

Tree automata techniques for the verification of infinite state-systems



Summer School VTSA 2011

Florent Jacquemard

INRIA Saclay & LSV (UMR CNRS/ENS Cachan)

`florent.jacquemard@inria.fr`

`http://www.lsv.ens-cachan.fr/~jacquema`

TATA book

<http://tata.gforge.inria.fr>

(chapters 1, 3, 7, 8)



**Tree
Automata
Techniques and
Applications**

HUBERT COMON

MAX DAUCHET

RÉMI GILLERON

FLORENT JACQUEMARD

DENIS LUGIEZ

CHRISTOF LÖDING

SOPHIE TISON MARC TOMMASI

Finite tree automata

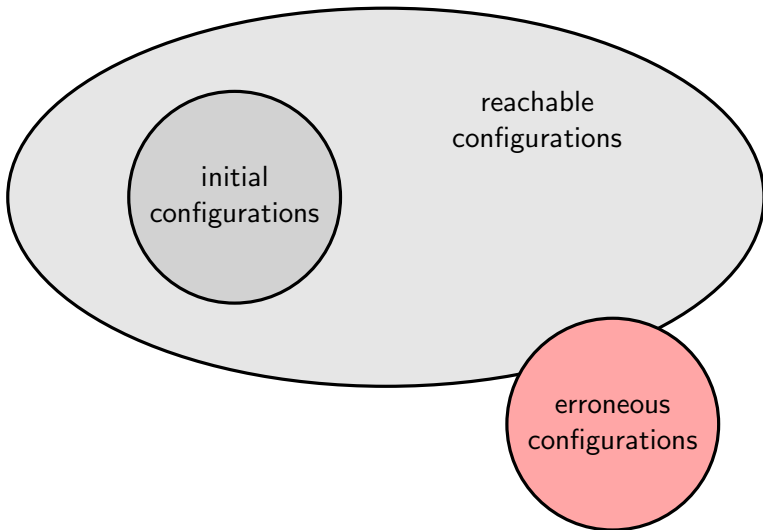
- ▶ tree recognizers
- ▶ generalize NFA from words to trees
- = finite representations of infinite set of labeled trees

are a useful tool for verification procedures

- ▶ composition results
 - ▶ closure under Boolean operations
 - ▶ closure under transformations
- ▶ decision results, efficient algorithms
- ▶ expressiveness, close relationship with logic

Verification of infinite state systems

regular model checking : static analysis of safety properties for infinite state systems, using symbolic reachability verification techniques.



Concurrent readers/writers

Example from [Clavel et al. LNCS 4350 2007]

1. $\text{state}(0, 0) = \text{state}(0, s(0))$
2. $\text{state}(r, 0) = \text{state}(s(r), 0)$
3. $\text{state}(r, s(w)) = \text{state}(r, w)$
4. $\text{state}(s(r), w) = \text{state}(r, w)$

- ▶ writers can access the file if nobody else is accessing it (1)
- ▶ readers can access the file if no writer is accessing it (2)
- ▶ readers and writers can leave the file at any time (3,4)

Properties expected:

- ▶ mutual exclusion between readers and writers
- ▶ mutual exclusion between writers

Concurrent readers/writers: reachable configurations

1. $\text{state}(0, 0) = \text{state}(0, s(0))$
2. $\text{state}(r, 0) = \text{state}(s(r), 0)$
3. $\text{state}(r, s(w)) = \text{state}(r, w)$
4. $\text{state}(s(r), w) = \text{state}(r, w)$

Initial configuration: $\text{state}(0, 0)$

Concurrent readers/writers: reachable configurations

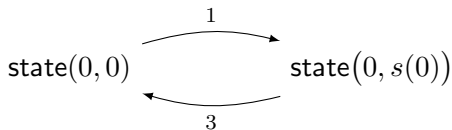
1. $\text{state}(0, 0) = \text{state}(0, s(0))$
2. $\text{state}(r, 0) = \text{state}(s(r), 0)$
3. $\text{state}(r, s(w)) = \text{state}(r, w)$
4. $\text{state}(s(r), w) = \text{state}(r, w)$

Reachable configurations: $\text{state}(0, 0)$

Concurrent readers/writers: reachable configurations

1. $\text{state}(0, 0) = \text{state}(0, s(0))$
2. $\text{state}(r, 0) = \text{state}(s(r), 0)$
3. $\text{state}(r, s(w)) = \text{state}(r, w)$
4. $\text{state}(s(r), w) = \text{state}(r, w)$

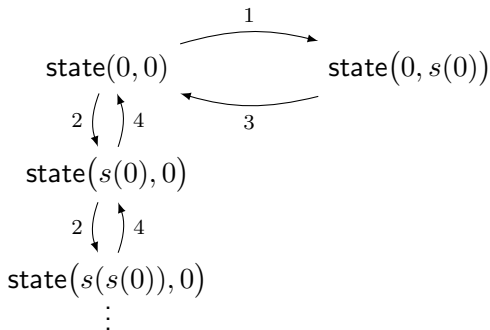
Reachable configurations:



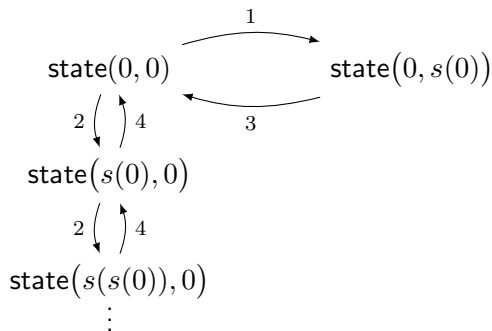
Concurrent readers/writers: reachable configurations

1. $\text{state}(0, 0) = \text{state}(0, s(0))$
2. $\text{state}(r, 0) = \text{state}(s(r), 0)$
3. $\text{state}(r, s(w)) = \text{state}(r, w)$
4. $\text{state}(s(r), w) = \text{state}(r, w)$

Reachable configurations:



Concurrent readers/writers: finite representation



$q_0 := 0$

$q := \text{state}(q_0, q_0) \mid \text{state}(q_0, q_1) \mid \text{state}(q_1, q_0) \mid \text{state}(q_2, q_0)$

$q_1 := s(q_0)$

$q_2 := s(q_1) \mid s(q_2)$

Concurrent readers/writers: automata construction

$$1. \text{ state}(0, 0) = \text{state}(0, s(0))$$

$$2. \text{ state}(r, 0) = \text{state}(s(r), 0)$$

$$3. \text{ state}(r, s(w)) = \text{state}(r, w)$$

$$4. \text{ state}(s(r), w) = \text{state}(r, w)$$

$$q_0 := 0$$

$$q := \text{state}(q_0, q_0)$$

Concurrent readers/writers: automata construction

1. $\text{state}(0, 0) = \text{state}(0, s(0))$
 $\text{state}(0, 0) \in q \Rightarrow \text{state}(0, s(0)) \in q$
2. $\text{state}(r, 0) = \text{state}(s(r), 0)$
3. $\text{state}(r, s(w)) = \text{state}(r, w)$
4. $\text{state}(s(r), w) = \text{state}(r, w)$

$$q_0 := 0$$

$$q := \text{state}(q_0, q_0)$$

Concurrent readers/writers: automata construction

1. $\text{state}(0, 0) = \text{state}(0, s(0))$
 $\text{state}(0, 0) \in q \Rightarrow \text{state}(0, s(0)) \in q$
2. $\text{state}(r, 0) = \text{state}(s(r), 0)$
3. $\text{state}(r, s(w)) = \text{state}(r, w)$
4. $\text{state}(s(r), w) = \text{state}(r, w)$

$$q_0 := 0$$

$$q := \text{state}(q_0, q_0) \mid \text{state}(q_0, q_1)$$

$$q_1 := s(q_0)$$

Concurrent readers/writers: automata construction

1. $\text{state}(0, 0) = \text{state}(0, s(0))$
2. $\text{state}(r, 0) = \text{state}(s(r), 0)$
 $\text{state}(q_0, 0) \in q \Rightarrow \text{state}(s(q_0), 0) \in q$
3. $\text{state}(r, s(w)) = \text{state}(r, w)$
4. $\text{state}(s(r), w) = \text{state}(r, w)$

$$q_0 := 0$$

$$q := \text{state}(q_0, q_0) \mid \text{state}(q_0, q_1)$$

$$q_1 := s(q_0)$$

Concurrent readers/writers: automata construction

1. $\text{state}(0, 0) = \text{state}(0, s(0))$
2. $\text{state}(r, 0) = \text{state}(s(r), 0)$
 $\text{state}(q_0, 0) \in q \Rightarrow \text{state}(s(q_0), 0) \in q$
3. $\text{state}(r, s(w)) = \text{state}(r, w)$
4. $\text{state}(s(r), w) = \text{state}(r, w)$

$$q_0 := 0$$

$$q := \text{state}(q_0, q_0) \mid \text{state}(q_0, q_1) \mid \text{state}(q_1, q_0)$$

$$q_1 := s(q_0)$$

Concurrent readers/writers: automata construction

1. $\text{state}(0, 0) = \text{state}(0, s(0))$
2. $\text{state}(r, 0) = \text{state}(s(r), 0)$
 $\text{state}(q_1, 0) \in q \Rightarrow \text{state}(s(q_1), 0) \in q$
3. $\text{state}(r, s(w)) = \text{state}(r, w)$
4. $\text{state}(s(r), w) = \text{state}(r, w)$

$$q_0 := 0$$

$$q := \text{state}(q_0, q_0) \mid \text{state}(q_0, q_1) \mid \text{state}(q_1, q_0)$$

$$q_1 := s(q_0)$$

Concurrent readers/writers: automata construction

1. $\text{state}(0, 0) = \text{state}(0, s(0))$
2. $\text{state}(r, 0) = \text{state}(s(r), 0)$
 $\text{state}(q_1, 0) \in q \Rightarrow \text{state}(s(q_1), 0) \in q$
3. $\text{state}(r, s(w)) = \text{state}(r, w)$
4. $\text{state}(s(r), w) = \text{state}(r, w)$

$$q_0 := 0$$

$$q := \text{state}(q_0, q_0) \mid \text{state}(q_0, q_1) \mid \text{state}(q_1, q_0) \mid \text{state}(q_2, q_0)$$

$$q_1 := s(q_0)$$

System Timbuk [Thomas Genet]. Automated construction, with guess of *acceleration* $q_2 := s(q_2)$ by user assistance.

Concurrent readers/writers: automata construction

1. $\text{state}(0, 0) = \text{state}(0, s(0))$
2. $\text{state}(r, 0) = \text{state}(s(r), 0)$
 $\text{state}(q_2, 0) \in q \Rightarrow \text{state}(s(q_2), 0) \in q$
3. $\text{state}(r, s(w)) = \text{state}(r, w)$
4. $\text{state}(s(r), w) = \text{state}(r, w)$

$$q_0 := 0$$

$$q := \text{state}(q_0, q_0) \mid \text{state}(q_0, q_1) \mid \text{state}(q_1, q_0) \mid \text{state}(q_2, q_0)$$

$$q_1 := s(q_0)$$

System Timbuk [Thomas Genet]. Automated construction, with guess of *acceleration* $q_2 := s(q_2)$ by user assistance.

Concurrent readers/writers: automata construction

1. $\text{state}(0, 0) = \text{state}(0, s(0))$
2. $\text{state}(r, 0) = \text{state}(s(r), 0)$
3. $\text{state}(r, s(w)) = \text{state}(r, w)$
 $\text{state}(q_0, s(q_0)) \in q \Rightarrow \text{state}(q_0, q_0) \in q$
4. $\text{state}(s(r), w) = \text{state}(r, w)$

$$q_0 := 0$$

$$q := \text{state}(q_0, q_0) \mid \text{state}(q_0, q_1) \mid \text{state}(q_1, q_0) \mid \text{state}(q_2, q_0)$$

$$q_1 := s(q_0)$$

$$q_2 := s(q_1) \mid s(q_2)$$

System Timbuk [Thomas Genet]. Automated construction, with guess of *acceleration* $q_2 := s(q_2)$ by user assistance.

Concurrent readers/writers: automata construction

1. $\text{state}(0, 0) = \text{state}(0, s(0))$
2. $\text{state}(r, 0) = \text{state}(s(r), 0)$
3. $\text{state}(r, s(w)) = \text{state}(r, w)$
4. $\text{state}(s(r), w) = \text{state}(r, w)$
 $\text{state}(s(q_0 \mid q_1 \mid q_2), q_0) \in q \Rightarrow \text{state}(q_0 \mid q_1 \mid q_2, q_0) \in q$

$$q_0 := 0$$

$$q := \text{state}(q_0, q_0) \mid \text{state}(q_0, q_1) \mid \text{state}(q_1, q_0) \mid \text{state}(q_2, q_0)$$

$$q_1 := s(q_0)$$

$$q_2 := s(q_1) \mid s(q_2)$$

System Timbuk [Thomas Genet]. Automated construction, with guess of *acceleration* $q_2 := s(q_2)$ by user assistance.

Concurrent readers/writers: verification

Properties expected:

1. mutual exclusion between readers and writers
forbidden pattern: $\text{state}(s(x), s(y))$
2. mutual exclusion between writers
forbidden pattern: $\text{state}(x, s(s(y)))$

The **red set**: union of

1. $\text{state}((q_1 \mid q_2), (q_1 \mid q_2))$
2. $\text{state}((q_0 \mid q_1 \mid q_2), (q_1 \mid q_2))$

with $q_0 := 0$, $q_1 := s(q_0)$, $q_2 := s(q_1) \mid s(q_2)$

Verification: The intersection between the set of reachable configurations and the **red set** is empty.

Functional program

Lists built with *constructor* symbols cons and nil.

$$\begin{aligned}\text{app}(\text{nil}, y) &= y \\ \text{app}(\text{cons}(x, y), z) &= \text{cons}(x, \text{app}(y, z))\end{aligned}$$

Functional program analysis

set of initial configurations q_{app} : terms of the form $\text{app}(\ell_1, \ell_2)$
where ℓ_1, ℓ_2 are lists of 0 and 1, defined by

$$q := 0 \mid 1$$

$$q\ell := \text{nil} \mid \text{cons}(q, q\ell)$$

$$q_{\text{app}} := \text{app}(q\ell, q\ell)$$

set of reachable configurations = the closure according to

$$\begin{aligned}\text{app}(\text{nil}, y) &= y \\ \text{app}(\text{cons}(x, y), z) &= \text{cons}(x, \text{app}(y, z))\end{aligned}$$

it is

$$q := 0 \mid 1$$

$$q\ell := \text{nil} \mid \text{cons}(q, q\ell)$$

$$q_{\text{app}} := \text{app}(q\ell, q\ell) \mid \text{cons}(q, q_{\text{app}})$$

Functional program : rev

[Thomas Genet, Valérie Viet Triem Tong, LPAR 01]. Timbuk.

$$\begin{aligned} \text{app}(\text{nil}, y) &= y \\ \text{app}(\text{cons}(x, y), z) &= \text{cons}(x, \text{app}(y, z)) \\ \text{rev}(\text{nil}) &= \text{nil} \\ \text{rev}(\text{cons}(x, y)) &= \text{app}(\text{rev}(y), \text{cons}(x, \text{nil})) \end{aligned}$$

set of initial config.:

$$\begin{aligned} q_0 &:= 0 \\ q_1 &:= 1 \\ q_{l_1} &:= \text{nil} \mid \text{cons}(q_1, q_{l_1}) \\ q_{l_{01}} &:= \text{nil} \mid \text{cons}(q_0, q_{l_1}) \mid \text{cons}(q_0, q_{l_{01}}) \\ q_{\text{rev}} &:= \text{rev}(q_{l_{01}}) \end{aligned}$$

Functional program : rev

[Thomas Genet, Valérie Viet Triem Tong, LPAR 01]. Timbuk.

$$\begin{aligned} \text{app}(\text{nil}, y) &= y \\ \text{app}(\text{cons}(x, y), z) &= \text{cons}(x, \text{app}(y, z)) \\ \text{rev}(\text{nil}) &= \text{nil} \\ \text{rev}(\text{cons}(x, y)) &= \text{app}(\text{rev}(y), \text{cons}(x, \text{nil})) \end{aligned}$$

set of initial config.: $\text{rev}(\ell)$ where $\ell \in q_{\ell_{01}}$, list of 0's followed by 1's

$$\begin{aligned} q_0 &:= 0 \\ q_1 &:= 1 \\ q_{\ell_1} &:= \text{nil} \mid \text{cons}(q_1, q_{\ell_1}) \\ q_{\ell_{01}} &:= \text{nil} \mid \text{cons}(q_0, q_{\ell_1}) \mid \text{cons}(q_0, q_{\ell_{01}}) \\ q_{\text{rev}} &:= \text{rev}(q_{\ell_{01}}) \end{aligned}$$

Functional program cntd

set of reachable configurations: by completion of equations for initial configurations

$$q_0 := 0$$

$$q_1 := 1$$

$$q_{\ell_1} := \text{nil} \mid \text{cons}(q_1, q_{\ell_1}) \mid \text{cons}(q_1, q_{\text{nil}}) \mid \text{app}(q_{\text{nil}}, q_{\ell_1})$$

$$q_{\ell_{01}} := \text{nil} \mid \text{cons}(q_0, q_{\ell_1}) \mid \text{cons}(q_0, q_{\ell_{01}})$$

$$q_{\text{rev}} := \text{rev}(q_{\ell_{01}}) \mid \text{nil} \mid \text{app}(q_{\ell_{10}}, q_{\text{nil}})$$

$$q_{\ell_{10}} := \text{rev}(q_{\ell_{01}}) \mid \text{app}(q_{\ell_1}, q_{\ell_0})$$

$$q_{\text{nil}} := \text{nil} \mid \text{rev}(q_{\text{nil}})$$

$$q_{\ell_0} := \text{cons}(q_0, q_{\text{nil}}) \mid \text{app}(q_{\text{nil}}, q_{\ell_0}) \mid \text{app}(q_{\ell_0}, q_{\ell_0})$$

property expected: $\text{rev}(\ell)$ not reachable when
 $\ell \models \exists x, y \ x < y \wedge 0(x) \wedge 1(y)$.

verification The intersection of q_{rev} and the above set is empty.

Imperative programs

$$p ::= 0 \mid X \mid p \cdot p \mid p \parallel p$$

- ▶ 0: null process (termination)
- ▶ X : program point
- ▶ $p \cdot p$: sequential composition
- ▶ $p \parallel p$: parallel composition

Transition rules

- ▶ procedure call: $X \rightarrow Y \cdot Z$ ($Z =$ return point)
- ▶ procedure call with global state: $Q \cdot X \rightarrow Q' \cdot Y \cdot Z$
- ▶ procedure return: $Q \cdot Y \rightarrow Q'$
- ▶ global state change: $Q \cdot X \rightarrow Q' \cdot X$
- ▶ dynamic thread creation: $X \rightarrow Y \parallel Z$
- ▶ handshake : $X \parallel Y \rightarrow X' \parallel Y'$

Imperative program

[Bouajjani Touili CAV 02]

<code>void X() {</code>	X	\rightarrow	$Y \cdot X$	(r_1)
<code>while(true) {</code>	Y	\rightarrow	t	(r_2)
<code>if Y() {</code>	Y	\rightarrow	f	(r_3)
<code>thread_create(&t1,Z)</code>	$t \cdot X$	\rightarrow	$X \parallel Z$	(r_4)
<code>} else { return }</code>	f	\rightarrow	0	(r_5)
<code>}</code>				
<code>}</code>				
<code>}</code>				

The set of reachable configurations is infinite but **regular**.

Related models of imperative programs

- ▶ Pushdown systems (sequential programs with procedure calls)

$$X_1 \cdot \dots \cdot X_n \rightarrow Y_1 \cdot \dots \cdot Y_m$$

- ▶ Petri nets (multi-threaded programs)

$$X_1 \parallel \dots \parallel X_n \rightarrow Y_1 \parallel \dots \parallel Y_m$$

- ▶ PA processes [Lugiez Schnoebelen TCS 02]

$$X_1 \rightarrow Y_1 \cdot \dots \cdot Y_m, \quad X_1 \rightarrow Y_1 \parallel \dots \parallel Y_m$$

- ▶ Process rewrite systems (PRS) [Mayr Rusinowitch IC 99],
[Bouajjani Touili RTA 05]

$$X_1 \cdot \dots \cdot X_n \rightarrow Y_1 \cdot \dots \cdot Y_m, \quad X_1 \parallel \dots \parallel X_n \rightarrow Y_1 \parallel \dots \parallel Y_m$$

- ▶ Dynamic pushdown networks [Seidl CIAA 09]

Tree languages modulo

In the above model,

- ▶ \cdot is associative,
- ▶ \parallel is associative and commutative.

The terms of the above algebra correspond to **unranked trees**,

- ▶ **ordered** (modulo A) and
- ▶ **unordered** (modulo AC).

(models for XML processing)

Overview

Verification of other infinite-states systems.

- ▶ configuration = tree (ranked or unranked)
 - ▶ process,
 - ▶ message exchanged in a protocol,
 - ▶ local network with a tree shape,
 - ▶ tree data structure in memory, with pointers (e.g. binary search trees)...
- ▶ (infinite) set of configurations = tree language L
- ▶ transition relation between configurations
- ▶ safety: $\text{transitive closure}(L_{\text{init}}) \cap L_{\text{error}} = \emptyset$.

Different kinds of trees

- ▶ finite ranked trees (terms in first order logic)
- ▶ finite unranked ordered trees
- ▶ finite unranked unordered trees
- ▶ infinite trees...

⇒ several classes of tree automata.

Overview: properties of automata

- ▶ determinism,
- ▶ Boolean closures,
- ▶ closures under transformations
(homomorphisms, transducers, rewrite systems...)
- ▶ minimization,
- ▶ decision problems, complexity,
 - ▶ membership,
 - ▶ emptiness,
 - ▶ universality,
 - ▶ inclusion, equivalence,
 - ▶ emptiness of intersection,
 - ▶ finiteness...
- ▶ pumping and star lemma,
- ▶ expressiveness, correspondence with logics.

Organization of the tutorial

1. finite ranked tree automata
 - ▶ properties
 - ▶ algorithms
 - ▶ closure under transformation, applications to program verification
2. correspondence with the monadic second order logic of the tree (Thatcher and Wright's theorem).
3. finite unranked tree automata
 - ▶ ordered = Hedge Automata
 - ▶ unordered = Presburger automata
 - ▶ closure modulo A and AC
 - ▶ XML typing and analysis of transformations
4. tree automata as Horn clause sets

Part I

Automata on Finite Ranked Trees

Terms in first order logic

Plan

Terms

TA: Definitions and Expressiveness

Determinism and Boolean Closures

Decision Problems

Minimization

Closure under Tree Transformations, Program Verification

Signature

Definition : Signature

A signature Σ is a finite set of function symbols each of them with an arity greater or equal to 0.

We denote Σ_i the set of symbols of arity i .

Example :

$\{+ : 2, s : 1, 0 : 0\}, \{\wedge : 2, \vee : 2, \neg : 1, \top, \perp : 0\}$.

We also consider a countable set \mathcal{X} of variable symbols.

