# 9. XML Schema Languages

↳ XML document

```
<lecture>
    <title> Applied Automata Theory </title>
    <block>
        <title> Finite Words </title>
        <topic>
            <title> WMSO </title>
            <goal> Satisfiability </goal>
            <approach> Buechi </approach>
        </topic>
    </block>
</lecture>
```
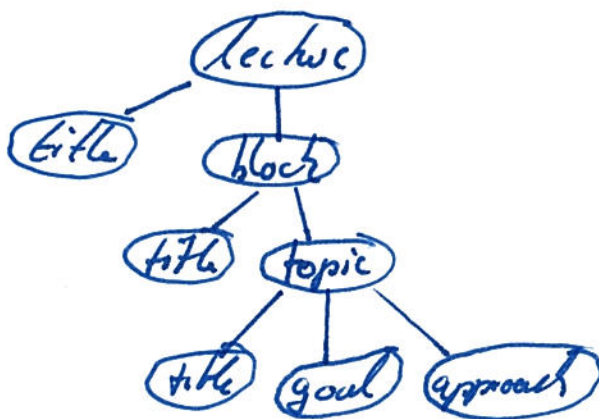
yields a tree that
- reflects structure of the document
- without data.



Goal: Pose requirements on the structure of documents
   ↳ Every lecture is split into blocks
   ↳ Blocks are divided into topics

Observation: ↳ Requirements describe a tree language
      over the alphabet of tags
     ↳ Such a description is called a _schema_
     ↳ Document is _valid_ wrt. a schema
       if it belongs to the tree language defined by the schema

Several XML schema languages exist:

↳ Document Type Definitions (here), XML Schema, Relax NG

↳ We are interested in

    ↳ Connection to automata theory

    (↳ Expressiveness)

    ↳ <u>Algorithmic problems</u>

        • Is a document <u>valid wrt. a schema</u>? (here)

        (membership in the language)

        • <u>Is there a document</u> that is valid for this schema?

          → Sanity check

          → Emptiness in language theory

            (used as subproblem for inclusion)

        • Are all documents <u>valid wrt. one schema</u>

                    <u>valid for another schema</u>?

          → Needed when merging archives/companies

          → Inclusion in language theory.

## 9.1 Document Type Definitions and Hedge Automata

↳ A <u>document type definition (DTD)</u> is a context-free grammar
with regular expressions on the right hand side

↳ Tree language of this grammar = all derivation trees.

Example:

```
<! DOCTYPE LECTURE [
    <!ELEMENT lecture    (title, (block + | (topic, exercise?)+))>
    <!ELEMENT block      (title, (topic, exercise?)+)>
    <!ELEMENT topic      (title, goal, problem?, approach)>
    <!ELEMENT title      (# PCDATA)>
    ...
]>
```

                                      <u>content model</u>

On the right hand side:

  | = choice       + = one or more occurrences

  , = sequence      ? = zero or one occurrences

  #PCDATA = arbitrary character sequence

        (parsed character data, for data inside document)

Corresponding CF-grammar:

  $lecture \longrightarrow title.(block^+ + (topic.(exercise + \varepsilon))^+)$

  $block \longrightarrow title.(topic.(exercise + \varepsilon))^+$

  $topic \longrightarrow title.goal.(problem + \varepsilon).approach$

  $title \longrightarrow \varepsilon$

  ...

To define the tree language described by a DTD, need hedge automata.

## 9.1.1 Unranked Trees and Hedge Automata

- Drop the restriction of ranked alphabets for trees
  - ↳ Remember • Letters in $\Sigma$ have ranks
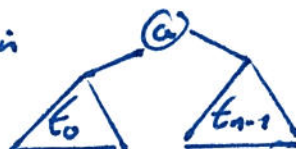    - • Trees $t: T \rightarrow \Sigma$ obey these ranks
- Consider underlined <u>unranked alphabet</u> $\Sigma$
  - ↳ Each node has arbitrarily but finitely many children
  - ↳ Tree $t: T \rightarrow \Sigma$ without further constraints

    called an <u>unranked tree</u>

- ↳ Call $t_0,...,t_{n-1}$ in  a <u>hedge</u>

- <u>Hedge automata</u> process unranked trees <u>bottom-up</u>

  Goal: Obtain similar properties as for ranked trees

  Problem: Number of successors of a node is not bounded (but finite)
  (unbounded branching)

↳ Transitions cannot be listed
↳ Represent _symbolically_ infinite number of transitions.

