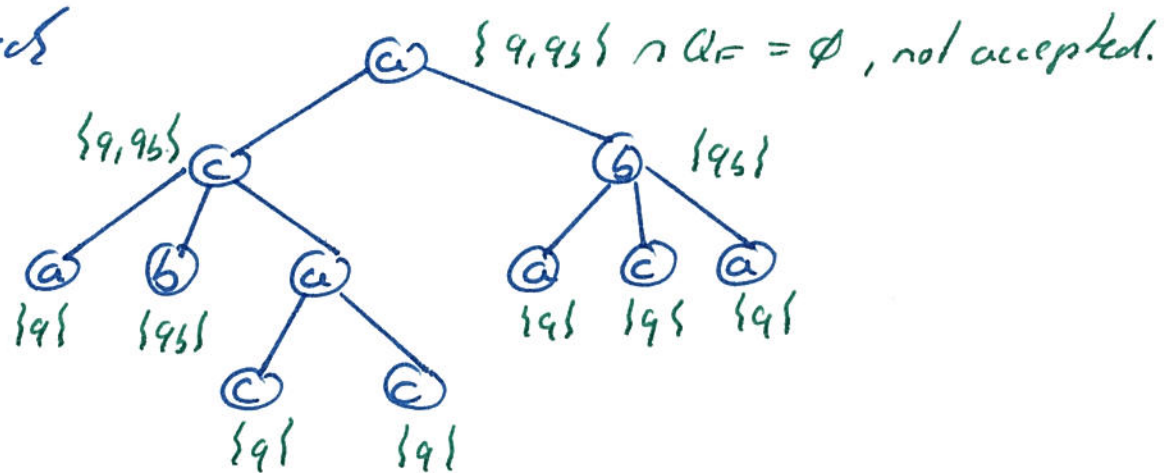
apply the above procedure:



Since  $s_2 \in \{s_0, s_1, s_2\}$ , add  $qc$  to  $M$ .

Example:

Type check



More on XML documents:

- ↳ Work on data inside documents
- ↳ Adapt schema representations of different companies
  - Deal with inclusion, renaming, structure changes
  - Algorithmic support is needed.