Terms

Definition : Term

The set of terms over the signature Σ and \mathcal{X} is the smallest set $\mathcal{T}(\Sigma, \mathcal{X})$ such that:

- $\Sigma_0 \subseteq \mathcal{T}(\Sigma, \mathcal{X})$,
- $\mathcal{X} \subseteq \mathcal{T}(\Sigma, \mathcal{X})$,
- if $f \in \Sigma_n$ and if $t_1, \dots, t_n \in \mathcal{T}(\Sigma, \mathcal{X})$, then $f(t_1, \dots, t_n) \in \mathcal{T}(\Sigma, \mathcal{X})$.

The set of ground terms (terms without variables, i.e. $\mathcal{T}(\Sigma, \emptyset)$) is denoted $\mathcal{T}(\Sigma)$.

Example :

$x, \neg(x), \wedge(\vee(x, \neg(y)), \neg(x))$.

Terms (2)

A term where each variable appears at most once is called **linear**.
A term without variable is called **ground**.

Depth $h(t)$:

- ▶ $h(a) = h(x) = 0$ if $a \in \Sigma_0, x \in \mathcal{X}$,
- ▶ $h(f(t_1, \dots, t_n)) = \max\{h(t_1), \dots, h(t_n)\} + 1$.

Positions

A term $t \in \mathcal{T}(\Sigma, \mathcal{X})$ can also be seen as a function from the set of its **positions** $\mathcal{Pos}(t)$ into $\Sigma \cup \mathcal{X}$.

The empty position (**root**) is denoted ε .

$\mathcal{Pos}(t)$ is a subset of \mathbb{N}^* satisfying the following properties:

- ▶ $\mathcal{Pos}(t)$ is closed under prefix,
- ▶ for all $p \in \mathcal{Pos}(t)$ such that $t(p) \in \Sigma_n$ ($n \geq 1$),
 $\{pj \in \mathcal{Pos}(t) \mid j \in \mathbb{N}\} = \{p1, \dots, pn\}$,
- ▶ every $p \in \mathcal{Pos}(t)$ such that $t(p) \in \Sigma_0 \cup \mathcal{X}$ is maximal in $\mathcal{Pos}(t)$ for the prefix ordering.

The **size** of t is defined by $\|t\| = |\mathcal{Pos}(t)|$.

Subterm $t|_p$ at position $p \in \mathcal{Pos}(t)$:

- ▶ $t|_\varepsilon = t$,
- ▶ $f(t_1, \dots, t_n)|_{ip} = t_i|_p$.

The **replacement** in t of $t|_p$ by s is denoted $t[s]_p$.

Positions (Example)

Example :

$$t = \wedge(\wedge(x, \vee(x, \neg(y))), \neg(x)),$$

$$t|_{11} = x, t|_{12} = \vee(x, \neg(y)), t|_2 = \neg(x),$$

$$t[\neg(y)]_{11} = \wedge(\wedge(\neg(y), \vee(x, \neg(y))), \neg(x)).$$

Contexts

Definition : Contexte

A *context* is a linear term.

The application of a context $C \in \mathcal{T}(\Sigma, \{x_1, \dots, x_n\})$ to n terms t_1, \dots, t_n , denoted $C[t_1, \dots, t_n]$, is obtained by the replacement of each x_i by t_i , for $1 \leq i \leq n$.

Plan

Terms

TA: Definitions and Expressiveness

Determinism and Boolean Closures

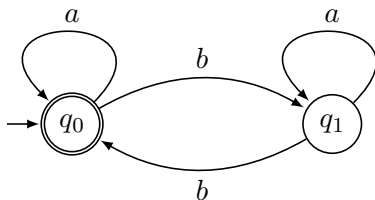
Decision Problems

Minimization

Closure under Tree Transformations, Program Verification

Bottom-up Finite Tree Automata

$(a + b a^* b)^*$



word. run on $aabba$: $q_0 \xrightarrow{a} q_0 \xrightarrow{a} q_0 \xrightarrow{b} q_1 \xrightarrow{b} q_0 \xrightarrow{a} q_0$.

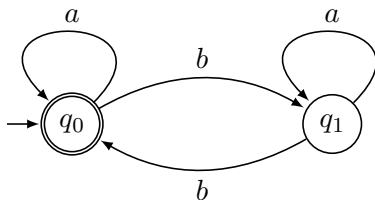
tree. run on $a(a(b(b(a(\varepsilon))))))$:

$q_0 \rightarrow a(q_0) \rightarrow a(a(q_0)) \rightarrow a(a(b(q_1))) \rightarrow a(a(b(b(q_0)))) \rightarrow a(a(b(b(a(q_0)))))) \rightarrow a(a(b(b(a(\varepsilon))))))$

with $q_0 := \varepsilon$, $q_0 := a(q_0)$, $q_1 := a(q_1)$, $q_1 := b(q_0)$, $q_0 := b(q_1)$.

Bottom-up Finite Tree Automata

$(a + b a^* b)^*$



word. run on $aabba$: $q_0 \xrightarrow{a} q_0 \xrightarrow{a} q_0 \xrightarrow{b} q_1 \xrightarrow{b} q_0 \xrightarrow{a} q_0$.

tree. run on $a(a(b(b(a(\varepsilon))))))$:

$a(a(b(b(a(\varepsilon)))))) \rightarrow a(a(b(b(a(q_0)))))) \rightarrow a(a(b(b(q_0)))) \rightarrow$
 $a(a(b(q_1))) \rightarrow a(a(q_0)) \rightarrow a(q_0) \rightarrow q_0$

with $\varepsilon \rightarrow q_0$, $a(q_0) \rightarrow q_0$, $a(q_1) \rightarrow q_1$, $b(q_0) \rightarrow q_1$, $b(q_1) \rightarrow q_0$.

Bottom-up Finite Tree Automata

Definition : Tree Automata

A *tree automaton* (TA) over a signature Σ is a tuple $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$ where Q is a finite set of *states*, $Q^f \subseteq Q$ is the subset of final states and Δ is a set of transition rules of the form: $f(q_1, \dots, q_n) \rightarrow q$ with $f \in \Sigma_n$ ($n \geq 0$) and $q_1, \dots, q_n, q \in Q$.

The state q is called the head of the rule.

The **language** of \mathcal{A} in state q is recursively defined by

$$L(\mathcal{A}, q) = \{a \in \Sigma_0 \mid a \rightarrow q \in \Delta\} \\ \cup \bigcup_{f(q_1, \dots, q_n) \rightarrow q \in \Delta} f(L(\mathcal{A}, q_1), \dots, L(\mathcal{A}, q_n))$$

with $f(L_1, \dots, L_n) := \{f(t_1, \dots, t_n) \mid t_1 \in L_1, \dots, t_n \in L_n\}$.

We say that $t \in L(\mathcal{A}, q)$ is **accepted**, or **recognized**, by \mathcal{A} in state q .

The **language** of \mathcal{A} is $L(\mathcal{A}) := \bigcup_{q^f \in Q^f} L(\mathcal{A}, q^f)$ (**regular** language).

Recognized Languages: Operational Definition

Rewrite Relation

The rewrite relation associated to Δ is the smallest binary relation, denoted $\xrightarrow{\Delta}$, containing Δ and closed under application of contexts.

The reflexive and transitive closure of $\xrightarrow{\Delta}$ is denoted $\xrightarrow{\Delta}^*$.

For $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$, it holds that

$$L(\mathcal{A}, q) = \{t \in \mathcal{T}(\Sigma) \mid t \xrightarrow{\Delta}^* q\}$$

and hence

$$L(\mathcal{A}) = \{t \in \mathcal{T}(\Sigma) \mid t \xrightarrow{\Delta}^* q \in Q^f\}$$

Tree Automata: Example 1

Example :

$$\Sigma = \{\wedge : 2, \vee : 2, \neg : 1, \top, \perp : 0\},$$

$$\mathcal{A} = \left(\Sigma, \{q_0, q_1\}, \{q_1\}, \left\{ \begin{array}{ll} \perp \rightarrow q_0 & \top \rightarrow q_1 \\ \neg(q_0) \rightarrow q_1 & \neg(q_1) \rightarrow q_0 \\ \vee(q_0, q_0) \rightarrow q_0 & \vee(q_0, q_1) \rightarrow q_1 \\ \vee(q_1, q_0) \rightarrow q_1 & \vee(q_1, q_1) \rightarrow q_1 \\ \wedge(q_0, q_0) \rightarrow q_0 & \wedge(q_0, q_1) \rightarrow q_0 \\ \wedge(q_1, q_0) \rightarrow q_0 & \wedge(q_1, q_1) \rightarrow q_1 \end{array} \right\} \right)$$

$$\begin{aligned} & \wedge(\wedge(\top, \vee(\top, \neg(\perp))), \neg(\top)) \xrightarrow{\mathcal{A}} \wedge(\wedge(\top, \vee(\top, \neg(\perp))), \neg(q_1)) \\ \xrightarrow{\mathcal{A}} & \wedge(\wedge(q_1, \vee(q_1, \neg(q_0))), \neg(q_1)) \xrightarrow{\mathcal{A}} \wedge(\wedge(q_1, \vee(q_1, \neg(q_0))), q_0) \\ \xrightarrow{\mathcal{A}} & \wedge(\wedge(q_1, \vee(q_1, q_1)), q_0) \xrightarrow{\mathcal{A}} \wedge(\wedge(q_1, q_1), q_0) \xrightarrow{\mathcal{A}} \wedge(q_1, q_0) \xrightarrow{\mathcal{A}} q_0 \end{aligned}$$

Tree Automata: Example 2

Example :

$\Sigma = \{\wedge : 2, \vee : 2, \neg : 1, \top, \perp : 0\}$,

TA recognizing the ground instances of $\neg(\neg(x))$:

$$\mathcal{A} = \left(\Sigma, \{q, q_{\neg}, q_f\}, \{q_f\}, \left\{ \begin{array}{ll} \perp \rightarrow q & \top \rightarrow q \\ \neg(q) \rightarrow q & \neg(q) \rightarrow q_{\neg} \\ \neg(q_{\neg}) \rightarrow q_f & \\ \vee(q, q) \rightarrow q & \wedge(q, q) \rightarrow q \end{array} \right\} \right)$$

Example :

Ground terms embedding the pattern $\neg(\neg(x))$: $\mathcal{A} \cup \{\neg(q_f) \rightarrow q_f, \vee(q_f, q_*) \rightarrow q_f, \vee(q_*, q_f) \rightarrow q_f, \dots\}$ (propagation of q_f).

Linear Pattern Matching

Proposition :

Given a linear term $t \in \mathcal{T}(\Sigma, \mathcal{X})$, there exists a TA \mathcal{A} recognizing the set of ground instances of t : $L(\mathcal{A}) = \{t\sigma \mid \sigma : \mathcal{X} \rightarrow \mathcal{T}(\Sigma)\}$.

e.g. in regular tree model checking, definition of error configurations by **forbidden patterns**.

Runs

Definition : Run

A *run* of a TA (Σ, Q, Q^f, Δ) on a term $t \in \mathcal{T}(\Sigma)$ is a function $r : \mathcal{Pos}(t) \rightarrow Q$ such that for all $p \in \mathcal{Pos}(t)$,
if $t(p) = f \in \Sigma_n$, $r(p) = q$ and $r(pi) = q_i$ for all $1 \leq i \leq n$,
then $f(q_1, \dots, q_n) \rightarrow q \in \Delta$.

The run r is *accepting* if $r(\varepsilon) \in Q^f$.

$L(\mathcal{A})$ is the set of ground terms of $\mathcal{T}(\Sigma)$ for which there exists an accepting run.

Pumping Lemma

Lemma : Pumping Lemma

Let $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$.

$L(\mathcal{A}) \neq \emptyset$ iff there exists $t \in L(\mathcal{A})$ such that $h(t) \leq |Q|$.

Lemma : Iteration Lemma

For all TA \mathcal{A} , there exists $k > 0$ such that for all term $t \in L(\mathcal{A})$ with $h(t) > k$, there exists 2 contexts $C, D \in \mathcal{T}(\Sigma, \{x_1\})$ with $D \neq x_1$ and a term $u \in \mathcal{T}(\Sigma)$ such that $t = C[D[u]]$ and for all $n \geq 0$, $C[D^n[u]] \in L(\mathcal{A})$.

usage: to show that a language is not regular.

Non Regular Languages

We show with the pumping and iteration lemmatas that the following tree languages are not regular:

- ▶ $\{f(t, t) \mid t \in \mathcal{T}(\Sigma)\}$,
- ▶ $\{f(g^n(a), h^n(a)) \mid n \geq 0\}$,
- ▶ $\{t \in \mathcal{T}(\Sigma) \mid |\mathcal{P}os(t)| \text{ is prime}\}$.

Epsilon-Transitions

We extend the class TA into TA_ϵ with the addition of another type of transition rules of the form $q \xrightarrow{\epsilon} q'$ (ϵ -transition).
with the same expressiveness as TA.

Proposition : Suppression of ϵ -transitions

For all $TA_\epsilon \mathcal{A}_\epsilon$, there exists a TA (without ϵ -transition) \mathcal{A}' such that $L(\mathcal{A}) = L(\mathcal{A}_\epsilon)$. The size of \mathcal{A} is polynomial in the size of \mathcal{A}_ϵ .

pr.: We start with \mathcal{A}_ϵ and we add $f(q_1, \dots, q_n) \rightarrow q'$ if there exists $f(q_1, \dots, q_n) \rightarrow q$ and $q \xrightarrow{\epsilon} q'$.

Top-Down Tree Automata

Definition : Top-Down Tree Automata

A top-down tree automaton over a signature Σ is a tuple $\mathcal{A} = (\Sigma, Q, Q^{\text{init}}, \Delta)$ where Q is a finite set of *states*, $Q^{\text{init}} \subseteq Q$ is the subset of initial states and Δ is a set of transition rules of the form: $q \rightarrow f(q_1, \dots, q_n)$ with $f \in \Sigma_n$ ($n \geq 0$) and $q_1, \dots, q_n, q \in Q$.

A ground term $t \in \mathcal{T}(\Sigma)$ is accepted by \mathcal{A} in the state q iff $q \xrightarrow{\Delta}^* t$.

The language of \mathcal{A} starting from the state q is $L(\mathcal{A}, q) := \{t \in \mathcal{T}(\Sigma) \mid q \xrightarrow{\Delta}^* t\}$.

The language of \mathcal{A} is $L(\mathcal{A}) := \bigcup_{q^i \in Q^{\text{init}}} L(Q, q^i)$.

Top-Down Tree Automata (Expressiveness)

Proposition : Expressiveness

The set of top-down tree automata languages is exactly the set of regular tree languages.

Remark: Notations

In the next slides

TA = Bottom-Up Tree Automata

Plan

Terms

TA: Definitions and Expressiveness

Determinism and Boolean Closures

Decision Problems

Minimization

Closure under Tree Transformations, Program Verification

Determinism

Definition : Determinism

A TA \mathcal{A} is *deterministic* if for all $f \in \Sigma_n$, for all states q_1, \dots, q_n of \mathcal{A} , there is at most one state q of \mathcal{A} such that \mathcal{A} contains a transition $f(q_1, \dots, q_n) \rightarrow q$.

If \mathcal{A} is deterministic, then for all $t \in \mathcal{T}(\Sigma)$, there exists at most one state q of \mathcal{A} such that $t \in L(\mathcal{A}, q)$. It is denoted $\mathcal{A}(t)$ or $\Delta(t)$.

Completeness

Definition : Completeness

A TA \mathcal{A} is *complete* if for all $f \in \Sigma_n$, for all states q_1, \dots, q_n of \mathcal{A} , there is at least one state q of \mathcal{A} such that \mathcal{A} contains a transition $f(q_1, \dots, q_n) \rightarrow q$.

If \mathcal{A} is complete, then for all $t \in \mathcal{T}(\Sigma)$, there exists at least one state q of \mathcal{A} such that $t \in L(\mathcal{A}, q)$.

Completion

Proposition : Completion

For all TA \mathcal{A} , there exists a complete TA \mathcal{A}_c such that $L(\mathcal{A}_c) = L(\mathcal{A})$. Moreover, if \mathcal{A} is deterministic, then \mathcal{A}_c is deterministic. The size of \mathcal{A}_c is polynomial in the size of \mathcal{A} , its construction is PTIME.

Completion

Proposition : Completion

For all TA \mathcal{A} , there exists a complete TA \mathcal{A}_c such that $L(\mathcal{A}_c) = L(\mathcal{A})$. Moreover, if \mathcal{A} is deterministic, then \mathcal{A}_c is deterministic. The size of \mathcal{A}_c is polynomial in the size of \mathcal{A} , its construction is PTIME.

pr.: add a trash state q_{\perp} .

Determinization

Proposition : Determinization

For all TA \mathcal{A} , there exists a deterministic TA \mathcal{A}_{det} such that $L(\mathcal{A}_{det}) = L(\mathcal{A})$. Moreover, if \mathcal{A} is complete, then \mathcal{A}_{det} is complete. The size of \mathcal{A}_{det} is exponential in the size of \mathcal{A} , its construction is EXPTIME.

pr.: subset construction. Transitions:

$$f(S_1, \dots, S_n) \rightarrow \{q \mid \exists q_1 \in S_1 \dots \exists q_n \in S_n f(q_1, \dots, q_n) \rightarrow q \in \Delta\}$$

for all $S_1, \dots, S_n \subseteq Q$.

Determinization (Example)

Exercise :

Determinise and complete the previous TA (pattern matching of $\neg(\neg(x))$):

$$\mathcal{A} = \left(\Sigma, \{q, q_{\neg}, q_f\}, \{q_f\}, \left\{ \begin{array}{ll} \perp \rightarrow q & \top \rightarrow q \\ \neg(q) \rightarrow q & \neg(q) \rightarrow q_{\neg} \\ \neg(q_{\neg}) \rightarrow q_f & \neg(q_f) \rightarrow q_f \\ \vee(q, q) \rightarrow q & \wedge(q, q) \rightarrow q \\ \vee(q_f, q_*) \rightarrow q_f & \vee(q_*, q_f) \rightarrow q_f \end{array} \right\} \right)$$

Top-Down Tree Automata and Determinism

Definition : Determinism

A top-down tree automaton $(\Sigma, Q, Q^{\text{init}}, \Delta)$ is *deterministic* if $|Q^{\text{init}}| = 1$ and for all state $q \in Q$ and $f \in \Sigma$, Δ contains at most one rule with left member q and symbol f .

The top-down tree automata are in general not determinizable .

Proposition :

There exists a regular tree language which is not recognizable by a deterministic top-down tree automaton.

Top-Down Tree Automata and Determinism

Definition : Determinism

A top-down tree automaton $(\Sigma, Q, Q^{\text{init}}, \Delta)$ is *deterministic* if $|Q^{\text{init}}| = 1$ and for all state $q \in Q$ and $f \in \Sigma$, Δ contains at most one rule with left member q and symbol f .

The top-down tree automata are in general not determinizable .

Proposition :

There exists a regular tree language which is not recognizable by a deterministic top-down tree automaton.

pr.: $L = \{f(a, b), f(b, a)\}$.

Boolean Closure of Regular Tree Languages

Proposition : Closure

The class of regular tree languages is closed under union, intersection and complementation.

op.	technique	computation time and size of automata
\cup	disjoint \cup	
\cap	Cartesian product	
\neg	determinization, completion, invert final / non-final states	(lower bound)

Remark :

For the deterministic TA, the construction for the complementation is polynomial.

Boolean Closure of Regular Tree Languages

Proposition : Closure

The class of regular tree languages is closed under union, intersection and complementation.

op.	technique	computation time and size of automata
\cup	disjoint \cup	linear
\cap	Cartesian product	
\neg	determinization, completion, invert final / non-final states	(lower bound)

Remark :

For the deterministic TA, the construction for the complementation is polynomial.

Boolean Closure of Regular Tree Languages

Proposition : Closure

The class of regular tree languages is closed under union, intersection and complementation.

op.	technique	computation time and size of automata
\cup	disjoint \cup	linear
\cap	Cartesian product	quadratic
\neg	determinization, completion, invert final / non-final states	(lower bound)

Remark :

For the deterministic TA, the construction for the complementation is polynomial.

Boolean Closure of Regular Tree Languages

Proposition : Closure

The class of regular tree languages is closed under union, intersection and complementation.

op.	technique	computation time and size of automata
\cup	disjoint \cup	linear
\cap	Cartesian product	quadratic
\neg	determinization, completion, invert final / non-final states	exponential (lower bound)

Remark :

For the deterministic TA, the construction for the complementation is polynomial.

Plan

Terms

TA: Definitions and Expressiveness

Determinism and Boolean Closures

Decision Problems

Minimization

Closure under Tree Transformations, Program Verification

Cleaning

Definition : Clean

A state q of a TA \mathcal{A} is called *inhabited* if there exists at least one $t \in L(\mathcal{A}, q)$. A TA is called *clean* if all its states are inhabited.

Proposition : Cleaning

For all TA \mathcal{A} , there exists a clean TA \mathcal{A}_{clean} such that $L(\mathcal{A}_{clean}) = L(\mathcal{A})$. The size of \mathcal{A}_{clean} is smaller than the size of \mathcal{A} , its construction is PTIME.

pr.: state marking algorithm, running time $O(|Q| \times \|\Delta\|)$.

State Marking Algorithm

We construct $M \subseteq Q$ containing all the inhabited states.

- ▶ start with $M = \emptyset$
- ▶ for all $f \in \Sigma$, of arity $n \geq 0$, and all $q_1, \dots, q_n \in M$ st there exists $f(q_1, \dots, q_n) \rightarrow q$ in Δ , add q to M (if it was not already).

We iterate the last step until a fixpoint M_* is reached.

Lemma :

$q \in M_*$ iff $\exists t \in L(\mathcal{A}, q)$.

Membership Problem

Definition : Membership

INPUT: a TA \mathcal{A} over Σ , a term $t \in \mathcal{T}(\Sigma)$.
QUESTION: $t \in L(\mathcal{A})$?

Proposition : Membership

The membership problem is decidable in polynomial time.

Exact complexity:

- ▶ non-deterministic bottom-up: LOGCFL-complete
- ▶ deterministic bottom-up: unknown (LOGDCFL)
- ▶ deterministic top-down: LOGSPACE-complete.

Emptiness Problem

Definition : Emptiness

INPUT: a TA \mathcal{A} over Σ .

QUESTION: $L(\mathcal{A}) = \emptyset$?

Proposition : Emptiness

The emptiness problem is decidable in linear time.

Emptiness Problem

Definition : Emptiness

INPUT: a TA \mathcal{A} over Σ .

QUESTION: $L(\mathcal{A}) = \emptyset?$

Proposition : Emptiness

The emptiness problem is decidable in linear time.

pr.:

quadratic: clean, check if the clean automaton contains a final state.

linear: reduction to propositional HORN-SAT.

linear bis: optimization of the data structures for the cleaning (exo).

Remark :

The problem of the emptiness is PTIME-complete.

Instance-Membership Problem

Definition : Instance-Membership (IM)

INPUT: a TA \mathcal{A} over Σ , a term $t \in \mathcal{T}(\Sigma, \mathcal{X})$.

QUESTION: does there exists $\sigma : \text{vars}(t) \rightarrow \mathcal{T}(\Sigma)$ s.t. $t\sigma \in L(\mathcal{A})$?