Example:

↳ Accept all trees of height 2 where
- root is labelled by "a"
- all children labelled by "b"
- number of children even



↳ Use · two states $q, q_b$
- transitions $\longrightarrow_b q_b$

$$\longrightarrow_a q, \quad (q_b, q_b) \longrightarrow_a q, \quad (q_b, q_b, q_b, q_b) \longrightarrow_a q, \dots$$

↳ Infinite set of transitions

can be represented by $(q_b q_b)^* \longrightarrow_a q$

More generally,

$$R \longrightarrow_a q \qquad \text{where } R \subseteq Q^*$$

is a _regular language_ over the states of the automaton.

Definition (Hedge automaton):

- A _(nondeterministic)_ hedge automaton _(NHA)_ over $\Sigma$

is a tuple $A = (Q, \longrightarrow, Q_F)$ where

$$\longrightarrow \; \subseteq \; \mathbb{P}(Q^*) \times \Sigma \times Q$$

with $R \subseteq Q^*$ on lhs of transitions _regular_.

These $R$ are called _horizontal languages_.

- A _run_ of _A_ on $t : T \to \Sigma$

is a function

$$r : T \to Q$$

so that for all $w \in T$ with $r(w) = q, t(w) = a$, and $n = \#$ successors of $w$
we have a transition $R \longrightarrow_a q$ with $r(w.0) \dots r(w.(n-1)) \in R$.

To apply a transition $R \to_a q$ at a leaf, we require $\varepsilon \in R$.

- Run is <u>accepting</u> if $r(\varepsilon) \in Q_F$

<u>Language</u> of $A$ is

$$L(A) := \{\, t: T \to \Sigma \mid \text{there is an accepting run} \atop \text{of } A \text{ on } t \,\}$$

<u>Example:</u>

- There are two nodes labelled "$b$" whose greatest common ancestor is a "$c$".
  
  (longest common prefix)

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

  $A = \{Q, \to, Q_F\}$, $Q = \{q, q_b, q_c\}$, $Q_F = \{q_c\}$

  Transitions:
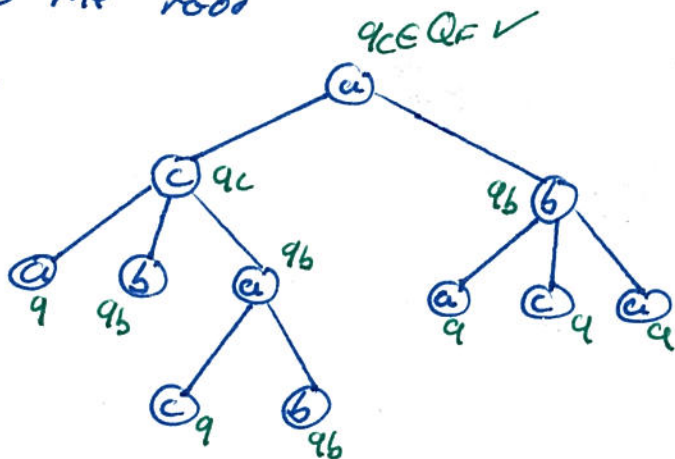
  $$Q^* \longrightarrow_{a/c} q \qquad Q^* q_b Q^* \longrightarrow_{a/c} q_b \qquad Q^* q_c Q^* \longrightarrow_{a/b/c} q_c$$

  $$Q^* \longrightarrow_b q_b \qquad Q^* q_b Q^* q_b Q^* \longrightarrow_c q_c$$

- Label $b$ nodes by $q_b$
- As long as condition has not been satisfied, pass $q_b$ upwards
- Until you find a $c$ node with two $q_b$ children
- Label this node by final state $q_c$
- Pass to the root