Proposition : Instance-Membership

1. The problem IM is decidable in polynomial time when t is linear.
2. The problem IM is NP-complet when \mathcal{A} is deterministic.
3. The problem IM is EXPTIME-complete in general.

Problem of the Emptiness of Intersection

Definition : Emptiness of Intersection

INPUT: n TA $\mathcal{A}_1, \dots, \mathcal{A}_n$ over Σ .

QUESTION: $L(\mathcal{A}_1) \cap \dots \cap L(\mathcal{A}_n) = \emptyset?$

Proposition : Emptiness of Intersection

The problem of the emptiness of intersection is EXPTIME-complete.

Problem of the Emptiness of Intersection

Definition : Emptiness of Intersection

INPUT: n TA $\mathcal{A}_1, \dots, \mathcal{A}_n$ over Σ .

QUESTION: $L(\mathcal{A}_1) \cap \dots \cap L(\mathcal{A}_n) = \emptyset?$

Proposition : Emptiness of Intersection

The problem of the emptiness of intersection is EXPTIME-complete.

pr.: EXPTIME: n applications of the closure under \cap and emptiness decision.

EXPTIME-hardness: APSPACE = EXPTIME

reduction of the problem of the existence of a successful run (starting from an initial configuration) of an alternating Turing machine (ATM) $M = (\Gamma, S, s_0, S_f, \delta)$.

[Seidl 94], [Veanes 97]

Let $M = (\Gamma, S, s_0, S_f, \delta)$ be a Turing Machine (Γ : input alphabet, S : state set, s_0 initial state, S_f final states, δ : transition relation).
First some notations.

- ▶ a *configuration* of M is a word of $\Gamma^* \Gamma_S \Gamma^*$ where $\Gamma_S = \{a^s \mid a \in \Gamma, s \in S\}$. In this word, the letter of Γ_S indicates both the current state and the current position of the head of M .
- ▶ a *final configuration* of M is a word of $\Gamma^* \Gamma_{S_f} \Gamma^*$.
- ▶ an *initial configuration* of M is a word of $\Gamma_{s_0} \Gamma^*$.
- ▶ a *transition* of M (following δ) between two configurations v and v' is denoted $v \triangleright v'$

The initial configuration v_0 is accepting iff there exists a final configuration v_f and a finite sequence of transitions $v_0 \triangleright \dots \triangleright v_f$?
This problem whether v_0 is accepting is undecidable in general.
If the tape is polynomially bounded (we are restricted to configurations of length $n = |v_0|^c$, for some fixed $c \in \mathbb{N}$), the problem is PSPACE complete.

M alternating: $S = S_{\exists} \uplus S_{\forall}$.

Definition accepting configurations:

- ▶ every final configuration (whose state is in S_f) is accepting
- ▶ a configuration c whose state is in S_{\exists} is accepting if it has at least one successor accepting
- ▶ a configuration c whose state is in S_{\forall} is accepting if all its successors are accepting

Theorem (Chandra, Kozen, Stockmeyer 81)

APSPACE = EXPTIME

In order to show EXPTIME-hardness, we reduce the problem of deciding whether v_0 is accepting for M alternating and polynomially bounded.

Hypotheses (non restrictive):

- ▶ $s_0 \in S_{\exists}$ or $s_0 \in S_{\forall} \cap S_f$
- ▶ s_0 is non reentering (it only occurs in v_0)
- ▶ every configuration with state in S_{\forall} has 0 or 2 successors
- ▶ final configurations are restricted to $b_{S_f}b^*$ where $b \in \Gamma$ is the blank symbol.
- ▶ S_f is a singleton.

2 technical definitions: for $k \leq n$,

$$\begin{aligned} \text{view}(v, k) = & \begin{array}{ll} v[k]v[k+1] & \text{if } k = 1 \\ v[k-1]v[k] & \text{if } k = n \\ v[k-1]v[k]v[k+1] & \text{otherwise} \end{array} \end{aligned}$$

$$\text{view}(v, v_1, v_2, k) = \langle \text{view}(v, k), \text{view}(v_1, k), \text{view}(v_2, k) \rangle$$

$v \triangleright_k \langle v_1, v_2 \rangle$ iff

1. if $v[k] \in \Gamma_S$, then $\exists w \triangleright w_1, w_2$ s.t.
 $\text{view}(v, v_1, v_2, k) = \text{view}(w, w_1, w_2, k)$
2. if $v[k] = a \in \Gamma$, then $v_1[k] \in \{a\} \cup a_S$ and $v_2 = \varepsilon$ or
 $v_2[k] \in \{a\} \cup a_S$.

first item: around position k , we have two correct transitions of M . This can be tested by the membership of $\text{view}(v, v_1, v_2, k)$ to a given set which only depends on M .

Lemma

$v \triangleright v_1, v_2$ iff $\forall k \leq n \ v \triangleright_k \langle v_1, v_2 \rangle$.

Term representations of runs:

rem. a run of M is not a sequence of configurations but a tree of configurations (because of alternation).

Signature Σ : \emptyset : constant, Γ : unary, S : unaires, p binary.

Notation: if $v = a_1 \dots a_n$, $v(x)$ denotes $a_n(a_{n-1}(\dots a_1(x)))$.

Term representations of runs:

- ▶ $v_f(p(\emptyset, \emptyset))$ with v_f final configuration,
- ▶ $v(p(t_1, t_2))$ with v \forall -configuration, $t_1 = v'_1(p(t_{1,1}, t_{1,2}))$, $t_2 = v'_2(p(t_{2,1}, t_{2,2}))$ are two term representations of runs, and $v_1 \triangleright v'_1$, $v_2 \triangleright v'_2$
- ▶ $v(p(t_1, \emptyset))$ with v \exists -configuration, $t_1 = v'_1(p(t_{1,1}, t_{1,2}))$ term representations of run, and $v_1 \triangleright v'_1$.

notations for $t_1 = v'_1(p(t_{1,1}, t_{1,2}))$:

- ▶ $\text{head}(t_1) = v_1$
- ▶ $\text{left}(t_1) = t_{1,1}$
- ▶ $\text{right}(t_1) = t_{1,2}$.

This recursive definition suggest the construction of a TA recognizing term representations of successful runs. The difficulty

is the conditions $v_1 \triangleright v'_1$, $v_2 \triangleright v'_2$, for which we use the above lemma.

We build $2n$ deterministic automata :

for all $1 < k < n$, \mathcal{A}_k recognizes

- ▶ $v_f(p(\emptyset, \emptyset))$ (recall there is only 1 final configuration by hyp.)
- ▶ $v(p(t_1, t_2))$ such that $t_1 \neq \emptyset$ and
 - ▶ $v \triangleright_k \langle \text{head}(t_1), \text{head}(t_2) \rangle$
 - ▶ $\text{left}(t_1) \in L(\mathcal{A}_k)$, $\text{right}(t_1) \in L(\mathcal{A}_k) \cup \{\emptyset\}$,
 - ▶ $t_2 = \emptyset$ or $\text{left}(t_2) \in L(\mathcal{A}_k)$, $\text{right}(t_2) \in L(\mathcal{A}_k) \cup \{\emptyset\}$

idea: \mathcal{A}_k memorizes $\text{view}(\text{head}(t_1), k)$ and $\text{view}(\text{head}(t_2), k)$ and compare with $\text{view}(v, k)$.

for all $1 < k < n$, \mathcal{A}'_k recognizes the terms $v_0(p(t_1, t_2))$ with $t_1 = t_2 = \emptyset$ (if s_0 universal and final) or $t_2 = \emptyset$ (if s_0 existential, not final) and $t_1, t_2 \in T$, minimal set of terms without s_0 containing

- ▶ \emptyset
- ▶ $v(p(t_1, t_2))$ such that $t_1 \neq \emptyset$ and
 - ▶ $v \triangleright_k \langle \text{head}(t_1), \text{head}(t_2) \rangle$
 - ▶ $\text{left}(t_1) \in T$, $\text{right}(t_1) \in T$,

- ▶ $t_2 = \emptyset$ or $\text{left}(t_2) \in T$, $\text{right}(t_2) \in T$

representations of successful runs = $\bigcap_{k=1}^n L(\mathcal{A}_k) \cap L(\mathcal{A}'_k)$.

Problem of Universality

Definition : Universality

INPUT: a TA \mathcal{A} over Σ .

QUESTION: $L(\mathcal{A}) = \mathcal{T}(\Sigma)$

Proposition : Universality

The problem of universality is EXPTIME-complete.

Problem of Universality

Definition : Universality

INPUT: a TA \mathcal{A} over Σ .

QUESTION: $L(\mathcal{A}) = \mathcal{T}(\Sigma)$

Proposition : Universality

The problem of universality is EXPTIME-complete.

pr.: EXPTIME: Boolean closure and emptiness decision.

EXPTIME-hardness: again APSPACE = EXPTIME.

Remark :

The problem of universality is decidable in polynomial time for the deterministic (bottom-up) TA.

pr.: completion and cleaning.

Problems of Inclusion and Equivalence

Definition : Inclusion

INPUT: two TA \mathcal{A}_1 and \mathcal{A}_2 over Σ .

QUESTION: $L(\mathcal{A}_1) \subseteq L(\mathcal{A}_2)$

Definition : Equivalence

INPUT: two TA \mathcal{A}_1 and \mathcal{A}_2 over Σ .

QUESTION: $L(\mathcal{A}_1) = L(\mathcal{A}_2)$

Proposition : Inclusion, Equivalence

The problems of inclusion and equivalence are EXPTIME-complete.

Problems of Inclusion and Equivalence

Definition : Inclusion

INPUT: two TA \mathcal{A}_1 and \mathcal{A}_2 over Σ .
QUESTION: $L(\mathcal{A}_1) \subseteq L(\mathcal{A}_2)$

Definition : Equivalence

INPUT: two TA \mathcal{A}_1 and \mathcal{A}_2 over Σ .
QUESTION: $L(\mathcal{A}_1) = L(\mathcal{A}_2)$

Proposition : Inclusion, Equivalence

The problems of inclusion and equivalence are EXPTIME-complete.

pr.: $L(\mathcal{A}_1) \subseteq L(\mathcal{A}_2)$ iff $L(\mathcal{A}_1) \cap \overline{L(\mathcal{A}_2)} = \emptyset$.

Problems of Inclusion and Equivalence

Definition : Inclusion

INPUT: two TA \mathcal{A}_1 and \mathcal{A}_2 over Σ .
QUESTION: $L(\mathcal{A}_1) \subseteq L(\mathcal{A}_2)$

Definition : Equivalence

INPUT: two TA \mathcal{A}_1 and \mathcal{A}_2 over Σ .
QUESTION: $L(\mathcal{A}_1) = L(\mathcal{A}_2)$

Proposition : Inclusion, Equivalence

The problems of inclusion and equivalence are EXPTIME-complete.

pr.: $L(\mathcal{A}_1) \subseteq L(\mathcal{A}_2)$ iff $L(\mathcal{A}_1) \cap \overline{L(\mathcal{A}_2)} = \emptyset$.

EXPTIME-hardness: universality is $\mathcal{T}(\Sigma) = L(\mathcal{A}_2)$?

Remark :

If \mathcal{A}_1 and \mathcal{A}_2 are deterministic, it is $O(\|\mathcal{A}_1\| \times \|\mathcal{A}_2\|)$.

Problem of Finiteness

Definition : Finiteness

INPUT: a TA \mathcal{A}

QUESTION: is $L(\mathcal{A})$ finite?

Proposition : Finiteness

The problem of finiteness is decidable in polynomial time.

Plan

Terms

TA: Definitions and Expressiveness

Determinism and Boolean Closures

Decision Problems

Minimization

Closure under Tree Transformations, Program Verification

Theorem of Myhill-Nerode

Definition :

A *congruence* \equiv on $\mathcal{T}(\Sigma)$ is an equivalence relation such that for all $f \in \Sigma_n$, if $s_1 \equiv t_1, \dots, s_n \equiv t_n$, then $f(s_1, \dots, s_n) \equiv f(t_1, \dots, t_n)$.

Given $L \subseteq \mathcal{T}(\Sigma)$, the congruence \equiv_L is defined by:

$s \equiv_L t$ if for all context $C \in \mathcal{T}(\Sigma, \{x\})$, $C[s] \in L$ iff $C[t] \in L$.

Theorem : Myhill-Nerode

The three following propositions are equivalent:

1. L is regular
2. L is a union of equivalence classes for a congruence \equiv of finite index
3. \equiv_L is a congruence of finite index

Proof Theorem of Myhill-Nerode

1 \Rightarrow 2. \mathcal{A} deterministic, def. $s \equiv_{\mathcal{A}} t$ iff $\mathcal{A}(s) = \mathcal{A}(t)$.

2 \Rightarrow 3. we show that if $s \equiv t$ then $s \equiv_L t$, hence the index of $\equiv_L \leq$ index of \equiv (since we have $\equiv \subseteq \equiv_L$).

If $s \equiv t$ then $C[s] \equiv C[t]$ for all $C[]$ (induction on C), hence $C[s] \in L$ iff $C[t] \in L$, i.e. $s \equiv_L t$.

3 \Rightarrow 1. we construct $\mathcal{A}_{\min} = (Q_{\min}, Q_{\min}^f, \Delta_{\min})$,

- ▶ $Q_{\min} =$ equivalence classes of \equiv_L ,
- ▶ $Q_{\min}^f = \{[s] \mid s \in L\}$,
- ▶ $\Delta_{\min} = \{f([s_1], \dots, [s_n]) \rightarrow [f(s_1, \dots, s_n)]\}$

Clearly, \mathcal{A}_{\min} is deterministic, and for all $s \in \mathcal{T}(\Sigma)$, $\mathcal{A}_{\min}(s) = [s]_L$, i.e. $s \in L(\mathcal{A}_{\min})$ iff $s \in L$.

Minimization

Corollary :

For all DTA $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$, there exists a unique DTA \mathcal{A}_{\min} whose number of states is the index of $\equiv_{L(\mathcal{A})}$ and such that $L(\mathcal{A}_{\min}) = L(\mathcal{A})$.

Minimization

Let $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$ be a DTA, we build a deterministic minimal automaton \mathcal{A}_{\min} as in the proof of $3 \Rightarrow 1$ of the previous theorem for $L(\mathcal{A})$ (i.e. Q_{\min} is the set of equivalence classes for $\equiv_{L(\mathcal{A})}$).

We build first an equivalence \approx on the states of Q :

▶ $q \approx_0 q'$ iff $q, q' \in Q^f$ ou $q, q' \in Q \setminus Q^f$.

▶ $q \approx_{k+1} q'$ iff $q \approx_k q'$ et $\forall f \in \Sigma_n$,
 $\forall q_1, \dots, q_{i-1}, q_{i+1}, \dots, q_n \in Q$ ($1 \leq i \leq n$),

$$\Delta(f(q_1, \dots, q_{i-1}, q, q_{i+1}, \dots, q_n)) \approx_k \Delta(f(q_1, \dots, q_{i-1}, q', q_{i+1}, \dots, q_n))$$

Let \approx be the fixpoint of this construction, \approx is $\equiv_{L(\mathcal{A})}$, hence

$\mathcal{A}_{\min} = (\Sigma, Q_{\min}, Q_{\min}^f, \Delta_{\min})$ with :

▶ $Q_{\min} = \{[q]_{\approx} \mid q \in Q\}$,

▶ $Q_{\min}^f = \{[q^f]_{\approx} \mid q^f \in Q^f\}$,

▶ $\Delta_{\min} = \{f([q_1]_{\approx}, \dots, [q_n]_{\approx}) \rightarrow [f(q_1, \dots, q_n)]_{\approx}\}$.

recognizes $L(\mathcal{A})$. and it is smaller than \mathcal{A} .

Algebraic Characterization of Regular Languages

Corollary :

A set $L \subseteq \mathcal{T}(\Sigma)$ is regular iff there exists

- ▶ a Σ -algebra \mathcal{Q} of finite domain Q ,
- ▶ an homomorphism $h : \mathcal{T}(\Sigma) \rightarrow \mathcal{A}$,
- ▶ a subset $Q^f \subseteq Q$ such that $L = h^{-1}(Q^f)$.

operations of \mathcal{Q} :

for each $f \in \Sigma_n$, there is a function $f^{\mathcal{Q}} : Q^n \rightarrow Q$.

Plan

Terms

TA: Definitions and Expressiveness

Determinism and Boolean Closures

Decision Problems

Minimization

Closure under Tree Transformations, Program Verification

- Tree Homomorphisms

- Tree Transducers

- Term Rewriting

- Tree Automata Based Program Verification

Tree Transformations, Verification

- ▶ formalisms for the **transformation** of terms (languages):
rewrite systems, tree homomorphisms, transducers...
 - = transitions in an infinite states system,
 - = evaluation of programs,
 - = transformation of XML documents, updates...
- ▶ problem of the **type checking**:
given:
 - ▶ $L_{\text{in}} \subseteq \mathcal{T}(\Sigma)$, (regular) input language
 - ▶ h transformation $\mathcal{T}(\Sigma) \rightarrow \mathcal{T}(\Sigma')$
 - ▶ $L_{\text{out}} \subseteq \mathcal{T}(\Sigma')$ (regular) output languagequestion: do we have $h(L_{\text{in}}) \subseteq L_{\text{out}}$?

Tree Homomorphisms

Tree Homomorphisms

Definition :

$$h : \mathcal{T}(\Sigma) \rightarrow \mathcal{T}(\Sigma')$$

$$h(f(t_1, \dots, t_n)) := t_f \{x_1 \leftarrow h(t_1), \dots, x_n \leftarrow h(t_n)\}$$

for $f \in \Sigma_n$, with $t_f \in \mathcal{T}(\Sigma', \{x_1, \dots, x_n\})$.

h is called

- ▶ *linear* if for all $f \in \Sigma$, t_f is linear,
- ▶ *complete* if for all $f \in \Sigma_n$, $\text{vars}(t_f) = \{x_1, \dots, x_n\}$,
- ▶ *symbol-to-symbol* if for all $f \in \Sigma_n$, $\text{height}(t_f) = 1$.

Homomorphisms: Examples

Example : ternary trees \rightarrow binary trees

Let $\Sigma = \{a : 0, b : 0, g : 3\}$, $\Sigma' = \{a : 0, b : 0, f : 2\}$ and $h : \mathcal{T}(\Sigma) \rightarrow \mathcal{T}(\Sigma')$ defined by

- ▶ $t_a = a$,
- ▶ $t_b = b$,
- ▶ $t_g = f(x_1, f(x_2, x_3))$.

$$h(g(a, g(b, b, b), a)) = f(a, f(f(f(b, f(b, b))), a))$$

Example : Elimination of the \wedge

Let $\Sigma = \{0 : 0, 1 : 0, \neg : 1, \vee : 2, \wedge : 2\}$, $\Sigma' = \{0 : 0, 1 : 0, \neg : 1, \vee : 2\}$ and $h : \mathcal{T}(\Sigma) \rightarrow \mathcal{T}(\Sigma')$ with $t_\wedge = \neg(\vee(\neg(x_1), \neg(x_2)))$.

Closure of Regular Languages under Linear Homomorphisms

Theorem :

If L is regular and h is a linear homomorphism, then $h(L)$ is regular.

Closure of Regular Languages under Linear Homomorphisms

Theorem :

If L is regular and h is a linear homomorphism, then $h(L)$ is regular.

let $\mathcal{A} = (Q, Q^f, \Delta)$ be clean, we build $\mathcal{A}' = (Q', Q'_f, \Delta')$.

For each $r = f(q_1, \dots, q_n) \rightarrow q \in \Delta$, with $t_f \in \mathcal{T}(\Sigma', \mathcal{X}_n)$ (linear),

let $Q^r = \{q_p^r \mid p \in \text{Pos}(t_f)\}$, and Δ_r defined as follows:

for all $p \in \text{Pos}(t_f)$:

- ▶ if $t_f(p) = g \in \Sigma'_m$, then $g(q_{p_1}^r, \dots, q_{p_m}^r) \rightarrow q_p^r \in \Delta_r$,
- ▶ if $t_f(p) = x_i$, then $q_i \xrightarrow{\varepsilon} q_p^r \in \Delta_r$,
- ▶ $q_\varepsilon^r \xrightarrow{\varepsilon} q \in \Delta_r$.

$$Q' = Q \cup \bigcup_{r \in \Delta} Q^r,$$

$$Q'_f = Q_f,$$

$$\Delta' = \bigcup_{r \in \Delta} \Delta_r.$$

It holds that $h(L(\mathcal{A})) = L(\mathcal{A}')$.

Closure of Regular Languages under Linear Homomorphisms

This is not true in general for the non-linear homomorphisms.

Closure of Regular Languages under Linear Homomorphisms

This is not true in general for the non-linear homomorphisms.

Example : Non-linear homomorphisms

$$\Sigma = \{a : 0, g : 1, f : 1\}, \Sigma' = \{a : 0, g : 1, f' : 2\},$$

$$h : \mathcal{T}(\Sigma) \rightarrow \mathcal{T}(\Sigma') \text{ with } t_a = a, t_g = g(x_1), t_f = f'(x_1, x_1).$$

$$\text{Let } L = \{f(g^n(a)) \mid n \geq 0\},$$

$$h(L) = \{f'(g^n(a), g^n(a)) \mid n \geq 0\} \text{ is not regular.}$$

Closure of Regular Languages under Inverse Homomorphisms

Theorem :

For all regular languages L and all homomorphisms h , $h^{-1}(L)$ is regular.

$\mathcal{A}' = (Q', Q'_f, \Delta')$ complete deterministic such that $L(\mathcal{A}') = L$.

We construct $\mathcal{A} = (Q, Q_f, \Delta)$ with $Q = Q' \uplus \{q_\forall\}$ $Q_f = Q'_f$ and Δ is defined by:

- ▶ for $a \in \Sigma_0$, if $t_a \xrightarrow{\mathcal{A}'}^* q$ then $a \rightarrow q \in \Delta$;
- ▶ for all $f \in \Sigma_n$ with $n > 0$, for $p_1, \dots, p_n \in Q$, if $t_f\{x_1 \mapsto p_1, \dots, x_n \mapsto p_n\} \xrightarrow{\mathcal{A}'}^* q$ then $f(q_1, \dots, q_n) \rightarrow q \in \Delta$ where $q_i = p_i$ if x_i occurs in t_f and $q_i = q_\forall$ otherwise;
- ▶ for $a \in \Sigma_0$, $a \rightarrow q_\forall \in \Delta$;
- ▶ for $f \in \Sigma_n$ where $n > 0$, $f(q_\forall, \dots, q_\forall) \rightarrow q_\forall \in \Delta$.

It holds that $t \xrightarrow{\mathcal{A}}^* q$ iff $h(t) \xrightarrow{\mathcal{A}'}^* q$ for all $q \in Q'$.

Closure under Homomorphisms

Theorem :

The class of regular tree languages is the smallest non trivial class of sets of trees closed under linear homomorphisms and inverse homomorphisms.

A problem whose decidability has been open for 35 years:

INPUT: a TA \mathcal{A} , an homomorphism h

QUESTION: is $h(L(\mathcal{A}))$ regular?

Tree Transducers

Tree Transducers

Definition : Bottom-up Tree Transducers

A *bottom-up tree transducer* (TT) is a tuple $U = (\Sigma, \Sigma', Q, Q^f, \Delta)$ where

- ▶ Σ, Σ' are the input, resp. output, signatures,
- ▶ Q is a finite set of *states*,
- ▶ $Q^f \subseteq Q$ is the subset of final states
- ▶ Δ is a set of transduction (rewrite) rules of the form:
 - ▶ $f(p_1(x_1), \dots, p_n(x_n)) \rightarrow p(u)$ with $f \in \Sigma_n$ ($n \geq 0$), $p_1, \dots, p_n, p \in Q$, x_1, \dots, x_n pairwise distinct and $u \in \mathcal{T}(\Sigma', \{x_1, \dots, x_n\})$, or
 - ▶ $p(x_1) \rightarrow p'(u)$ with $q, q' \in Q$, $u \in \mathcal{T}(\Sigma', \{x_1\})$.

A TT is *linear* if all the u in transduction rules are linear.

The transduction relation of U is the binary relation:

$$L(U) = \{ \langle t, t' \rangle \mid t \xrightarrow[U]{*} q(t'), t \in \mathcal{T}(\Sigma), t' \in \mathcal{T}(\Sigma'), q \in Q^f \}$$

Example 1

$$U_1 = (\{f : 1, a : 0\}, \{g : 2, f, f' : 1, a : 0\}, \{q, q'\}, \{q'\}, \Delta_1),$$

$$\Delta_1 = \left\{ \begin{array}{l} a \rightarrow q(a) \\ f(q(x_1)) \rightarrow q(f(x_1)) \mid q(f'(x_1)) \mid q'(g(x_1, x_1)) \end{array} \right\}$$

Example 2

$$\Sigma_{in} = \{f : 2, g : 1, a : 0\},$$

$$U_2 = (\Sigma_{in}, \Sigma_{in} \cup \{f' : 1\}, \{q, q', q_f\}, \{q_f\}, \Delta_2),$$

$$\Delta_2 = \left\{ \begin{array}{l} a \rightarrow q(a) \mid q'(a) \\ g(q(x_1)) \rightarrow q(g(x_1)) \\ g(q'(x_1)) \rightarrow q'(g(x_1)) \\ f(q'(x_1), q'(x_2)) \rightarrow q'(f(x_1, x_2)) \\ f(q'(x_1), q'(x_2)) \rightarrow q_f(f'(x_1)) \end{array} \right\}$$

$$L(U_2) = \{ \langle f(t_1, t_2), f'(t_1) \mid t_2 = g^m(a), m \geq 0 \rangle \}$$

Tree Transducers, Example

Token tree protocol [Abdulla et al CAV02]

$$\begin{aligned} \underline{n} &\rightarrow q_0(\underline{n'}) \\ \underline{t} &\rightarrow q_1(\underline{n'}) \\ n(q_0(x_1), q_0(x_2)) &\rightarrow q_0(n(x_1, x_2)) \\ t(q_0(x_1), q_0(x_2)) &\rightarrow q_1(n(x_1, x_2)) \\ n(q_1(x_1), q_0(x_2)) &\rightarrow q_2(t(x_1, x_2)) \\ n(q_0(x_1), q_1(x_2)) &\rightarrow q_2(t(x_1, x_2)) \\ n(q_2(x_1), q_0(x_2)) &\rightarrow q_2(n(x_1, x_2)) \\ n(q_0(x_1), q_2(x_2)) &\rightarrow q_2(n(x_1, x_2)) \end{aligned}$$

property: mutual exclusion (for every network)

initial: terms of $\mathcal{T}(\{t, n, \underline{t}, \underline{n}\})$, containing exactly one token.

verification: the intersection of his closure with the set $\{q_2(t) \mid t \in \mathcal{T}(\{t, n, \underline{t}, \underline{n}\}), t \text{ contains at least 2 tokens}\}$ (regular) is empty.

Languages

- ▶ Linear bottom-up TT are closed under composition.
- ▶ Deterministic bottom-up TT are closed under composition.

Theorem :

- ▶ The domain of a TT is a regular tree language.
- ▶ The image of a regular tree language by a linear TT is a regular tree language.

Transducers and Homomorphisms

An homomorphism is called *delabeling* if it is linear, complete, symbol-to-symbol.

Definition : Bimorphisms

A *bimorphism* is a triple $B = (h, h', L)$ where h, h' are homomorphisms and L is a regular tree language.

$$L(B) = \{ \langle h(t), h'(t) \rangle \mid t \in L \}$$

Theorem :

TT \equiv bimorphisms (h, h', L) where h delabeling.

Term Rewriting Systems

Term Rewriting

Definition : Substitution

A *substitution* is a function of finite domain from \mathcal{X} into $\mathcal{T}(\Sigma, \mathcal{X})$. We extend the definition to $\mathcal{T}(\Sigma, \mathcal{X}) \rightarrow \mathcal{T}(\Sigma, \mathcal{X})$ by:

$$f(t_1, \dots, t_n)\sigma = f(t_1\sigma, \dots, t_n\sigma) \quad (n \geq 0)$$

The application $C[t_1, \dots, t_n]$ of a context $C \in \mathcal{T}(\Sigma, \{x_1, \dots, x_n\})$ to n terms t_1, \dots, t_n , is $C\sigma$ with $\sigma = \{x_1 \mapsto t_1, \dots, x_n \mapsto t_n\}$.

Term Rewriting

A *rewrite system* \mathcal{R} is a finite set of rewrite rules of the form $\ell \rightarrow r$ with $\ell, r \in \mathcal{T}(\Sigma, \mathcal{X})$.

The relation $\xrightarrow{\mathcal{R}}$ is the smallest binary relation containing \mathcal{R} , and closed under application of contexts and substitutions.

i.e. $s \xrightarrow{\mathcal{R}} t$ iff $\exists p \in \mathcal{Pos}(s), \ell \rightarrow r \in \mathcal{R}, \sigma, s|_p = \ell\sigma$ and $t = s[r\sigma]_p$.

We note $\xrightarrow{\mathcal{R}^*}$ the reflexive and transitive closure of $\xrightarrow{\mathcal{R}}$.

Example :

$\mathcal{R} = \{+(0, x) \rightarrow x, +(s(x), y) \rightarrow s(+ (x, y))\}$.

$$\begin{array}{lcl} +(s(s(0)), +(0, s(0))) & \xrightarrow{\mathcal{R}} & +(s(s(0)), s(0)) \\ & \xrightarrow{\mathcal{R}} & s(+ (s(0), s(0))) \\ & \xrightarrow{\mathcal{R}} & s(s(+ (0, s(0)))) \\ & \xrightarrow{\mathcal{R}} & s(s(s(0))) \end{array}$$

TRS Preserving Regularity

For a TRS \mathcal{R} over Σ and $L \subseteq \mathcal{T}(\Sigma)$,

$$\mathcal{R}^*(L) = \{t \in \mathcal{T}(\Sigma) \mid \exists s \in L, s \xrightarrow{*}_{\mathcal{R}} t\}$$

Regularity Preservation

Identify a class \mathcal{C} of TRS such that
for all $\mathcal{R} \in \mathcal{C}$, $\mathcal{R}^*(L)$ is regular if L is regular.

Theorem : [Gilleron STACS 91]

It is undecidable in general whether a given TRS is preserving regularity.

Ground TRS

Theorem : [Brainerd 69]

Ground TRS are preserving regularity.

Given: TA \mathcal{A}_{in} and ground TRS \mathcal{R} . We start with

$$\mathcal{A}_{in} \cup (\Sigma, Q_{\mathcal{R}}, \emptyset, \{f(q_{r_1}, \dots, q_{r_n}) \rightarrow q_r \mid r = f(r_1, \dots, r_n) \in Q_{\mathcal{R}}\})$$

where $Q_{\mathcal{R}} = \text{strict subterms}(\text{rhs}(\mathcal{R}))$,

and add transitions according to the schema:

$$\begin{array}{ccc} \text{lhs}(\mathcal{R}) \ni \ell & \xrightarrow{\mathcal{A}} & q \\ \downarrow \mathcal{R} & & \uparrow \mathcal{A} \\ f(r_1, \dots, r_n) & \xrightarrow{\mathcal{A}} & f(q_{r_1}, \dots, q_{r_n}) \end{array}$$

no states are added \rightarrow termination.

The TA obtained recognizes $\mathcal{R}^*(L(\mathcal{A}_{in}))$.

Ground TRS (examples)

$$\begin{array}{ccc} \text{lhs}(\mathcal{R}) \ni \ell & \xrightarrow{\mathcal{A}} & q \\ \downarrow \mathcal{R} & & \uparrow \mathcal{A} \\ f(r_1, \dots, r_n) & \xrightarrow{\mathcal{A}} & f(q_{r_1}, \dots, q_{r_n}) \end{array}$$

$s(s(0)) \rightarrow 0$	$\perp + 1 \rightarrow s(\perp)$
$\begin{array}{ccc} s(s(0)) & \xrightarrow[*]{\mathcal{A}} & q \\ \downarrow \mathcal{R} & \nearrow \mathcal{A} & \\ 0 & & \end{array}$	$\begin{array}{ccc} \perp + 1 & \xrightarrow{\mathcal{A}} & q \\ \downarrow \mathcal{R} & & \uparrow \mathcal{A} \\ s(\perp) & \xrightarrow{\mathcal{A}} & s(q_\perp) \end{array}$

Linear and Right-Shallow TRS

right-shallow: variables at depth at most 1 in rhs of rules.

Theorem : [Salomaa 88]

Linear and right-shallow TRS preserve regularity.

Given: TA \mathcal{A}_{in} and linear and right-shallow TRS \mathcal{R} .

The construction is similar to the ground TRS case: We start with

$$\mathcal{A}_{in} \cup (\Sigma, Q_{\mathcal{R}}, \emptyset, \{f(q_{r_1}, \dots, q_{r_n}) \rightarrow q_r \mid r = f(r_1, \dots, r_n) \in Q_{\mathcal{R}}\})$$

where $Q_{\mathcal{R}} = \text{strict subterms}(\text{rhs}(\mathcal{R})) \setminus \mathcal{X}$,

and add transitions according to the schema:

$$\begin{array}{ccc} \ell\sigma & \xrightarrow{\mathcal{A}} & q \\ \downarrow \mathcal{R} & & \uparrow \mathcal{A} \\ f(r_1, \dots, r_n)\sigma & \xrightarrow{\mathcal{A}} & f(q_1, \dots, q_n) \end{array}$$

where $\ell \in \text{lhs}(\mathcal{R})$, substitution $\sigma : \text{vars}(\ell) \rightarrow Q$, for all $i \leq n$, if $r_i \notin \mathcal{X}$ then $q_i = q_{r_i}$ and $q_i = r_i\sigma$ otherwise.

Linear and Right-Shallow TRS (Examples)

$$\begin{array}{ccc}
 \ell\sigma & \xrightarrow{\mathcal{A}} & q \\
 \downarrow \mathcal{R} & & \uparrow \mathcal{A} \\
 f(r_1, \dots, r_n)\sigma & \xrightarrow{\mathcal{A}} & f(q_1, \dots, q_n)
 \end{array}$$

where $\ell \in lhs(\mathcal{R})$, substitution $\sigma : vars(\ell) \rightarrow Q$, for all $i \leq n$, if $r_i \notin \mathcal{X}$ then $q_i = q_{r_i}$ and $q_i = r_i\sigma$ otherwise.

$s(x) - s(y) \rightarrow x - y$	$s(x) \rightarrow s(0) + x$
$ \begin{array}{ccc} s(q_1) - s(q_2) & \xrightarrow{\mathcal{A}} & q'_1 - q'_2 \xrightarrow{\mathcal{A}} q \\ \downarrow \mathcal{R} & & \nearrow \mathcal{A} \\ q_1 - q_2 & & \end{array} $	$ \begin{array}{ccc} s(q_1) & \xrightarrow{\mathcal{A}} & q \\ \downarrow \mathcal{R} & & \uparrow \mathcal{A} \\ s(0) + q_1 & \xrightarrow{\mathcal{A}} & q_{s(0)} + q_1 \end{array} $

Linear and Right-Shallow TRS: Extensions

Other classes of TRS preserving regularity

- ▶ [Coquide et al 94] *semi-monadic* or *inverse-growing* TRS:
for all $\ell \rightarrow r \in \mathcal{R}$, $\text{vars}(r) \cap \text{vars}(\ell)$ at depth at most 1 in r .
- ▶ [Nagaya Toyama RTA 02] right-linear and right-shallow TRS.
NOT left-linear.
- ▶ [Gyenize Vagvolgyi GSMTRS 98]
linear and *generalized semi-monadic* TRS
- ▶ [Takai Kaji Seki RTA 00]
right-linear *finite path overlapping* TRS

Right-Linearity and Right-Shalowness Conditions

Relaxing these conditions generally breaks regularity preservation.

Example : Right-Linearity

let $\mathcal{R} = \{f(x) \rightarrow g(x, x)\}$ (flat and left-linear), $L_{\text{in}} = \{f(\dots f(c))\}$.
 $\mathcal{R}^*(L_{\text{in}}) \cap \mathcal{T}(\{g, c\})$ is the set of balanced binary trees of $\mathcal{T}(\{g, c\})$,
which is not regular.

Example : Right-Shalowness

With rewrite rules whose left and right hand-side have height at most two, it is possible simulate Turing machine computations, even in the case of words (symbols of arity 0 or 1).

Exceptions (for the right-shalowness)

- ▶ [Rety LPAR 99] constructor based (with restrictions on L_{in}).
ex: $\text{app}(\text{nil}, y) \rightarrow y$, $\text{app}(\text{cons}(x, y), z) \rightarrow \text{cons}(x, \text{app}(y, z))$.
- ▶ [Seki et al RTA 02] Layered Transducing TRS

Linear I/O Separated Layered Transducing TRS

[Seki et al RTA 02]

This class corresponds to linear tree transducers.

over $\Sigma = \Sigma_i \uplus \Sigma_o \uplus Q$, rewrite rules of the form

$$\begin{aligned} f_i(p_1(x_1), \dots, p_n(x_n)) &\rightarrow p(t) \\ p'_1(x_1) &\rightarrow p'(t') \end{aligned}$$

where $f_i \in \Sigma_i$, $p_1, \dots, p_n, p, p'_1, p' \in Q$ x_1, \dots, x_n are disjoint variables, $t, t' \in \mathcal{T}(\Sigma_o, \mathcal{X})$ such that $\text{vars}(t) \subseteq \{x_1, \dots, x_n\}$ and $\text{vars}(t') \subseteq \{x_1\}$.

To know more

Further results closure of tree automata languages:

- ▶ closure of extended tree automata languages, modulo
[Gallagher Rosendahl 08], [JRV JLAP 08], [JKV LATA 09],
[JKV IC 11]
- ▶ rewrite strategies (bottom-up, context-sensitive, innermost,
outermost...) [Durand et al RTA 07,10,11],
[Kojima Sakai RTA 08], [Rety Vuotto JSC 05], [GGJ WRS 08]
- ▶ constrained/controlled rewriting
[Sénizergues *French Spring School of TCS 93*],
[JKS FroCoS 11]
- ▶ unranked tree rewriting (XML updates)
[JR RTA 08], [JR PPDP 10]

Tree Automata Based Program Verification

Some Techniques and Tools

Program Analysis with Tree Automata / Grammars

(very partial list) focus on 3 approaches

- ▶ [Reynolds IP 68] LISP programs → lfp solutions of equations
- ▶ [Jones Muchnick POPL 79] LISP programs → tree grammars
- ▶ [Jones 87] lazy higher-order functional programs
- ▶ [Heintze Jaffar 90] logic programs → set constraints
- ▶ [Lugiez Schnoebelen CONCUR 98], [Bouajjani Touili 03+] imperative programs w. prefix rewriting: PA-processes, PAD systems, PRS...
- ▶ [Genet et al 98+] functional programs, security protocols, Java Bytecode
- ▶ [Jones Andersen TCS 07] functional programs

Timbuk

[Genet et al] (IRISA)

<http://www.irisa.fr/celtique/genet/timbuk>

Computation of rewrite closure by tree automata completion, with *over-approximations*. User defined or inferred accelerations.

- ▶ analysis of security protocols
SmartRight, Copy Protection Technology for DVB, Thomson
- ▶ analysis of Java Bytecode with Copster

Timbuk library, used in other tools like

- ▶ TA4SP, one of the proof back-ends of the AVISPA tool for security protocol verification
- ▶ SPADE

[Tayssir Touili et al CAV 07] (LIAFA).

<http://www.liafa.jussieu.fr/~touili/spade.html>

Reachability analysis for multithreaded dynamic and recursive programs.

- ▶ (PAD) Systems [Touili VISSAS 05]

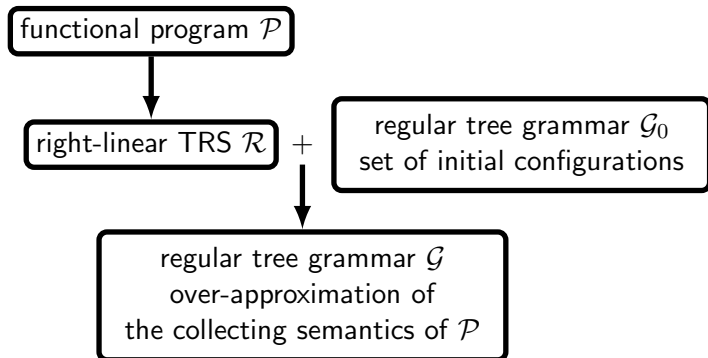
$$X_1 \cdot \dots \cdot X_n \rightarrow Y_1 \cdot \dots \cdot Y_m, \quad X_1 \rightarrow Y_1 \parallel \dots \parallel Y_m$$

Case studies

- ▶ Windows Bluetooth driver
- ▶ multithreaded program based on the class `java.util.Vector` from the Java Standard Collection Framework
- ▶ concurrent insertions on a binary search tree

Approximations of Collecting Semantics

[Jones Andersen TCS 07]



collecting semantics [Cousot²] (roughly): mapping associating to each program point p the set of configurations reachable at p .

[Kochems Ong RTA 11] finer approximation using indexed linear tree grammars (instead of regular grammars).

Regular Tree Grammars

Definition : Regular Tree Grammars

A is a tuple $\mathcal{G} = \langle \mathcal{N}, S, \Sigma, P \rangle$ where \mathcal{N} is a finite set of nullary *non-terminal* symbols, $S \in \mathcal{N}$ (*axiom* of \mathcal{G}), Σ is a signature disjoint from \mathcal{N} and P is a set of *production rules* of the form $X := r$ with $r \in \mathcal{T}(\Sigma \cup \mathcal{N})$.

Example :

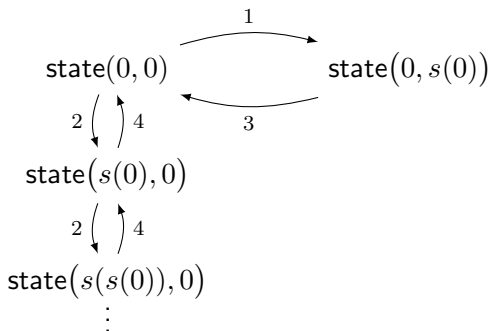
$\Sigma = \{\wedge : 2, \vee : 2, \neg : 1, \top, \perp : 0\}$, $\mathcal{G} = (\{X_0, X_1\}, X_1, \Sigma, P)$.

$$P = \left\{ \begin{array}{ll} X_0 := \perp & X_1 := \top \\ X_1 := \neg(X_0) & X_0 := \neg(X_1) \\ X_0 := \vee(X_0, X_0) & X_1 := \vee(X_0, X_1) \\ X_1 := \vee(X_1, X_0) & X_1 := \vee(X_1, X_1) \\ X_0 := \wedge(X_0, X_0) & X_0 := \wedge(X_0, X_1) \\ X_0 := \wedge(X_1, X_0) & X_1 := \wedge(X_1, X_1) \end{array} \right\}$$

Approximations of Collecting Semantics: Example

Concurrent readers/writers: reachable configurations

$$\begin{aligned}\mathcal{R} = R_1 : & \quad \text{state}(0, 0) \rightarrow \text{state}(0, s(0)) \\ R_2 : & \quad \text{state}(X_2, 0) \rightarrow \text{state}(s(X_2), 0) \\ R_3 : & \quad \text{state}(X_3, s(Y_3)) \rightarrow \text{state}(X_3, Y_3) \\ R_4 : & \quad \text{state}(s(X_4), Y_4) \rightarrow \text{state}(X_4, Y_4)\end{aligned}$$



Approximations of Collecting Semantics: Example

$$\begin{aligned} \mathcal{R} = R_1 : & \quad \text{state}(0, 0) \rightarrow \text{state}(0, s(0)) \\ R_2 : & \quad \text{state}(X_2, 0) \rightarrow \text{state}(s(X_2), 0) \\ R_3 : & \quad \text{state}(X_3, s(Y_3)) \rightarrow \text{state}(X_3, Y_3) \\ R_4 : & \quad \text{state}(s(X_4), Y_4) \rightarrow \text{state}(X_4, Y_4) \end{aligned}$$

$$R_0 := \text{state}(0, 0)$$

$R_0 := R_1$ $R_1 := \text{state}(0, s(0))$	$\text{state}(0, 0) = \text{lhs}(R_1)$
$R_0 := R_2$ $R_2 := \text{state}(s(X_2), 0)$ $X_2 := 0$	$\text{state}(0, 0) = \text{state}(X_2, 0)\{X_2 \mapsto 0\}$
$X_2 := s(X_2)$	$\text{state}(s(X_2), 0) =$ $\text{state}(X_2, 0)\{X_2 \mapsto s(X_2)\}$
$R_1 := R_3$ $R_3 := \text{state}(X_3, Y_3)$ $X_3 := 0, Y_3 := 0$	$\text{state}(0, s(0)) =$ $\text{state}(X_3, s(Y_3))\{X_3 \mapsto 0, Y_3 \mapsto 0\}$
$R_2 := R_4$ $R_4 := \text{state}(s(X_4), Y_4)$ $X_4 := X_2, Y_4 := 0$	$\text{state}(s(X_2), 0) =$ $\text{state}(s(X_4), Y_4)\{X_4 \mapsto X_2, Y_4 \mapsto 0\}$

Approximations of Collecting Semantics: Example

$$\begin{aligned} \mathcal{R} = R_1 : & \quad \text{state}(0, 0) \rightarrow \text{state}(0, s(0)) \\ R_2 : & \quad \text{state}(X_2, 0) \rightarrow \text{state}(s(X_2), 0) \\ R_3 : & \quad \text{state}(X_3, s(Y_3)) \rightarrow \text{state}(X_3, Y_3) \\ R_4 : & \quad \text{state}(s(X_4), Y_4) \rightarrow \text{state}(X_4, Y_4) \end{aligned}$$

$$R_0 := \text{state}(0, 0)$$

$$R_0 := R_1$$

$$R_1 := \text{state}(0, s(0))$$

$$R_0 := R_2$$

$$R_2 := \text{state}(s(X_2), 0)$$

$$X_2 := 0$$

$$X_2 := s(X_2)$$

$$R_1 := R_3$$

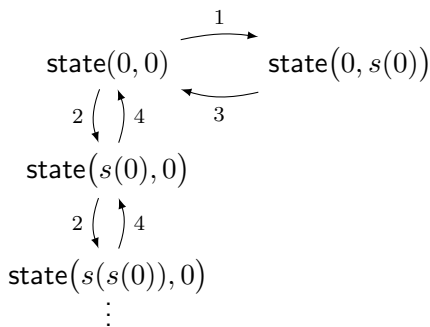
$$R_3 := \text{state}(X_3, Y_3)$$

$$X_3 := 0, Y_3 := 0$$

$$R_2 := R_4$$

$$R_4 := \text{state}(s(X_4), Y_4)$$

$$X_4 := X_2, Y_4 := 0$$



Approximations of Collecting Semantics: Example 2

[Jones Andersen TCS 07]

```
let rec first l1 l2 =  
  match l1, l2 with  
  [], - → []  
  l::m, x::xs → x::(first m xs);
```

$R_2 : \text{first}(\text{nil}, X_s) \rightarrow \text{nil}$

$R_3 : \text{first}(\text{cons}(1, M), \text{cons}(X, X_s)) \rightarrow \text{cons}(X, \text{first}(M, X_s))$

```
let rec sequence y =  
  y::(sequence (1::y));
```

$R_4 : \text{sequence}(Y) \rightarrow \text{cons}(Y, \text{sequence}(\text{cons}(1, Y)))$

```
let g n =  
  first n (sequence []);
```

$R_1 : g(N) \rightarrow \text{first}(N, \text{sequence}(\text{nil}))$

Part II

Weak Second Order Monadic Logic with k Successors

Logic and Automata

- ▶ logic for expressing properties of labeled binary trees
= specification of tree languages,

Logic and Automata

- ▶ **logic** for expressing properties of labeled binary trees
= **specification** of tree languages, example:

$$t \models \forall x a(x) \Rightarrow \exists y y > x \wedge b(y)$$

- ▶ compilation of formulae into **automata**
= decision **algorithms**.
- ▶ equivalence between both formalisms
[Thatcher & Wright's theorem].