Run on

## 9.1.2 Document Type Definitions

### Definition (DTD):

A <u>document type definition (DTD)</u> is a tuple $D = (\Sigma, s, \delta)$ with

- <u>start symbol</u> $s$
- function $\delta$ assigning each element $a \in \Sigma$ a regular expression $\delta(a)$ over $\Sigma$.

To define the language of a DTD, understand it as a hedge automaton

$$H_D := (Q, \to, Q_F)$$

with

$$Q := \{ q_a \mid a \in \Sigma \} \qquad // \text{ State for each letter (tag)}$$
$$Q_F := \{ q_s \}$$

To define transitions, understand

$$L(\delta(a)) \subseteq \Sigma^* \quad \text{as subset of } Q^*$$

(by taking $a_1 \ldots a_n$ as $q_{a_1} \ldots q_{a_n}$).

Then

$$L(\delta(a)) \to_a q_a \qquad f.a. \; a \in \Sigma$$

The <u>language of the DTD</u> is $L(D) := L(H_D)$.

Note that $H_D$ is <u>deterministic</u> in the following sense:

for all $R_1 \to_a q_1$ and $R_2 \to_a q_2$ we have

$$R_1 \cap R_2 \neq \emptyset \quad \text{implies} \quad q_1 = q_2.$$

### Example:

Reconsider $\quad$ lecture $\longrightarrow$ title. $(block^* + (topic.(exercise + \epsilon))^*)$

$\ldots$

The corresponding hedge automaton is
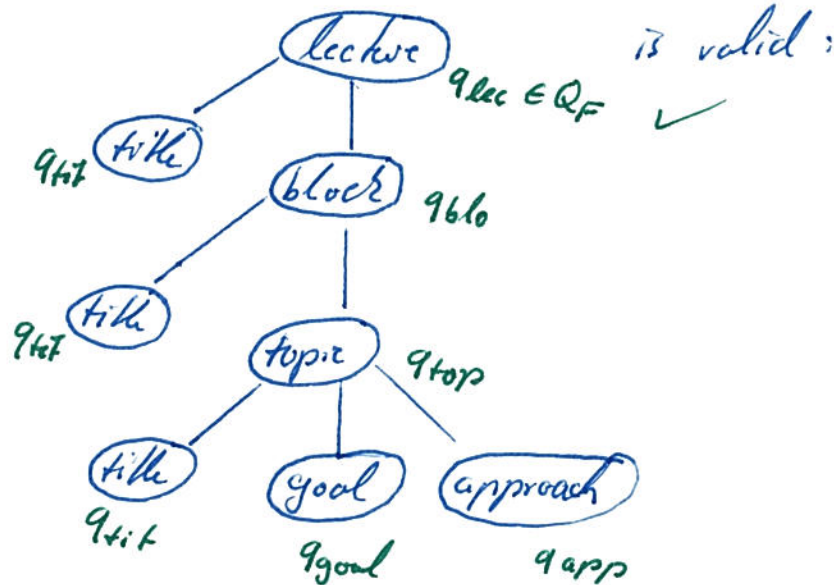
$\mathcal{A}_{LEC} = (Q, \to, \{q_{lec}\})$

with

$Q = \{q_{lec}, q_{tit}, q_{blo}, q_{ex}, q_{top}, q_{goal}, q_{pro\delta}, q_{app}\}$,

and

$q_{tit} \cdot (q_{blo}{}^+ + (q_{top} \cdot (q_{ex} + \varepsilon))^+) \to_{lec} q_{lec}$.

...

Check that



is valid:

$q_{lec} \in Q_F$ ✓

How to check validity systematically?
→ Algorithm
→ With low complexity.