Plan

WSkS: Definition

Automata \rightarrow Logic

Logic \rightarrow Automata

Fragments and Extensions of WSkS

Interpretation Structures

$\mathcal{L} :=$ set of **predicate** symbols P_1, \dots, P_n with arity.

A **structure** \mathcal{M} over \mathcal{L} is a tuple

$$\mathcal{M} := \langle \mathcal{D}, P_1^{\mathcal{M}}, \dots, P_n^{\mathcal{M}} \rangle$$

where

- ▶ \mathcal{D} is the **domain** of \mathcal{M} ,
- ▶ every $P_i^{\mathcal{M}}$ (**interpretation** of P_i) is a subset of $\mathcal{D}^{\text{arity}(P_i)}$ (relation).

Term as Structure

Σ signature, $k = \text{maximal arity}$.

$$\mathcal{L}_\Sigma := \{=, <, S_1, \dots, S_k, L_a \mid a \in \Sigma\}.$$

to $t \in \mathcal{T}(\Sigma)$, we associate a **structure** \underline{t} over \mathcal{L}_Σ

$$\underline{t} := \langle \mathcal{Pos}(t), =, <, S_1, \dots, S_k, L_a^t, L_b^t, \dots \rangle$$

where

- ▶ domain = positions of t ($\mathcal{Pos}(t) \subset \{1, \dots, k\}^*$)
- ▶ $=$ equality over $\mathcal{Pos}(t)$,
- ▶ $<$ prefix ordering over $\mathcal{Pos}(t)$,
- ▶ $S_i = \{\langle p, p \cdot i \rangle \mid p, p \cdot i \in \mathcal{Pos}(t)\}$ (i^{th} successor position),
- ▶ $L_a^t = \{p \in \mathcal{Pos}(t) \mid t(p) = a\}$.

FOL with k Successors

- ▶ first order variables x, y, \dots
- ▶ form ::= $x = y \mid x < y$
| $S_1(x, y) \mid \dots \mid S_k(x, y) \mid L_a(x) \quad a \in \Sigma$
| form \wedge form | form \vee form | \neg form
| $\exists x$ form | $\forall x$ form

Notation: $\phi(x_1, \dots, x_m)$,

where x_1, \dots, x_m are the free variables of ϕ .

WSkS: Syntax

- ▶ first order variables x, y, \dots
- ▶ second order variables X, Y, \dots
- ▶ form ::= $x = y \mid x < y \mid x \in X$
 $\mid S_1(x, y) \mid \dots \mid S_k(x, y) \mid L_a(x) \quad a \in \Sigma$
 $\mid \text{form} \wedge \text{form} \mid \text{form} \vee \text{form} \mid \neg \text{form}$
 $\mid \exists x \text{ form} \mid \exists X \text{ form} \mid \forall x \text{ form} \mid \forall X \text{ form}$

Notation: $\phi(x_1, \dots, x_m, X_1, \dots, X_n)$,

where $x_1, \dots, x_m, X_1, \dots, X_n$ are the free variables of ϕ .

WSkS: Semantics

- ▶ $t \in \mathcal{T}(\Sigma)$,
- ▶ valuation σ of first order variables into $\mathcal{Pos}(t)$,
- ▶ valuation δ of second order variables into subsets of $\mathcal{Pos}(t)$,
- ▶ $\underline{t}, \sigma, \delta \models x = y$ iff $\sigma(x) = \sigma(y)$,
- ▶ $\underline{t}, \sigma, \delta \models x < y$ iff $\sigma(x) <_{\text{prefix}} \sigma(y)$,
- ▶ $\underline{t}, \sigma, \delta \models x \in X$ iff $\sigma(x) \in \delta(X)$,
- ▶ $\underline{t}, \sigma, \delta \models S_i(x, y)$ iff $\sigma(y) = \sigma(x) \cdot i$,
- ▶ $\underline{t}, \sigma, \delta \models L_a(x)$ iff $t(\sigma(x)) = a$ i.e. $\sigma(x) \in L_a^{\underline{t}}$,
- ▶ $\underline{t}, \sigma, \delta \models \phi_1 \wedge \phi_2$ iff $\underline{t}, \sigma, \delta \models \phi_1$ and $\underline{t}, \sigma, \delta \models \phi_2$,
- ▶ $\underline{t}, \sigma, \delta \models \phi_1 \vee \phi_2$ iff $\underline{t}, \sigma, \delta \models \phi_1$ or $\underline{t}, \sigma, \delta \models \phi_2$,
- ▶ $\underline{t}, \sigma, \delta \models \neg\phi$ iff $\underline{t}, \sigma, \delta \not\models \phi$,

WSkS: Semantics (Quantifiers)

- ▶ $\underline{t}, \sigma, \delta \models \exists x \phi$ iff $x \notin \text{dom}(\sigma)$, x free in ϕ
and exists $p \in \text{Pos}(t)$ s.t. $\underline{t}, \sigma \cup \{x \mapsto p\}, \delta \models \phi$,
- ▶ $\underline{t}, \sigma, \delta \models \forall x \phi$ iff $x \notin \text{dom}(\sigma)$, x free in ϕ
and for all $p \in \text{Pos}(t)$, $\underline{t}, \sigma \cup \{x \mapsto p\}, \delta \models \phi$,
- ▶ $\underline{t}, \sigma, \delta \models \exists X \phi$ iff $X \notin \text{dom}(\delta)$, X free in ϕ
and exists $P \subseteq \text{Pos}(t)$ s.t. $\underline{t}, \sigma, \delta \cup \{X \mapsto P\} \models \phi$,
- ▶ $\underline{t}, \sigma, \delta \models \forall X \phi$ iff $X \notin \text{dom}(\delta)$, X free in ϕ
and for all $P \subseteq \text{Pos}(t)$, $\underline{t}, \sigma, \delta \cup \{X \mapsto P\} \models \phi$.

WSkS: Languages

Definition : WSkS-definability

For $\phi \in \text{WSkS}$ closed (without free variables) over \mathcal{L}_Σ ,

$$L(\phi) := \{t \in \mathcal{T}(\Sigma) \mid \underline{t} \models \phi\}.$$

Example :

$\Sigma = \{a : 2, b : 2, c : 0\}$. Language of terms in $\mathcal{T}(\Sigma)$

- ▶ containing the pattern $a(b(x_1, x_2), x_3)$:
 $\exists x \exists y S_1(x, y) \wedge L_a(x) \wedge L_b(y)$
- ▶ such that every a -labelled node has a b -labelled child.
 $\forall x \exists y L_a(x) \Rightarrow \bigvee_{i=1}^2 S_i(x, y) \wedge L_b(y)$
- ▶ such that every a -labelled node has a b -labelled descendant.
 $\forall x \exists y L_a(x) \Rightarrow x < y \wedge L_b(y)$

WS k S: Examples

- ▶ root position:

WSkS: Examples

- ▶ root position: $\text{root}(x) \equiv \neg \exists y \ y < x$
- ▶ inclusion:

WSkS: Examples

- ▶ root position: $\text{root}(x) \equiv \neg \exists y \ y < x$
- ▶ inclusion: $X \subseteq Y \equiv \forall x (x \in X \Rightarrow x \in Y)$
- ▶ intersection:

WSkS: Examples

- ▶ root position: $\text{root}(x) \equiv \neg \exists y \ y < x$
- ▶ inclusion: $X \subseteq Y \equiv \forall x (x \in X \Rightarrow x \in Y)$
- ▶ intersection: $Z = X \cap Y \equiv \forall x (x \in Z \Leftrightarrow (x \in X \wedge x \in Y))$
- ▶ emptiness:

WSkS: Examples

- ▶ root position: $\text{root}(x) \equiv \neg \exists y y < x$
- ▶ inclusion: $X \subseteq Y \equiv \forall x (x \in X \Rightarrow x \in Y)$
- ▶ intersection: $Z = X \cap Y \equiv \forall x (x \in Z \Leftrightarrow (x \in X \wedge x \in Y))$
- ▶ emptiness: $X = \emptyset \equiv \forall x x \notin X$
- ▶ finite union:

WSkS: Examples

- ▶ root position: $\text{root}(x) \equiv \neg \exists y y < x$
- ▶ inclusion: $X \subseteq Y \equiv \forall x (x \in X \Rightarrow x \in Y)$
- ▶ intersection: $Z = X \cap Y \equiv \forall x (x \in Z \Leftrightarrow (x \in X \wedge x \in Y))$
- ▶ emptiness: $X = \emptyset \equiv \forall x x \notin X$
- ▶ finite union:
$$X = \bigcup_{i=1}^n X_i \equiv \left(\bigwedge_{i=1}^n X_i \subseteq X \right) \wedge \forall x (x \in X \Rightarrow \bigvee_{i=1}^n x \in X_i)$$
- ▶ partition:

WSkS: Examples

- ▶ root position: $\text{root}(x) \equiv \neg \exists y y < x$
- ▶ inclusion: $X \subseteq Y \equiv \forall x (x \in X \Rightarrow x \in Y)$
- ▶ intersection: $Z = X \cap Y \equiv \forall x (x \in Z \Leftrightarrow (x \in X \wedge x \in Y))$
- ▶ emptiness: $X = \emptyset \equiv \forall x x \notin X$

- ▶ finite union:

$$X = \bigcup_{i=1}^n X_i \equiv \left(\bigwedge_{i=1}^n X_i \subseteq X \right) \wedge \forall x (x \in X \Rightarrow \bigvee_{i=1}^n x \in X_i)$$

- ▶ partition:

$$X_1, \dots, X_n \text{ partition } X \equiv X = \bigcup_{i=1}^n X_i \wedge \bigwedge_{i=1}^{n-1} \bigwedge_{j=i+1}^n X_i \cap X_j = \emptyset$$

WS k S: Examples (2)

- ▶ singleton:

WSkS: Examples (2)

- ▶ singleton:

$$\text{sing}(X) \equiv X \neq \emptyset \wedge \forall Y (Y \subseteq X \Rightarrow (Y = X \vee Y = \emptyset))$$

- ▶ \leq (without $<$)

WSkS: Examples (2)

- ▶ singleton:

$$\text{sing}(X) \equiv X \neq \emptyset \wedge \forall Y (Y \subseteq X \Rightarrow (Y = X \vee Y = \emptyset))$$

- ▶ \leq (without $<$)

$$x \leq y \equiv \forall X \left(\begin{array}{l} y \in X \\ \wedge \forall z \forall z' (z' \in X \wedge \bigvee_{i \leq k} S_i(z, z')) \Rightarrow z \in X \end{array} \right) \\ \Rightarrow x \in X$$

or

$$x \leq y \equiv \exists X (\forall z z \in X \Rightarrow (\exists z' \bigvee_{i \leq k} S_i(z', z) \wedge z' \in X) \vee z = x) \\ \wedge y \in X$$

Thatcher & Wright's Theorem

Theorem : Thatcher and Wright

Languages of $WSkS$ formulae = regular tree languages.

pr.: 2 directions (2 constructions):

- ▶ $TA \rightarrow WSkS$,
- ▶ $WSkS \rightarrow TA$.

Plan

WSkS: Definition

Automata \rightarrow Logic

Logic \rightarrow Automata

Fragments and Extensions of WSkS

Regular Languages \rightarrow $WSkS$ Languages

Let $\Sigma = \{a_1, \dots, a_n\}$.

Theorem :

For all tree automaton \mathcal{A} over Σ , there exists $\phi_{\mathcal{A}} \in WSkS$ such that $L(\phi_{\mathcal{A}}) = L(\mathcal{A})$.

$\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$ with $Q = \{q_0, \dots, q_m\}$.

$\phi_{\mathcal{A}}$: existence of an accepting run of \mathcal{A} on $t \in \mathcal{T}(\Sigma)$.

$$\phi_{\mathcal{A}} := \exists Y_0 \dots \exists Y_m \phi_{\text{lab}}(\bar{Y}) \wedge \phi_{\text{acc}}(\bar{Y}) \wedge \phi_{\text{tr}_0}(\bar{Y}) \wedge \phi_{\text{tr}}(\bar{Y})$$

Regular Languages \rightarrow WS_kS Languages

$\phi_{\text{lab}}(\overline{Y})$: every position is labeled with one state exactly.

Regular Languages \rightarrow WS_kS Languages

$\phi_{\text{lab}}(\bar{Y})$: every position is labeled with one state exactly.

$$\phi_{\text{lab}}(\bar{Y}) \equiv \forall x \bigvee_{0 \leq i \leq m} x \in Y_i \wedge \bigwedge_{\substack{0 \leq i, j \leq m \\ i \neq j}} (x \in Y_i \Rightarrow \neg x \in Y_j)$$

Regular Languages \rightarrow WS $_k$ S Languages

$\phi_{\text{lab}}(\overline{Y})$: every position is labeled with one state exactly.

$$\phi_{\text{lab}}(\overline{Y}) \equiv \forall x \bigvee_{0 \leq i \leq m} x \in Y_i \wedge \bigwedge_{\substack{0 \leq i, j \leq m \\ i \neq j}} (x \in Y_i \Rightarrow \neg x \in Y_j)$$

$\phi_{\text{acc}}(\overline{Y})$: the root is labeled with a final state

Regular Languages \rightarrow WS k S Languages

$\phi_{\text{lab}}(\overline{Y})$: every position is labeled with one state exactly.

$$\phi_{\text{lab}}(\overline{Y}) \equiv \forall x \bigvee_{0 \leq i \leq m} x \in Y_i \wedge \bigwedge_{\substack{0 \leq i, j \leq m \\ i \neq j}} (x \in Y_i \Rightarrow \neg x \in Y_j)$$

$\phi_{\text{acc}}(\overline{Y})$: the root is labeled with a final state

$$\phi_{\text{acc}}(\overline{Y}) \equiv \forall x_0 \text{ root}(x_0) \Rightarrow \bigvee_{q_i \in Q^f} x_0 \in Y_i$$

Regular Languages \rightarrow WS_kS Languages

$\phi_{\text{tr}_0}(\overline{Y})$: transitions for constants symbols

Regular Languages \rightarrow WSKS Languages

$\phi_{\text{tr}_0}(\overline{Y})$: transitions for constants symbols

$$\phi_{\text{tr}_0}(\overline{Y}) \equiv \bigwedge_{a \in \Sigma_0} \left(\forall x L_a(x) \Rightarrow \bigvee_{a \rightarrow q_i \in \Delta} x \in Y_i \right)$$

Regular Languages \rightarrow WS k S Languages

$\phi_{\text{tr}_0}(\overline{Y})$: transitions for constants symbols

$$\phi_{\text{tr}_0}(\overline{Y}) \equiv \bigwedge_{a \in \Sigma_0} \left(\forall x L_a(x) \Rightarrow \bigvee_{a \rightarrow q_i \in \Delta} x \in Y_i \right)$$

$\phi_{\text{tr}}(\overline{Y})$: transitions for non-constant symbols

Regular Languages \rightarrow WS $_k$ S Languages

$\phi_{\text{tr}_0}(\overline{Y})$: transitions for constants symbols

$$\phi_{\text{tr}_0}(\overline{Y}) \equiv \bigwedge_{a \in \Sigma_0} \left(\forall x L_a(x) \Rightarrow \bigvee_{a \rightarrow q_i \in \Delta} x \in Y_i \right)$$

$\phi_{\text{tr}}(\overline{Y})$: transitions for non-constant symbols

$$\begin{aligned} \phi_{\text{tr}}(\overline{Y}) &\equiv \bigwedge_{f \in \Sigma_j, 0 < j \leq k} \forall x \forall y_1 \dots \forall y_j \\ &\quad (L_f(x) \wedge S_1(x, y_1) \wedge \dots \wedge S_j(x, y_j)) \\ &\quad \Downarrow \\ &\quad \bigvee_{f(q_{i_1}, \dots, q_{i_j}) \rightarrow q_i \in \Delta} x \in Y_i \wedge y_1 \in Y_{i_1} \wedge \dots \wedge y_j \in Y_{i_j} \end{aligned}$$

Plan

WSkS: Definition

Automata \rightarrow Logic

Logic \rightarrow Automata

Fragments and Extensions of WSkS

Thatcher & Wright's Theorem

Theorem :

Every $WSkS$ language is regular.

For all formula $\phi \in WSkS$ over Σ (without free variables) there exists a tree automaton \mathcal{A}_ϕ over Σ , such that $L(\mathcal{A}_\phi) = L(\phi)$.

Corollary :

$WSkS$ is decidable.

pr.: reduction to emptiness decision for \mathcal{A}_ϕ .

Thatcher & Wright's Theorem

\mathcal{A}_ϕ is effectively constructed from ϕ , by induction.

- ▶ automata for atoms
⇒ need of automata for formula **with** free variables.
it will characterize
- ▶ Boolean closures for Boolean connectors.
- ▶ \exists quantifier: projection.

Thatcher & Wright's Theorem

When ϕ contains free variables, \mathcal{A}_ϕ will characterize both terms AND valuations satisfying ϕ : $L(\mathcal{A}_\phi) \equiv \{\langle t, \sigma, \delta \rangle \mid \underline{t}, \sigma, \delta \models \phi\}$.
Below we define the product $\langle t, \sigma, \delta \rangle$.

✓ for free second order variables:

$$\begin{array}{ccc} t \in \mathcal{T}(\Sigma) & & \\ \delta : \{X_1, \dots, X_n\} \rightarrow 2^{\mathcal{P}os(t)} & \mapsto & t \times \delta \in \mathcal{T}(\Sigma \times \{0, 1\}^n) \end{array}$$

arity of $\langle a, \bar{b} \rangle$ in $\Sigma \times \{0, 1\}^n = \text{arity of } a \text{ in } \Sigma$.

for all $p \in \mathcal{P}os(t)$, $(t \times \delta)(p) = \langle t(p), b_1, \dots, b_n \rangle$ where for all $i \leq n$,

- ▶ $b_i = 1$ if $p \in \delta(X_i)$,
- ▶ $b_i = 0$ otherwise.

✓ free first order variables are interpreted as singletons.

We consider a simplified language (wlog).

- ▶ no first order variables,
- ▶ only second order variables $X, Y \dots$,
- ▶ form ::= $X \subseteq Y \mid Y = X \cdot 1 \mid \dots \mid Y = X \cdot k$
 $\left| \begin{array}{l} X \subseteq L_a \quad a \in \Sigma \\ \text{form} \wedge \text{form} \mid \text{form} \vee \text{form} \mid \neg \text{form} \\ \exists X \text{ form} \mid \forall X \text{ form} \end{array} \right.$

interpretation $Y = X \cdot i$: $X = \{x\}$, $Y = \{y\}$ and $y = x \cdot i$.

ex: singleton

We consider a simplified language (wlog).

- ▶ no first order variables,
- ▶ only second order variables $X, Y \dots$,
- ▶ form ::= $X \subseteq Y \mid Y = X \cdot 1 \mid \dots \mid Y = X \cdot k$
 $\left| \begin{array}{l} X \subseteq L_a \quad a \in \Sigma \\ \text{form} \wedge \text{form} \mid \text{form} \vee \text{form} \mid \neg \text{form} \\ \exists X \text{ form} \mid \forall X \text{ form} \end{array} \right.$

interpretation $Y = X \cdot i$: $X = \{x\}$, $Y = \{y\}$ and $y = x \cdot i$.

ex: singleton

$$\text{singleton}(X) \equiv \exists Y \left(Y \subseteq X \wedge Y \neq X \wedge \neg \exists Z (Z \subseteq X \wedge Z \neq X \wedge Z \neq Y) \right)$$

$WSkS \rightarrow WSkS_0$

Lemma :

For all formula $\phi(x_1, \dots, x_m, X_1, \dots, X_n) \in WSkS$,
there exists a formula $\phi'(X'_1, \dots, X'_m, X_1, \dots, X_n) \in WSkS_0$
s.t. $\underline{t}, \sigma, \delta \models \phi(x_1, \dots, x_m, X_1, \dots, X_n)$
iff $\underline{t}, \sigma' \cup \delta \models \phi'(X'_1, \dots, X'_m, X_1, \dots, X_n)$, with $\sigma' : X'_i \mapsto \{\sigma(x_i)\}$.

pr.: several steps of formula rewriting:

1. elimination of $<$,
2. elimination of $S_i(x, y)$ ($i \leq k$), $L_a(x)$ ($a \in \Sigma$),
elimination of first order variables (use singleton(X)).

Compilation of $WSkS_0$ into Automata

notation: $\Sigma_{[m]} := \Sigma \times \{0, 1\}^m$.

For all $\phi(X_1, \dots, X_n) \in WSkS_0$ and $m \geq n$,
we construct a tree automaton $[[\phi]]_m$ over $\Sigma_{[m]}$ recognizing

$$\{t \times \delta \mid \delta : \{X_1, \dots, X_m\} \rightarrow 2^{\mathcal{P}os(t)}, \underline{t}, \delta \models \phi(X_1, \dots, X_n)\}$$

Projection, Cylindrification

projection

$proj_n : \bigcup_{m \geq n} \mathcal{T}(\Sigma_{[m]}) \rightarrow \mathcal{T}(\Sigma_{[n]})$
delete components $n + 1, \dots, m$.

Lemma : projection

For all $n \leq m$, if $L \subseteq \mathcal{T}(\Sigma_{[m]})$ is regular then $proj_n(L)$ is regular.

cylindrification ($m \geq n$)

$cyl_{n,m} : L \subseteq \mathcal{T}(\Sigma_{[n]}) \mapsto \{t \in \mathcal{T}(\Sigma_{[m]}) \mid proj_n(t) \in L\}$

Lemma : cylindrification

For all $n \leq m$, if $L \subseteq \mathcal{T}(\Sigma_{[n]})$ is regular, then $cyl_{n,m}(L)$ is regular.

Compilation: $X_1 \subseteq X_2$

Automaton $\llbracket X_1 \subseteq X_2 \rrbracket_2$:

- ▶ signature $\Sigma_{[2]} = \Sigma \times \{0, 1\}^2$.

Compilation: $X_1 \subseteq X_2$

Automaton $\llbracket X_1 \subseteq X_2 \rrbracket_2$:

- ▶ signature $\Sigma_{[2]} = \Sigma \times \{0, 1\}^2$.
- ▶ states: q_0
- ▶ final states: q_0
- ▶ transitions:

$$\langle a, 0, 0 \rangle (q_0, \dots, q_0) \rightarrow q_0$$

$$\langle a, 0, 1 \rangle (q_0, \dots, q_0) \rightarrow q_0$$

$$\langle a, 1, 1 \rangle (q_0, \dots, q_0) \rightarrow q_0$$

For $m \geq 2$,

$$\llbracket X_1 \subseteq X_2 \rrbracket_m := \text{cyl}_{2,m}(\llbracket X_1 \subseteq X_2 \rrbracket_2)$$

Compilation: $X_1 = X_2 \cdot 1$

Automaton $\llbracket X_1 = X_2 \cdot 1 \rrbracket_2$:

- ▶ signature $\Sigma_{[2]} = \Sigma \times \{0, 1\}^2$.

Compilation: $X_1 = X_2 \cdot 1$

Automaton $\llbracket X_1 = X_2 \cdot 1 \rrbracket_2$:

- ▶ signature $\Sigma_{[2]} = \Sigma \times \{0, 1\}^2$.
- ▶ states: q_0, q_1, q_2
- ▶ final states: q_2
- ▶ transitions:

$$\langle a, 0, 0 \rangle (q_0, \dots, q_0) \rightarrow q_0$$

$$\langle a, 1, 0 \rangle (q_0, \dots, q_0) \rightarrow q_1$$

$$\langle a, 0, 1 \rangle (q_1, q_0, \dots, q_0) \rightarrow q_2$$

$$\langle a, 0, 0 \rangle (q_0, \dots, q_0, q_2, q_0, \dots, q_0) \rightarrow q_2$$

For $m \geq 2$,

$$\llbracket X_2 = X_1 \cdot 1 \rrbracket_m := \text{cyl}_{2,m}(\llbracket X_2 = X_1 \cdot 1 \rrbracket_2)$$

Compilation: $X_1 \subseteq L_a$

Automate $\llbracket X_1 \subseteq L_a \rrbracket_1$:

- ▶ signature $\Sigma_{[2]} = \Sigma \times \{0, 1\}^2$.

Compilation: $X_1 \subseteq L_a$

Automate $\llbracket X_1 \subseteq L_a \rrbracket_1$:

- ▶ signature $\Sigma_{[2]} = \Sigma \times \{0, 1\}^2$.
- ▶ states: q_0
- ▶ final states: q_0
- ▶ transitions:

$$\langle a, 0 \rangle (q_0, \dots, q_0) \rightarrow q_0$$

$$\langle b, 0 \rangle (q_0, \dots, q_0) \rightarrow q_0 \quad (b \neq a)$$

$$\langle a, 1 \rangle (q_0, \dots, q_0) \rightarrow q_0$$

For $m \geq 1$,

$$\llbracket X_1 \subseteq L_a \rrbracket_m := \text{cyl}_{1,m}(\llbracket X_1 \subseteq L_a \rrbracket_1)$$

Compilation: Boolean Connectors

- ▶ $\llbracket \phi(X_1, \dots, X_n) \vee \phi(X_1, \dots, X_{n'}) \rrbracket_m :=$
 $\llbracket \phi(X_1, \dots, X_n) \rrbracket_m \cup \llbracket \phi(X_1, \dots, X_{n'}) \rrbracket_m$
with $m \geq \max(n, n')$
- ▶ $\llbracket \phi(X_1, \dots, X_n) \wedge \phi(X_1, \dots, X_{n'}) \rrbracket_m :=$
 $\llbracket \phi(X_1, \dots, X_n) \rrbracket_m \cap \llbracket \phi(X_1, \dots, X_{n'}) \rrbracket_m$
with $m \geq \max(n, n')$
- ▶ $\llbracket \neg \phi(X_1, \dots, X_n) \rrbracket_m := \mathcal{T}(\Sigma_{[m]}) \setminus \llbracket \phi(X_1, \dots, X_n) \rrbracket_m$
for $m \geq n$.

Compilation: Quantifiers

- ▶ $\llbracket \exists X_{n+1} \phi(X_1, \dots, X_{n+1}) \rrbracket_n := \text{proj}_n(\llbracket \phi(X_1, \dots, X_{n+1}) \rrbracket_{n+1})$
- ▶ NB: this construction does **not** preserve **determinism**.
- ▶ $\llbracket \exists X_{n+1} \phi(X_1, \dots, X_{n+1}) \rrbracket_m := \text{cyl}_{n,m}(\llbracket \exists X_{n+1} \phi(X_1, \dots, X_{n+1}) \rrbracket_n)$ for $m \geq n$.
- ▶ $\forall = \neg \exists \neg$

Thatcher & Wright's Theorem

Theorem :

For all formula $\phi \in WSkS_0$ over Σ without free variables, there exists a tree automaton \mathcal{A}_ϕ over Σ , such that $L(\mathcal{A}_\phi) = L(\phi)$.

$\mathcal{A}_\phi = \llbracket \phi \rrbracket_0$ can be computed **explicitly!**

Corollary :

For all formula $\phi \in WSkS$ over Σ without free variables there exists a tree automaton \mathcal{A}_ϕ over Σ , such that $L(\mathcal{A}_\phi) = L(\phi)$.

using translation of $WSkS$ into $WSkS_0$ first.

Size of \mathcal{A}_ϕ

Theorem : Stockmeyer and Meyer 1973

For all n there exists $\exists x_1 \neg \exists y_1 \exists x_2 \neg \exists y_2 \dots \exists x_n \neg \exists y_n \phi \in \text{FOL}$ such that for every automaton \mathcal{A} recognizing the same language

$$\text{size}(\mathcal{A}) \geq \left. 2^{2^{\dots^{2^{\text{size}(\phi)}}}} \right\} n$$

Plan

WS k S: Definition

Automata \rightarrow Logic

Logic \rightarrow Automata

Fragments and Extensions of WS k S

WSkS and FO

Using the 2 directions of the Thatcher & Wright theorem:

$$\text{WSkS} \ni \phi \mapsto \mathcal{A} \mapsto \exists Y_1 \dots \exists Y_n \psi$$

with $\psi \in \text{FOL}$.

Corollary :

Every WSkS formula is equivalent to a formula $\exists Y_1 \dots \exists Y_n \psi$ with ψ first order.

Proposition :

The language L of terms with an even number of nodes labeled by a is regular (hence WSkS-definable) but not FO-definable.

pr.: with Ehrenfeucht-Fraïssé games.

Ehrenfeucht-Fraïssé Games

goal: prove FO equivalence of finite structures
(wrt finite set of predicates \mathcal{L}).

Definition

for two finite \mathcal{L} -structures \mathfrak{A} and \mathfrak{B} $\mathfrak{A} \equiv_m \mathfrak{B}$ iff for all ϕ closed, of quantifier depth m , $\mathfrak{A} \models \phi$ iff $\mathfrak{B} \models \phi$

Ehrenfeucht-Fraïssé Games

$\mathcal{G}_m(\mathfrak{A}, \mathfrak{B})$

1 Spoiler chooses $a_1 \in \text{dom}(\mathfrak{A})$ or $b_1 \in \text{dom}(\mathfrak{B})$

1' Duplicator chooses $b_1 \in \text{dom}(\mathfrak{B})$ or $a_1 \in \text{dom}(\mathfrak{A})$

⋮

m' Duplicator chooses $b_m \in \text{dom}(\mathfrak{B})$ or $a_m \in \text{dom}(\mathfrak{A})$

Duplicator wins if $\{a_1 \mapsto b_1, \dots, a_m \mapsto b_m\}$ is an injective partial function compatible with the relations of \mathfrak{A} and \mathfrak{B} ($\forall P \in \mathcal{P}$, $P^{\mathfrak{A}}(a_{i_1}, \dots, a_{i_n})$ iff $P^{\mathfrak{B}}(b_{i_1}, \dots, b_{i_n})$)

= partial isomorphism.

Otherwise Spoiler wins.

Theorem : Ehrenfeucht-Fraïssé

$\mathfrak{A} \equiv_m \mathfrak{B}$ iff Duplicator has a winning strategy for $\mathcal{G}_m(\mathfrak{A}, \mathfrak{B})$.

Ehrenfeucht-Fraïssé Theorem

more generally: equivalence of finite structures + valuation of n free variables.

for two finite \mathcal{L} -structures \mathfrak{A} and \mathfrak{B} and
 $\alpha_1, \dots, \alpha_n \in \text{dom}(\mathfrak{A}), \beta_1, \dots, \beta_n \in \text{dom}(\mathfrak{B}), m \geq 0,$

$$\mathfrak{A}, \alpha_1, \dots, \alpha_n \equiv_m \mathfrak{B}, \beta_1, \dots, \beta_n$$

iff for all $\phi(x_1, \dots, x_n)$ of quantifier depth $m,$

$$\mathfrak{A}, \sigma_a \models \phi(\bar{x}) \text{ iff } \mathfrak{B}, \sigma_b \models \phi(\bar{x})$$

where $\sigma_a = \{x_1 \mapsto \alpha_1, \dots, x_n \mapsto \alpha_n\},$
 $\sigma_b = \{x_1 \mapsto \beta_1, \dots, x_n \mapsto \beta_n\}.$

Games: the partial isomorphisms must extend
 $\{\alpha_1 \mapsto \beta_1, \dots, \alpha_n \mapsto \beta_n\}.$

FO $\not\subseteq$ WSKS

let $\Sigma = \{a : 1, \perp : 0\}$.

Lemma :

For all $m \geq 3$ and all $i, j \geq 2^m - 1$,
Duplicator has a winning strategy for $\mathcal{G}_m(a^i(\perp), a^j(\perp))$.

Corollary :

The language $L \subseteq \mathcal{T}(\Sigma)$ of terms with an even number of nodes labeled by a is not FO-definable.

- ▶ Star-free languages = FO definable holds for words [McNaughton Papert] but not for trees.
- ▶ It is an active field of research to characterize regular tree languages definable in FO.
e.g. [Benedikt Segoufin 05] \approx locally threshold testable.

Restriction to Antichains

Definition :

An **antichain** is a subset $P \subseteq \mathcal{P}os(t)$ s.t. $\forall p, p' \in P$,
 $p \not\prec p'$ and $p' \not\prec p$.

antichain-WSkS: second-order quantifications are restricted to antichains.

Theorem :

If $\Sigma_1 = \emptyset$, the classes of antichain-WSkS languages and regular languages over Σ coincide.

Theorem :

chain-WSkS is strictly weaker than WSkS.

MSO on Graphs

Weak second-order monadic theory of the grid

Σ finite alphabet,

$$\mathcal{L}_{\text{grid}} := \{=, S_{\rightarrow}, S_{\uparrow}, L_a \mid a \in \Sigma\}$$

Grid $G : \mathbb{N} \times \mathbb{N} \rightarrow \Sigma$; Interpretation structure:

$$\underline{G} := \langle \mathbb{N} \times \mathbb{N}, =, x + 1, y + 1, L_a^G, L_b^G, \dots \rangle.$$

Proposition :

The weak monadic second-order theory of the grid is undecidable.

csq: weak MSO of graphs is undecidable.

MSO on Graphs (Remarks)

- ▶ algebraic framework [Courcelle]:
MSO decidable on graphs generated by a hedge replacement graph grammar = least solutions of equational systems based on graph operations: $\parallel : 2$, $exch_{i,j} : 1$, $forget_i : 1$, $edge : 0$, $ver : 0$.
- ▶ related notion: graphs with bounded *tree width*.
- ▶ FO-definable sets of graphs of bounded degree = locally threshold testable graphs (some local neighborhood appears n times with $n < \text{threshold} - \text{fixed}$).

Undecidable Extensions

Left concatenation: new predicate

$$S'_1 = \{ \langle p, 1 \cdot p \rangle \mid p, 1 \cdot p \in \mathcal{Pos}(t) \}$$

Proposition :

WS2S + left concatenation predicate is undecidable.

Predicate of equal length.

Proposition :

WS2S + $|x| = |y|$ is undecidable.

MONA

[Klarlund et al 01]

<http://www.brics.dk/mona/>

- ▶ decision procedures for WS1S and WS2S
- ▶ by translation of formulas into automata

Part III

Automata for Unranked Trees

Plan

Unranked Trees and Reasoning Tasks over XML Documents XML Processing

Automata for Unranked Ordered Trees

Automata for Unranked Unordered Trees

Automata for Unranked Mixed Trees

Regular Languages modulo Associativity and Commutativity

Verification of XML Updates

Automata for

- ▶ **ranked terms** = first order terms over a signature
 - every symbols has a fixed arity
- functional program analysis
- ▶ **unranked terms** = finite trees (directed, rooted) labelled over a finite alphabet
 - one node can have arbitrarily (though finitely) many childrens
 - the number of children of a node does not depend on its label
- Web data, regular term languages modulo A and AC

A Brief History of Tree Automata

- ▶ 60's 70's, logic for computer science
[Thatcher 67], [Takachi 75]: unranked labeled trees
- ▶ end of 80's: application to automated deduction
[Dauchet Tison et al] (ranked trees = terms)
- ▶ 90's feature trees over infinite set of features: unranked trees
[Smolka 92] – applications to computational linguistic
[Blackburn 94], [Carpenter 92].
- ▶ 2000 and later: XML processing - unranked trees
[Vianu CSL 01], [Schwentick 07].

Imperative Program

[Bouajjani Touili CAV 02]

void X() {	X	→	Y · X	(r ₁)
while(true) {	Y	→	t	(r ₂)
if Y() {	Y	→	f	(r ₃)
thread_create(&t1,Z)	t · X	→	X Z	(r ₄)
} else { return }	f	→	0	(r ₅)
}				
}				

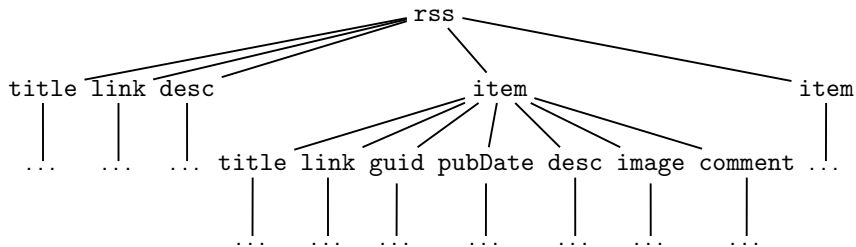
Reachability analysis:

- ▶ The set of reachable terms is **regular**, but
 - ▶ we want \cdot Associative,
 - ▶ and \parallel Associative and Commutative,
 - ▶ and regular term languages are not closed modulo A and AC.
- consider unranked trees as representative.

Web data (XML Document)

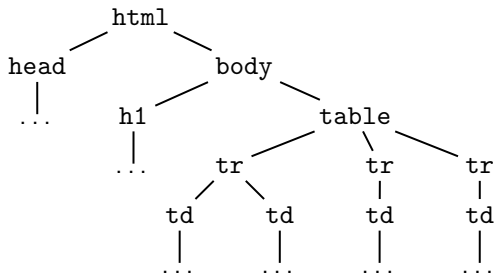
```
<rss version="2.0">
  <title>My blog</title>
  <link>http://myblog.blogspot.com</link>
  <description>bla bla bla</description>
  <item>
    <title>Concert</title>
    <link>http://myblog.blogspot.com/me/Mon blog/...</link>
    <guid>5f7da0aa-a593-4a2e</guid>
    <pubDate>Fri, 21 Mar 2009 14:40:02 +0100</pubDate>
    <description>...</description>
    <image href="..."></image>
    <comment link="..." count="0" enabled="0">...</comment>
  </item>
  <item>
    <title>Journée de surf</title>
    ...
  </item>
</rss>
```

Web Data



HTML Document

```
<html>
  <head>...</head>
  <body>
    <h1>...</h1>
    <table>
      <tr>
        <td>...</td>
        <td>...</td>
      </tr>
      <tr>
        <td>...</td>
      </tr>
      <tr>
        <td>...</td>
      </tr>
    </table>
  </body>
</html>
```



XML Document Types

documents = unranked trees

conformity / validation

- ▶ class of documents with a predefined structure (valid documents)

```
<table>
```

```
  <tr> <td> c11 </td> <td> c12 </td> </tr>
```

```
  <tr> <td> c21 </td> <td> c22 </td> </tr>
```

```
</table>
```

- ▶ defined by a schema (DTD, XML schema, Relax NG...) = tree language
- ▶ All the schema formalisms in use currently correspond to tree automata [Schwentick JCSS 07] [Murata et al 05].

Reasoning Tasks over XML Documents

- ▶ type definitions (DTD, XML schema, Relax NG...) \subseteq automata
 - validation = membership problem
 - schema entailment = inclusion problems
- ▶ querying
 - query satisfiability = emptiness problem
- ▶ integrity constraints (keys, inclusion):
 - consistency [Fan Libkin JACM 02]
- ▶ unranked tree transformations
 - type checking
 - given L_{in} , regular input language
 - T transformation (in XSLT...)
 - L_{out} regular output language
 - question: do we have $T(L_{in}) \subseteq L_{out}$?

Plan

Unranked Trees and Reasoning Tasks over XML Documents

Automata for Unranked Ordered Trees

- Unranked Ordered Trees

- Hedge Automata, Determinism

- Decision Problems

- Binary Encoding

- Boolean Closure

- Minimization

Automata for Unranked Unordered Trees

Automata for Unranked Mixed Trees

Regular Languages modulo Associativity and Commutativity

Modifications to XML Models

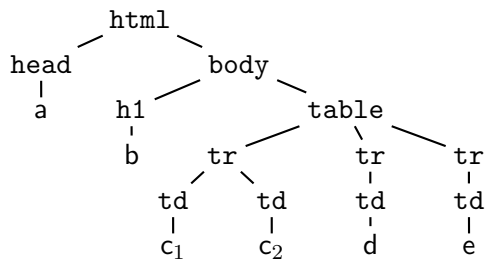
Unranked Ordered Trees

- ▶ Σ is a finite alphabet.

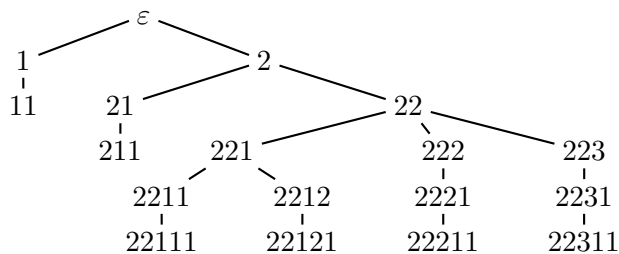
$$\begin{aligned}\text{tree} &:= a(\text{hedge}) \quad (a \in \Sigma) \\ \text{hedge} &:= \text{tree}^*\end{aligned}$$

- ▶ a hedge can be empty. $a()$ is denoted by a .
- ▶ The set of all unranked ordered trees over Σ is denoted $\mathcal{O}(\Sigma)$.
- ▶ The set of hedges over Σ is denoted $\mathcal{H}(\Sigma)$.
- ▶ set of positions $\subset \mathbb{N}^*$: as for terms.

Example: Tree of $\mathcal{O}(\Sigma)$



positions:



Example: Language $\subseteq \mathcal{O}(\Sigma)$

- ▶ $\Sigma = \{a, b\}$.
- ▶ $L :=$ terms of $\mathcal{O}(\Sigma)$
 - ▶ height at most 1,
 - ▶ root is labelled by a ,
 - ▶ even number of leaves, all leaves labelled by b .
- ▶ $L = \{a, a(bb), a(bbbb), \dots\}$.
- ▶ **finite description:** $L = a((bb)^*)$

Hedge Automata (HA)

Definition : Hedge Automata

A **Hedge Automaton** (HA) over an alphabet Σ is a tuple $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$ where Q is a finite set of **states**, $Q^f \subseteq Q$ is the subset of final states and Δ is a set of transition rules of the form: $a(L) \rightarrow q$ with $a \in \Sigma$ and $L \subseteq Q^*$ is a regular language.

A **run** of \mathcal{A} on $t \in \mathcal{O}(\Sigma)$ is a tree $r \in \mathcal{O}(Q)$ such that

- ▶ r and t have the same domain,
- ▶ for all $p \in \mathcal{Pos}(t)$, with $t(p) = a$, $r(p) = q$, there exists $a(L) \rightarrow q \in \Delta$ such that $r(p1) \dots r(pn) \in L$, where n is the number of successors of p in $\mathcal{Pos}(t)$.

The run r is **accepting** (**successful**) iff $r(\varepsilon) \in Q^f$.

HA Languages

- ▶ language of \mathcal{A} : $L(\mathcal{A})$ is the set of terms on which there exists an accepting run of \mathcal{A} ,
- ▶ language of \mathcal{A} in state $q \in Q$: $L(\mathcal{A}, q)$ is the set of terms t such that there exists a run r of \mathcal{A} on t with $r(\varepsilon) = q$,
- ▶ $L(\mathcal{A}) = \bigcup_{q \in Q_f} L(\mathcal{A}, q)$.
- ▶ equivalently, $L(\mathcal{A}, q)$ is the smallest set of terms $a(t_1, \dots, t_n) \in \mathcal{O}(\Sigma)$ ($n \geq 0$) such that there exists a transition $a(L) \rightarrow q$, and states $q_1, \dots, q_n \in Q$ with $t_i \in L(\mathcal{A}, q_i)$ s.t. $i \leq n$ and $q_1 \dots q_n \in L$.

HA Language: Example 1

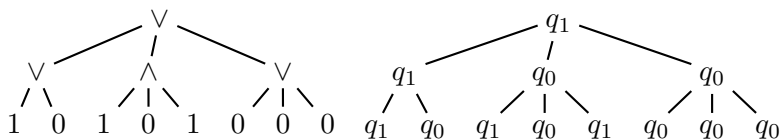
- ▶ $\Sigma = \{a, b\}$.
- ▶ $L :=$ terms of $\mathcal{O}(\Sigma)$
 - ▶ height 1,
 - ▶ root is labelled by a ,
 - ▶ even number of leaves, all leaves labelled by b .
- ▶ $L = \{a, a(bb), a(bbbb), \dots\}$.
- ▶ **finite description:** $L = L(\mathcal{A})$ with
 $\mathcal{A} := (\Sigma, \{q_a, q_b\}, \{q_a\}, \{b \rightarrow q_b, a((q_bq_b)^*) \rightarrow q_a\})$.

Boolean Expressions with Variadic \wedge and \vee

- ▶ $\Sigma = \{\wedge, \vee, 0, 1\}$,
- ▶ states $\{q_0, q_1\}$,
- ▶ transitions:

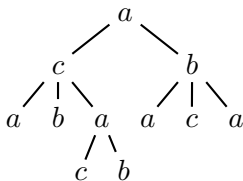
$$\begin{array}{ll} 0 & \rightarrow q_0 \\ \wedge(q_1^* q_0 (q_0 \mid q_1)^*) & \rightarrow q_0 \\ \vee(q_0 q_0^*) & \rightarrow q_0 \\ \neg(q_0) & \rightarrow q_1 \end{array} \qquad \begin{array}{ll} 1 & \rightarrow q_1 \\ \wedge(q_1 q_1^*) & \rightarrow q_1 \\ \vee(q_0^* q_1 (q_0 \mid q_1)^*) & \rightarrow q_1 \\ \neg(q_1) & \rightarrow q_0 \end{array}$$

- ▶ example: Boolean expression and associated run



HA Language: Example 3

- ▶ $\Sigma = \{a, b, c\}$.
- ▶ $L :=$ terms of $\mathcal{O}(\Sigma)$
 - ▶ with 2 b 's at positions p_1 and p_2 , and
 - ▶ one c on the smallest common ancestor of p_1 and p_2 .

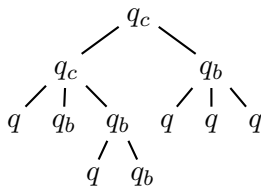
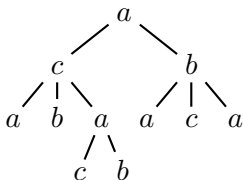


HA Language: Example 3

- ▶ $\Sigma = \{a, b, c\}$.
- ▶ $L :=$ terms of $\mathcal{O}(\Sigma)$
 - ▶ with 2 b 's at positions p_1 and p_2 , and
 - ▶ one c on the smallest common ancestor of p_1 and p_2 .

$\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$, avec $Q = \{q, q_b, q_c\}$, $Q^f = \{q_c\}$, $\Delta =$

$$\begin{array}{llll} a(Q^*) & \rightarrow & q & a(Q^*q_bQ^*) & \rightarrow & q_b & a(Q^*q_cQ^*) & \rightarrow & q_c \\ b(Q^*) & \rightarrow & q_b & c(Q^*q_bQ^*) & \rightarrow & q_b & b(Q^*q_cQ^*) & \rightarrow & q_c \\ c(Q^*) & \rightarrow & q & c(Q^*q_bQ^*q_bQ^*) & \rightarrow & q_c & c(Q^*q_cQ^*) & \rightarrow & q_c \end{array}$$



Normalized Hedge Automata

Definition :

A HA $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$ over Σ is called **normalized** if for all $a \in \Sigma$ and $q \in Q$, there is at most one transition of the form $a(L) \rightarrow q$ in Δ .

When \mathcal{A} is normalized, we denote $a(L_{a,q}) \rightarrow q$ the unique transition with a and q .

Proposition :

For all HA \mathcal{A} , there exists a normalized HA \mathcal{A}_n recognizing the same language.

The size of \mathcal{A}_n is linear in the size of \mathcal{A} .

Complete and Deterministic Hedge Automata

Semantical definitions

Definition :

A HA $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$ over Σ is **complete** if for all $t \in \mathcal{O}(\Sigma)$, there exists at least one state $q \in Q$ s.t. $t \in L(\mathcal{A}, q)$.

Definition :

A HA $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$ over Σ is **deterministic** if for all $t \in \mathcal{O}(\Sigma)$, there exists at most one state $q \in Q$ s.t. $t \in L(\mathcal{A}, q)$.

Complete and Deterministic Hedge Automata

Syntactical definitions

Definition :

A HA $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$ over Σ is **complete** if for all $a \in \Sigma$ and all finite sequence $q_1, \dots, q_n \in Q^*$, there exists a transition $a(L) \rightarrow q \in \Delta$ with $q_1 \dots q_n \in L$.

Definition :

A HA $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$ over Σ is **deterministic** if for all transitions $a(L_1) \rightarrow q_1$ and $a(L_2) \rightarrow q_2$ in Δ , either $L_1 \cap L_2 = \emptyset$, or $q_1 = q_2$.

Determinism: Examples

HA for Boolean expressions evaluation : **deterministic**.

$\Sigma = \{\wedge, \vee, 0, 1\}$, $Q = \{q_0, q_1\}$ et Δ :

$$\begin{array}{ll} 0 \rightarrow q_0 & 1 \rightarrow q_1 \\ \wedge(q_1^* q_0 (q_0 \mid q_1)^*) \rightarrow q_0 & \wedge(q_1 q_1^*) \rightarrow q_1 \\ \vee(q_0 q_0^*) \rightarrow q_0 & \vee(q_0^* q_1 (q_0 \mid q_1)^*) \rightarrow q_1 \end{array}$$

Language with 2 b 's and common ancestor c : **not deterministic**.

$\Sigma = \{a, b, c\}$, $Q = \{q, q_b, q_c\}$, $Q^f = \{q_c\}$, Δ :

$$\begin{array}{lll} a(Q^*) \rightarrow q & a(Q^* q_b Q^*) \rightarrow q_b & a(Q^* q_c Q^*) \rightarrow q_c \\ b(Q^*) \rightarrow q_b & c(Q^* q_b Q^*) \rightarrow q_b & b(Q^* q_c Q^*) \rightarrow q_c \\ c(Q^*) \rightarrow q & c(Q^* q_b Q^* q_b Q^*) \rightarrow q_c & c(Q^* q_c Q^*) \rightarrow q_c \end{array}$$

HA Completion

Proposition :

For all HA \mathcal{A} , there exists a complete HA \mathcal{A}_c recognizing the same language.

The size of \mathcal{A}_c is linear is the size of \mathcal{A} .

HA Completion

Proposition :

For all HA \mathcal{A} , there exists a complete HA \mathcal{A}_c recognizing the same language.

The size of \mathcal{A}_c is linear in the size of \mathcal{A} .

pr.: add a *trash* state q_{\perp} and transitions :

$$a\left(\bigcap_{q \in Q} Q^* \setminus L_{a,q} \cup Q_{\perp}^* q_{\perp} Q_{\perp}^*\right) \rightarrow q_{\perp}$$

HA Determinization

Proposition :

For all HA \mathcal{A} , there exists a deterministic HA \mathcal{A}_d recognizing the same language.

The size of \mathcal{A}_d is exponential in the size of \mathcal{A} (lower bound).

pr.: subset construction

HA Determinization

Proposition :

For all HA \mathcal{A} , there exists a deterministic HA \mathcal{A}_d recognizing the same language.

The size of \mathcal{A}_d is exponential in the size of \mathcal{A} (lower bound).

pr.: subset construction for $\mathcal{A} = (Q, Q_f, \Delta)$ (normalised):

- ▶ $\mathcal{A}_d = (2^Q, Q^f, \Delta_d)$,
- ▶ $Q^f = \{S \subseteq Q \mid S \cap Q_f \neq \emptyset\}$
- ▶ $\Delta_d: a(L_{a,S}) \rightarrow S$ ($S \subseteq Q$)

$$L_{a,S} = \bigcap_{q \in S} S_{a,q} \setminus \bigcup_{q \notin S} S_{a,q}$$

with

$$S_{a,q} = \{S_1 \dots S_n \in Q_d^* \mid \exists q_1 \in S_1, \dots, \exists q_n \in S_n, q_1 \dots q_n \in L_{a,q}\}$$

HA: Membership Decision

Proposition : \in

The problem of membership is decidable in polynomial time for the HAs whose languages $L_{a,q}$ are given by NFAs.

- ▶ linear time for DHA whose languages $L_{a,q}$ are given by DFA,
- ▶ NP-complete if the languages $L_{a,q}$ are given by alternating automata.

HA: Emptiness Decision

Proposition : \emptyset

The problem of emptiness is decidable in polynomial time for HAs whose languages $L_{a,q}$ are given by NFAs.

pr.: state marking. $M \subseteq Q$,

- ▶ initially, $M = \emptyset$.
- ▶ at each step , if $a(L_{a,q}) \rightarrow q \in \Delta$ and $L_{a,q} \cap M^* \neq \emptyset$ then $M := M \cup \{q\}$.
- ▶ PSPACE-complete if the languages $L_{a,q}$ are given by alternating automata.

HA: Decision of Inclusion and Equivalence

Proposition : \subseteq, \equiv

The problem of inclusion (resp. equivalence) is EXPTIME-complete for HAs whose languages $L_{a,q}$ are given by NFAs.

pr.: $L(\mathcal{A}_1) \subseteq L(\mathcal{A}_2)$ iff $L(\mathcal{A}_1) \cap (\mathcal{O}(\Sigma) \setminus L(\mathcal{A}_2)) = \emptyset$.

- ▶ in PTIME for DHAs whose languages $L_{a,q}$ are given by DFAs.
- ▶ PSPACE-complete for DHAs whose languages $L_{a,q}$ are given by alternating automata.

Currying

Transformation into binary trees with @ and constants.

We associate the following signature to an alphabet Σ :

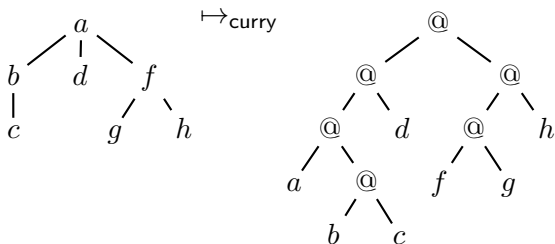
$$\Sigma_{@} := \{a : 0 \mid a \in \Sigma\} \cup \{@ : 2\}$$

The function $\text{curry} : \mathcal{O}(\Sigma) \rightarrow \mathcal{T}(\Sigma_{@})$, is defined recursively:

- ▶ $\text{curry}(a) := a$,
- ▶ $\text{curry}(a(t_1, \dots, t_n)) := @(\text{curry}(a(t_1, \dots, t_{n-1})), \text{curry}(t_n))$.

Currying: Example 1

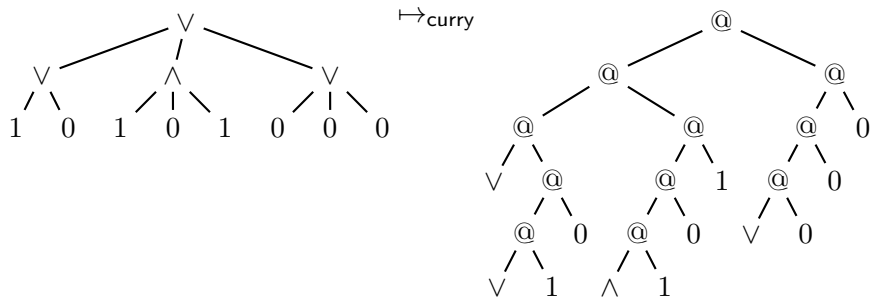
- ▶ $\text{curry}(a) = a,$
- ▶ $\text{curry}(a(t_1, \dots, t_n)) = @(\text{curry}(a(t_1, \dots, t_{n-1})), \text{curry}(t_n))$



Currying: Example 2

- ▶ $\text{curry}(a) = a,$
- ▶ $\text{curry}(a(t_1, \dots, t_n)) = @(\text{curry}(a(t_1, \dots, t_{n-1})), \text{curry}(t_n))$

transforming unranked Boolean expressions into binary.



Currying: Properties

Lemma :

curry is a bijection from $\mathcal{O}(\Sigma)$ into $\mathcal{T}(\Sigma_{@})$.

Proposition :

$L \subseteq \mathcal{O}(\Sigma)$ is a HA language iff $\text{curry}(L)$ is regular.

pr.:: Let $\mathcal{A} = \langle Q, Q^f, \Delta \rangle$ normalized HA, and $B_{a,q} = \langle P_{a,q}, \text{init}_{a,q}, F_{a,q}, R_{a,q} \rangle$ be a NFA recognizing $L_{a,q}$ f.a. $a \in \Sigma, q \in Q, a(L_{a,q}) \rightarrow q \in \Delta$.

$\text{curry}(\mathcal{A}) = \langle Q \cup \bigcup_{a,q} P_{a,q}, Q^f, \Delta' \rangle$ where Δ' contains:

$a \rightarrow q$ if $B_{a,q}$ recognizes the empty word,

$a \rightarrow \text{init}_{a,q}$ for all $q \in Q$,

$@(p, q) \rightarrow p'$ if there is a transition $p \xrightarrow{q} p'$ in some $B_{a',q'}$, and

$@(p, q) \rightarrow q'$ if there is a transition $p \xrightarrow{q} p'$ in some $B_{a',q'}$ and $p' \in F_{a',q'}$.

HA: Boolean Operations

Proposition :

The class of HA languages is closed under union, intersection and complement.

pr.: $\text{curry}(L_1 \cup L_2) = \text{curry}(L_1) \cup \text{curry}(L_2)$

Hence we can reuse the construction for ranked TA

$$L_1 \cup L_2 = \text{curry}^{-1}(\text{curry}(L_1) \cup \text{curry}(L_2))$$

HA: Closure under Morphisms

projection $h : \mathcal{O}(\Sigma) \rightarrow \mathcal{O}(\Sigma')$, defined by extension to trees of an application $h : \Sigma \rightarrow \Sigma'$.

$$h(L) = \{h(t) \mid t \in L\} \quad \text{and} \quad h^{-1}(L') = \{t \in \mathcal{O}(\Sigma) \mid h(t) \in L'\}$$

Proposition :

The class of HA languages is closed under projections and inverse projections.

pr.: it follows from the closure of regular ranked tree languages under linear morphisms and inverse morphisms.

Minimization

Definition of minimal deterministic HA:

2 questions must be addressed

1. for which definition of determinism?
(minimization makes sense only for deterministic automata)
2. what to minimize?

Minimization

For ranked tree automata, the answer to both questions is clear:

1. ranked DTA: every step of computation is deterministic
2. we want to minimize the number of states

For unranked tree automata, this is not so clear:

1. even for DHA[DFA] (DHA whose horizontal languages are defined by DFA), if we have $a(L) \rightarrow q$ and $a(L') \rightarrow q'$ and $L \cap L' = \emptyset$, in configuration $a(q_1 \dots q_n)$ we must test both $q_1 \dots q_n \in L$ and $q_1 \dots q_n \in L'$ before firing the right transition. Which one is tested first?
In the construction of a TA for $\text{curry}(\text{HA})$, we choose ND.
2. there are states for the DHA and for the DFAs for the horizontal languages.

Minimization of DHA

First approach: we ignore the formalism for horizontal languages, i.e. we chose

1. DFAs (whatever for the horizontal automata)
2. number of states of the DFA

Congruence of a language $L \subseteq \mathcal{O}(\Sigma)$:

$$s \equiv_L t \quad \text{iff} \quad \forall C \ C[s] \in L \Leftrightarrow C[t] \in L$$

Minimal DHA for the HA language L :

- ▶ states: $\{[t]_{\equiv_L} \mid t \in \mathcal{O}(\Sigma)\}$,
- ▶ final states: $\{[t]_{\equiv_L} \mid t \in L\}$ (we simply write $[t]$ below),
- ▶ transitions: $\{a(L_{a,[t]}) \rightarrow [t] \mid L_{a,[t]} = \{[t_1] \dots [t_n] \mid a(t_1 \dots t_n) \equiv_L t\}\}$,

Minimization of DHA (2)

2 drawbacks for the first approach

- ▶ the complexity of the effective construction depends on the formalism for horizontal languages.
- ▶ no analogous of Myhill-Nerode theorem for DFA or DTA (ranked):

L is an HA language $\Rightarrow \equiv_L$ has finite index

L is an HA language $\not\Leftarrow \equiv_L$ has finite index

Minimization of DHA[DTA]

Second approach: we consider both vertical and horizontal states and transitions, i.e. we chose

1. DFA[DTA] (DHA whose horizontal language are defined by disjoint DFAs)
2. number of states of the DFA + number of states of the horizontal (disjoint) DTAs

Minimization of DHA[DTA] (2)

first idea: use the curry encoding and minimize the ranked TA.

problem: the TA associated to an HA wrt curry is not deterministic; we have $a \rightarrow \text{init}_{a,q}$ for all $q \in Q$.

other question: uniqueness of minimal automaton?

→ stepwise automata [Niehren et al 04]:

one unique transition from each $a \in \Sigma$ to the start state of a deterministic machine that will read the state sequence below a and output a state.

Moreover, vertical states = horizontal states.

Deterministic Stepwise Automata

Definition : stepwise automata

A **deterministic stepwise hedge automaton** (DSHA) is a tuple $\mathcal{A} = (\Sigma, Q, Q_f, \delta_0, \delta)$, where Σ , Q , and Q_f are as usual, $\delta_0 : \Sigma \rightarrow Q$ is a function assigning to each letter of the alphabet an initial state, and $\delta : Q \times Q \rightarrow Q$ is the transition function.

$$\begin{aligned} \text{For } a \in \Sigma, \quad \delta_a : \quad & Q^* \rightarrow Q \\ \delta_a(\varepsilon) &= \delta_0(a) \\ \delta_a(w \cdot q) &= \delta(\delta_a(w), q) \end{aligned}$$

A **run** of \mathcal{A} on $t \in \mathcal{O}(\Sigma)$ is a tree $r \in \mathcal{O}(Q)$ such that

- ▶ r and t have the same domain,
- ▶ for all $p \in \mathcal{Pos}(t)$, $r(p) = \delta_{t(p)}(r(p1) \dots r(pn))$, where n is the number of successors of p in $\mathcal{Pos}(t)$.

The run r is **accepting** (**successful**) iff $r(\varepsilon) \in Q^f$.

Stepwise Automata & Ranked TA

stepwise DSHA \mathcal{A}	ranked DTA $\text{curry}(\mathcal{A})$
$\delta_0(a) = q$	$a \rightarrow q$
$\delta(q_1, q_2) = q$	$@(q_1, q_2) \rightarrow q$

Lemma :

For all $t, t' \in \mathcal{O}(\Sigma)$ and $q, q' \in Q$,
if $t \in L(\mathcal{A}, q)$ and $t' \in L(\mathcal{A}, q')$, then $t @ t' \in L(\mathcal{A}, \delta(q, q'))$.

Lemma :

For all DSHA \mathcal{A} , $\text{curry}(L(\mathcal{A})) = L(\text{curry}(\mathcal{A}))$.

Minimal Stepwise Automata

Corollary :

DSHA recognize all HA unranked tree languages.

Corollary :

For each HA language $L \subseteq \mathcal{O}(\Sigma)$ there is a unique (up to renaming of states) minimal DSHA accepting L .

Plan

Unranked Trees and Reasoning Tasks over XML Documents

Automata for Unranked Ordered Trees

Automata for Unranked Unordered Trees

- Unranked Unordered Trees

- Presburger Arithmetic

- Presburger Automata (PA)

- Determinism

- Boolean Closure

- Decision Problems

- Weak Second Order Monadic Logic PMSO

Automata for Unranked Mixed Trees

Regular Languages modulo Associativity and Commutativity

- ▶ previous part: unranked ordered trees
 - XML documents
 - hedge automata (HA)
 - see TATA book <http://tata.gforge.inria.fr> chapter 8
- ▶ this part: unranked **unordered** trees
 - web data
 - Presburger automata (PA)
 - see [Seidl, Schwentick, Muscholl. *Numerical Document Queries*. PODS 03]

Unranked Unordered Trees

- ▶ Σ is a finite alphabet.

$$\begin{aligned}\text{tree} &:= a(\text{multiset}) \quad (a \in \Sigma) \\ \text{multiset} &:= \{\text{tree}, \dots, \text{tree}\}\end{aligned}$$

- ▶ $a(\{t_1, \dots, t_n\})$ is denoted $a(t_1, \dots, t_n)$.
- ▶ rem: the multiset can be empty. $a(\emptyset)$ is denoted a .
- ▶ The set of unranked unordered trees over Σ is denoted $\mathcal{U}(\Sigma)$.

Examples of Languages of Trees of $\mathcal{U}(\Sigma)$

- ▶ $\Sigma = \{a, b\}$.
- ▶ terms of $\mathcal{U}(\Sigma)$
 - ▶ of height 1,
 - ▶ with one b at the root,
 - ▶ the leaves are a or b :
 - i.* with an even number of a ($\in \text{HA}$),
 - ii.* with the same number of a than b ($\notin \text{HA}$).

Presburger Arithmetic

Presburger Formulae:

$$\begin{array}{l} \text{term} ::= \\ \quad \left| \begin{array}{l} x \text{ (first order variables)} \\ n \text{ (natural number)} \\ \text{term} + \text{term} \end{array} \right. \\ \\ \text{form} ::= \\ \quad \left| \begin{array}{l} \text{term} = \text{term} \\ \neg \text{form} \mid \text{form} \vee \text{form} \mid \text{form} \wedge \text{form} \\ \forall x \text{ form} \mid \exists x \text{ form} \end{array} \right. \end{array}$$

Interpretation in the domain of natural numbers.

$(n_1, \dots, n_p) \models \phi(x_1, \dots, x_p)$ (x_1, \dots, x_p free variables)

iff $\phi(n_1, \dots, n_p)$ is evaluated to *true*.

Presburger Arithmetic

- ▶ notation: terms nx for $\underbrace{x + \dots + x}_n$,
- ▶ the natural can be restricted to 0 and 1,
- ▶ the atoms can be restricted to $x = n$ and $x = y + z$.

Examples:

- ▶ $x \leq y$

Presburger Arithmetic

- ▶ notation: terms nx for $\underbrace{x + \dots + x}_n$,
- ▶ the natural can be restricted to 0 and 1,
- ▶ the atoms can be restricted to $x = n$ and $x = y + z$.

Examples:

- ▶ $x \leq y \equiv \exists x' y = x + x'$.
- ▶ $odd(x)$

Presburger Arithmetic

- ▶ notation: terms nx for $\underbrace{x + \dots + x}_n$,
- ▶ the natural can be restricted to 0 and 1,
- ▶ the atoms can be restricted to $x = n$ and $x = y + z$.

Examples:

- ▶ $x \leq y \equiv \exists x' y = x + x'$.
- ▶ $odd(x) \equiv \exists y x = y + y + 1$.

Presburger Arithmetic

Theorem : Presburger Arithmetic

Presburger Arithmetic is decidable.

- ▶ lower bound 2-EXPTIME (Fischer and Rabin 1974)
- ▶ upper bound 3-EXPTIME (Klaedtke 2004, with an automata construction)
- ▶ NP-complete for the existential fragment.

Presburger Arithmetic

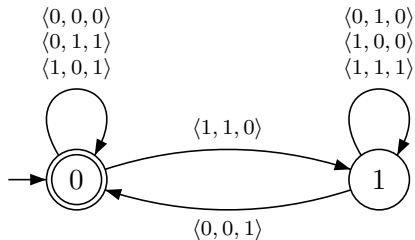
Decidability of Presburger Arithmetic.

We can associate to a formula $\phi(x_1, \dots, x_p)$ a finite automaton over the alphabet $\{0, 1\}^p$ recognizing the set of $\langle b_{1,1}, \dots, b_{p,1} \rangle \dots \langle b_{1,k}, \dots, b_{p,k} \rangle$ such that $b_{1,1} \dots b_{1,k}, \dots, b_{p,1} \dots b_{p,k}$ are the binary representations of integers n_1, \dots, n_p satisfying ϕ .

Hence we can decide whether there exists n_1, \dots, n_p such that $(n_1, \dots, n_p) \models \phi(x_1, \dots, x_p)$.

Presburger Arithmetic and Automata

finite automaton for $x_1 + x_2 = x_3$



Parikh Projection

$$\Sigma = \{a_1, \dots, a_p\}.$$

Definition : Parikh Projection

The Parikh projection of a word $w \in \Sigma^*$ is the tuple $\#(w) := (m_1, \dots, m_p)$ where m_i ($i \leq p$) is the number of occurrences of a_i in w .

For a set $L \subseteq \Sigma^*$, we denote $\#(L) := \{\#(w) \mid w \in L\}$.

Theorem :

For all context-free language $L \subseteq \Sigma^*$, there exists a Presburger formula $\phi(x_1, \dots, x_p)$ such that $\#(L) := \{(n_1, \dots, n_p) \mid \phi(x_1, \dots, x_p)\}$.

When L is regular, the Presburger formula is computed in linear time (in the size of the NFA defining L).

Theorem :

For all Presburger formula $\phi(x_1, \dots, x_p)$, one can build a NFA \mathcal{A} such that $\#(L(\mathcal{A})) = \{(n_1, \dots, n_k) \mid \phi(x_1, \dots, x_k)\}$.

Semi-linear sets

Definition :

- ▶ A *linear* set is a subset of \mathbb{N}^p of the form
$$\{\overline{v}_0 + \overline{v}_1 + \dots + \overline{v}_m \mid m \geq 0, \overline{v}_1, \dots, \overline{v}_m \in B\},$$

for $\overline{v}_0 \in \mathbb{N}^p$ and $B \subset \mathbb{N}^p$ finite (fixed).
- ▶ A *semi-linear* set is a finite union of linear sets.

Theorem : Parikh

Models of Presburger formulae \equiv semi-linear sets.

Presburger Automata (PA)

Definition : Presburger Automata

A *Presburger Automaton* (PA) over an alphabet Σ is a tuple $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$ where $Q = \{q_1, \dots, q_p\}$ is a finite set of *states*, $Q^f \subseteq Q$ is the subset of final states and Δ is a set of transition rules of the form: $a(\phi) \rightarrow q$ with $a \in \Sigma$, $q \in Q$, and $\phi = \phi(x_1, \dots, x_p)$ is a Presburger formula with one free variable x_i for each state q_i .

The language of \mathcal{A} in state $q \in Q$, denoted $L(\mathcal{A}, q)$, is the smallest subset of terms $a(t_1, \dots, t_n) \in \mathcal{U}(\Sigma)$ such that

- ▶ there exists $i_1, \dots, i_n \leq p$ such that for all $j \leq n$, $t_j \in L(\mathcal{A}, q_{i_j})$,
- ▶ there exists a transition $a(\phi) \rightarrow q \in \Delta$ such that $\#(q_{i_1}, \dots, q_{i_n}) \models \phi(x_1, \dots, x_p)$.

The language of \mathcal{A} is $L(\mathcal{A}) = \bigcup_{q \in Q^f} L(\mathcal{A}, q)$.

PA: Example 1

$$\Sigma = \{a, b, f\}.$$

Set of trees of $\mathcal{U}(\Sigma)$ where all the a and b label the leaves:

PA: Example 1

$$\Sigma = \{a, b, f\}.$$

Set of trees of $\mathcal{U}(\Sigma)$ where all the a and b label the leaves:

$$\mathcal{A} = (\{q\}, \{q\}, \{a(x_q = 0) \rightarrow q, b(x_q = 0) \rightarrow q, f(true) \rightarrow q\})$$

PA: Example 2

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

Set of trees of $\mathcal{U}(\Sigma)$ with the same number of a and of b under each node.

PA: Example 2

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

Set of trees of $\mathcal{U}(\Sigma)$ with the same number of a and of b under each node.

$$\mathcal{A} = (\{q_a, q_b, q\}, \{q_a, q_b, q\}, \{a(\phi) \rightarrow q_a, b(\phi) \rightarrow q_b, c(\phi) \rightarrow q\})$$

with $\phi \equiv x_{q_a} = x_{q_b}$.

PA: Example 3

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

Set of trees of $\mathcal{U}(\Sigma)$ where all internal nodes are labeled by a and every node labeled by a has at least as many sons without b than sons containing b .

PA: Example 3

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

Set of trees of $\mathcal{U}(\Sigma)$ where all internal nodes are labeled by a and every node labeled by a has at least as many sons without b as sons containing b .

$$Q = Q_f = \{q_a, q_b\}.$$

The state q_b accepts the trees containing a b and q_a accepts the others.

$$\Delta = \left\{ \begin{array}{ll} a(x_{q_a} \geq x_{q_b} = 0) & \rightarrow q_a \\ a(x_{q_a} \geq x_{q_b} > 0) & \rightarrow q_b \\ b(x_{q_a} = x_{q_b} = 0) & \rightarrow q_b \end{array} \right\}$$

Normalized Presburger Automata

Definition :

A PA $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$ sur Σ is called *normalized* if for all $a \in \Sigma$ and $q \in Q$, there is at most one transition of the form $a(\phi) \rightarrow q$ in Δ .

When \mathcal{A} is normalized, we note $a(\phi_{a,q}) \rightarrow q$ the unique transition with a and q .

Proposition :

For all PA \mathcal{A} , there exists a normalized PA \mathcal{A}_n recognizing the same language.

The size of \mathcal{A}_n is linear in the size of \mathcal{A} .

Complete Presburger Automata

Definition : Complete PA

A PA $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$ over Σ is *complete* if for all $a \in \Sigma$ and all $q_{i_1}, \dots, q_{i_n} \in Q^*$ there exists at least one transition $a(\phi) \rightarrow q \in \Delta$ such that $\#(q_{i_1}, \dots, q_{i_n}) \models \phi$.

Proposition : Completion

For all PA \mathcal{A} , there exists a complete PA \mathcal{A}_c recognizing the same language.

The size of \mathcal{A}_c is linear in the size of \mathcal{A} .

Complete Presburger Automata

Definition : Complete PA

A PA $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$ over Σ is *complete* if for all $a \in \Sigma$ and all $q_{i_1}, \dots, q_{i_n} \in Q^*$ there exists at least one transition $a(\phi) \rightarrow q \in \Delta$ such that $\#(q_{i_1}, \dots, q_{i_n}) \models \phi$.

Proposition : Completion

For all PA \mathcal{A} , there exists a complete PA \mathcal{A}_c recognizing the same language.

The size of \mathcal{A}_c is linear in the size of \mathcal{A} .

pr.: Let $\mathcal{A} = (Q, Q_f, \Delta)$ be a PA (normalized).

We add the state q_\perp : $\mathcal{A}_c = (Q \cup \{q_\perp\}, Q_f, \Delta_c)$ with

$$\begin{aligned} \Delta_c &:= a(\phi_{a,q} \wedge x_{q_\perp} = 0) \rightarrow q \text{ s.t. } a(\phi_{a,q}) \rightarrow q \in \Delta \\ &\cup a\left(\bigwedge_{q \in Q} \neg \phi_{a,q} \vee x_{q_\perp} > 0\right) \rightarrow q_\perp \end{aligned}$$

Deterministic Presburger Automata

Definition : Deterministic PA

A PA $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$ over Σ is *deterministic* if for all $a \in \Sigma$ and all $q_{i_1}, \dots, q_{i_n} \in Q^*$, if $a(\phi) \rightarrow q \in \Delta$ and $a(\phi') \rightarrow q' \in \Delta$ are such that $\#(q_{i_1}, \dots, q_{i_n}) \models \phi$ and $\#(q_{i_1}, \dots, q_{i_n}) \models \phi'$, then $q = q'$.

Lemma :

The determinism of PA is decidable.

pr.: Determinism is expressed the Presburger formula

$$\bigwedge \quad \neg(\phi \wedge \phi')$$
$$a(\phi) \rightarrow q$$
$$a(\phi') \rightarrow q'$$
$$q \neq q'$$

Determinization of PA

Proposition : Determinization

For all PA \mathcal{A} , there exists a deterministic PA \mathcal{A}_d recognizing the same language.

The size of \mathcal{A}_d is exponential in the size of \mathcal{A} (lower bound).

pr.: Let $\mathcal{A} = (Q, Q_f, \Delta)$, normalized with $Q = \{q_1, \dots, q_b\}$.

$$\mathcal{A}_d = (2^Q, \{S \subseteq Q \mid S \cap Q_f \neq \emptyset\}, \Delta_d)$$

Presburger formulae in Δ_d : a free variable x_S for each $S \subseteq Q$.

$$a\left(\bigwedge_{q \in S} \psi_{a,q} \wedge \bigwedge_{q \notin S} \neg \psi_{a,q}\right) \rightarrow S$$

with

$$\psi_{a,q} \equiv \exists_{p \in Q} x_p \quad \exists_{\substack{P \subseteq Q \\ p \in P}} x_{P,p} \phi_{a,q} \wedge \bigwedge_{p \in Q} x_p = \sum_{\substack{P \subseteq Q \\ p \in P}} x_{P,p} \wedge \bigwedge_{P \subseteq Q} x_P = \sum_{p \in P} x_{P,p}$$

ε -Transitions and Elimination

Remark :

PA with transitions $q \rightarrow q' \equiv$ PA.

PA: Boolean Operations

Proposition :

The class of languages of PA is closed under intersection and complementation.

- ∪ disjoint union (linear) or product (quadratic, preserves determinism).
- ∩ de Morgan law or product (quadratic).
- ¬ complete, determinize, invert final states (exponential).

PA: Membership Decision

Proposition : \in

The problem of membership is decidable for PA.

- ▶ in polynomial time for DPA,
- ▶ NP-complete for NPA.

PA: Emptiness Decision

Proposition : \emptyset

The problem of emptiness is decidable in polynomial time for PA.

pr.: states marking, construction of $M_i \subseteq Q$.

▶ initially, $M_0 = \emptyset$.

▶ at each step, if for $q \in Q \setminus M_i$, $\bigwedge_{p \in Q \setminus M_i} x_p = 0 \wedge \bigvee_{a \in \Sigma} \phi_{a,q}$ is satisfiable, then $M_{i+1} := M_i \cup \{q\}$.

PA: Decision of Inclusion, Equivalence

Proposition : \subseteq, \equiv

The problems of inclusion and equivalence are decidable for PA.

pr.: $L(\mathcal{A}_1) \subseteq L(\mathcal{A}_2)$ iff $L(\mathcal{A}_1) \cap (\mathcal{U}(\Sigma) \setminus L(\mathcal{A}_2)) = \emptyset$.

Weak Second Order Monadic Logic of Presburger

Syntax of formulae of PMSO.

- ▶ first order variables $x \dots$
- ▶ second order variables $X \dots$
- ▶

$$\begin{array}{l} \text{form} ::= \\ \text{pres} ::= \\ \text{term} ::= \end{array} \left| \begin{array}{l} x = y \mid x \downarrow y \mid a(x) \mid x \in X \mid x/\text{pres} \quad (a \in \Sigma) \\ \text{form} \wedge \text{form} \mid \text{form} \vee \text{form} \mid \neg \text{form} \\ \exists x \text{ form} \mid \exists X \text{ form} \mid \forall x \text{ form} \mid \forall X \text{ form} \\ \text{term} = \text{term} \\ \neg \text{pres} \mid \text{pres} \vee \text{pres} \mid \text{pres} \wedge \text{pres} \\ \forall z \text{ pres} \mid \exists z \text{ pres} \\ z \quad (\text{first order integer variable}) \\ n \quad (\text{natural number}) \\ [X] \quad (X \text{ second order var.}) \\ \text{term} + \text{term} \end{array} \right.$$

formulae pres such that the variables z are bounded.
rem. no atoms $x \rightarrow y$ as for ordered trees.

Weak Second Order Monadic Logic of Presburger

Semantics of PMSO.

- ▶ interpretation domain : set $\|t\|$ of nodes of a tree $t \in \mathcal{U}(\Sigma)$,
- ▶ σ : first order variables $\rightarrow \|t\|$,
- ▶ ρ : second order variable $\rightarrow 2^{\|t\|}$,
- ▶ $t, \sigma, \rho \models x = y$ iff $\sigma(x)$ is $\sigma(y)$,
- ▶ $t, \sigma, \rho \models x \downarrow y$ iff $\sigma(x)$ is the father of $\sigma(y)$ in t ,
- ▶ $t, \sigma, \rho \models a(x)$ iff $\sigma(x)$ labeled by a in t ,
- ▶ $t, \sigma, \rho \models x \in X$ iff $\sigma(x) \in \rho(X)$,
- ▶ $t, \sigma, \rho \models x/\phi$ iff $(n_1, \dots, n_p) \models \phi$ where n_i is the number of sons (in t) of $\sigma(x)$ in $\rho(X_i)$ (with $\text{dom}(\rho) = \{X_1, \dots, X_p\}$),
- ▶ $t, \sigma, \rho \models \psi_1 \vee \psi_2$ iff $t, \sigma, \rho \models \psi_1$ or $t, \sigma, \rho \models \psi_2$,
- ▶ $t, \sigma, \rho \models \psi_1 \wedge \psi_2$ iff $t, \sigma, \rho \models \psi_1$ and $t, \sigma, \rho \models \psi_2$,
- ▶ $t, \sigma, \rho \models \neg\psi$ iff $t, \sigma, \rho \not\models \psi$,
- ▶ $t, \sigma, \rho \models \exists x \psi$ iff there exists $p \in \|t\|$ s.t.
 $t, \sigma \cup \{p \rightarrow x\}, \rho \models \psi$,
- ▶ $t, \sigma, \rho \models \exists X \psi$ iff there exists $P \subseteq \|t\|$ s.t.
 $t, \sigma, \rho \cup \{P \rightarrow X\} \models \psi$,

PMSO: Examples

- ▶ root:

$$x = \varepsilon \equiv \neg \exists y y \downarrow x$$

- ▶ leaf:

$$\text{leaf}(x) \equiv \neg \exists y x \downarrow y$$

- ▶ $x \downarrow y \equiv \exists Y Y = \{y\} \wedge x/[Y]=1$

- ▶ prefix ordering = transitive closure of \downarrow :

$$x \downarrow^* y \equiv \forall X (x \in X \wedge \forall z \forall z' (z \in X \wedge z \downarrow z' \Rightarrow z' \in X)) \Rightarrow y \in X$$

PMSO Languages

Definition : language

The language defined by the closed PMSO formula ψ over Σ is the set of terms $t \in \mathcal{U}(\Sigma)$ s.t. $t \models \psi$.

PMSO Language: example 1

The set of trees of the form $f(a, \dots, a, b, \dots, b)$ with the same number of a and b .

PMSO Language: example 1

The set of trees of the form $f(a, \dots, a, b, \dots, b)$ with the same number of a and b .

$$\begin{aligned} \exists X_a \exists X_b f(\varepsilon) \quad & \wedge \forall y (y \in X_a \Leftrightarrow a(y)) \wedge (y \in X_b \Leftrightarrow b(y)) \\ & \wedge \forall y \varepsilon \downarrow y \Rightarrow (\text{leaf}(y) \wedge (y \in X_a \vee y \in X_b)) \\ & \wedge \varepsilon / [X_a]=[X_b] \end{aligned}$$

PMSO: PA example 2

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

The set of trees of $\mathcal{U}(\Sigma)$ s.t. every internal node is labeled by a and every node labeled by a has at least as much sons without b than sons containing b .

PMSO: PA example 2

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

The set of trees of $\mathcal{U}(\Sigma)$ s.t. every internal node is labeled by a and every node labeled by a has at least as much sons without b than sons containing b .

$$\begin{aligned} \exists X_a \exists X_b \forall x & \quad (b(x) \Rightarrow \text{leaf}(x)) \\ & \wedge (a(x) \wedge x/[X_a] \geq [X_b] > 0 \Rightarrow x \in X_b) \\ & \wedge (a(x) \wedge x/[X_a] \geq [X_b] = 0 \Rightarrow x \in X_a) \\ & \wedge (b(x) \wedge x/[X_a] = [X_b] = 0 \Rightarrow x \in X_b) \end{aligned}$$

$$Q = Q_f = \{q_a, q_b\}.$$

The state q_b accepts the trees containing a b and q_a accepts the others.

$$\Delta = \left\{ \begin{array}{ll} a(x_{q_a} \geq x_{q_b} = 0) & \rightarrow q_a \\ a(x_{q_a} \geq x_{q_b} > 0) & \rightarrow q_b \\ b(x_{q_a} = x_{q_b} = 0) & \rightarrow q_b \end{array} \right\}$$

PMSO: Examples of Queries

Base of clients of an online music store, stored in an unordered unranked tree.

A client is a subtree:

- ▶ root labeled by `client`
- ▶ informations in the sons (purchase, labeled at root by the kind).

Query for clients x who have purchased more jazz than blues:

$$\begin{aligned} \text{query}_1(x) \equiv & (\exists X_{\text{jazz}} \forall y y \in X_{\text{jazz}} \Leftrightarrow \text{jazz}(y)) \\ & \wedge (\exists X_{\text{blues}} \forall y y \in X_{\text{blues}} \Leftrightarrow \text{blues}(y)) \\ & \wedge \text{client}(x) \wedge x / [X_{\text{jazz}}] > [X_{\text{blues}}] \end{aligned}$$

Query for clients x who have purchased more jazz than anything else:

$$\begin{aligned} \text{query}_1(x) \equiv & (\exists X_{\text{jazz}} \forall y y \in X_{\text{jazz}} \Leftrightarrow \text{jazz}(y)) \\ & \wedge (\exists X_{\text{other}} \forall y y \in X_{\text{other}} \Leftrightarrow \neg \text{jazz}(y)) \\ & \wedge \text{client}(x) \wedge x / [X_{\text{jazz}}] > [X_{\text{other}}] \end{aligned}$$

PMSO and PA

Theorem

$L \subseteq \mathcal{O}(\Sigma)$ is definable in PMSO iff L is a PA language.

Plan

Unranked Trees and Reasoning Tasks over XML Documents

Automata for Unranked Ordered Trees

Automata for Unranked Unordered Trees

Automata for Unranked Mixed Trees

Presburger Constraints and Unranked Ordered Trees
Mixed Trees

Regular Languages modulo Associativity and Commutativity

Verification of XML Updates

Presburger Hedge Automata (PHA)

Definition : Presburger Hedge Automata

A *Presburger Hedge Automaton* (PHA) over an alphabet Σ is a tuple $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$ where $Q = \{q_1, \dots, q_p\}$ is a finite set of *states*, $Q^f \subseteq Q$ is the subset of final states and Δ is a set of transition rules of the form: $a(\bigvee_i (L_i \wedge \phi_i)) \rightarrow q$ with $a \in \Sigma$, $q \in Q$, $L_i \subseteq Q^*$ regular language and $\phi_i = \phi(x_1, \dots, x_p)$ is a Presburger formula with a free variable for each state.

for all $w \in Q^*$ we define $w \models L_i \wedge \phi_i$ if $w \in L_i$ and $\#(w) \models \phi_i$.

PHA: Languages

The language $L(\mathcal{A}, q)$ of \mathcal{A} in state $q \in Q$ is the smallest set of **ordered** unranked trees $a(t_1, \dots, t_n) \in \mathcal{O}(\Sigma)$ s.t.

- ▶ there exists $i_1, \dots, i_n \leq p$ such that for all $j \leq n$,
 $t_j \in L(\mathcal{A}, q_{i_j})$,
- ▶ there exists a transition $a(\bigvee_i (L_i \wedge \phi_i)) \rightarrow q \in \Delta$ such that
 $q_{i_1} \dots q_{i_n} \models \bigvee_i (L_i \wedge \phi_i)$, i.e. there exists i s.t.
 - ▶ $q_{i_1} \dots q_{i_n} \in L_i$,
 - ▶ $\#(q_{i_1} \dots q_{i_n}) \models \phi_i(x_1, \dots, x_p)$.

The language of \mathcal{A} is $L(\mathcal{A}) = \bigcup_{q \in Q_f} L(\mathcal{A}, q)$.

PHA: Properties

Proposition :

- ▶ The class of PHA languages is closed under union and intersection,
 - ▶ The class of PHA languages is not closed under complementation.
-
- ▶ \cup, \cap : product
 - ▶ csq undecidability of the problem of universality.

Proposition :

$DPHA \neq NPHA$.

PHA: Decision Problems

Lemma :

Given a finite set Q , $L \subseteq Q^*$, regular, and $\phi = \phi(x_1, \dots, x_p)$ a Presburger formula ($p = |Q|$), it is decidable whether there exists $w \in L$ such that $\#(w) \models \phi$.

Proposition : \in

The membership is decidable for PHA.

Proposition : \emptyset

The emptiness is decidable for PHA.

Proposition : \forall

Universality is undecidable for PHA.

PHA: Logic

Theorem :

A set of trees of $\mathcal{O}(\Sigma)$ is recognizable by a PHA iff it is defined by a PMSO formula of the form $\exists X_1 \dots \exists X_k \phi$ where ϕ is first order.

Corollary :

- ▶ EPMSO (existential fragment) is decidable in $\mathcal{O}(\Sigma)$.
- ▶ PMSO is undecidable over $\mathcal{O}(\Sigma)$.

Mixed Trees

- ▶ $\Sigma = \Sigma_A \cup \Sigma_{AC}$.

$$\begin{aligned} \text{tree} &:= a(\text{hedge}) \mid c(\text{multiset}) \quad (a \in \Sigma_A, c \in \Sigma_{AC}) \\ \text{hedge} &:= \text{tree}, \dots, \text{tree} \\ \text{multiset} &:= \{\text{tree}, \dots, \text{tree}\} \end{aligned}$$

- ▶ $c(\{t_1, \dots, t_n\})$ is denoted $c(t_1, \dots, t_n)$.
- ▶ The set of mixed unranked trees over Σ is denoted $\mathcal{M}(\Sigma)$.

Presburger m-Tree Automata (PMA)

Definition : Presburger m-Tree Automata

A *Presburger m-Tree Automaton* (PMA) over an alphabet $\Sigma = \Sigma_A \cup \Sigma_{AC}$ is a tuple $\mathcal{A} = (\Sigma, Q, Q^f, \Delta)$ where $Q = \{q_1, \dots, q_p\}$ is a finite set of *states*, $Q^f \subseteq Q$ is the subset of final states and Δ is a set of transition rules of the form:

- ▶ $a(L) \rightarrow q$ with $a \in \Sigma_A$, $q \in Q$, $L \subseteq Q^*$ is a regular language
or
- ▶ $c(\phi) \rightarrow q$ with $c \in \Sigma$, $q \in Q$, $\phi = \phi(x_1, \dots, x_p)$ is a Presburger formula with one free variable for each state.

PMA: Languages

The language $L(\mathcal{A}, q)$ of the PMA \mathcal{A} in state $q \in Q$, is the smallest set of mixed trees

- ▶ $a(t_1, \dots, t_n)$, $a \in \Sigma_A$, such that
 - ▶ there exists $i_1, \dots, i_n \leq p$ with $t_j \in L(\mathcal{A}, q_{i_j})$ for all $j \leq n$,
 - ▶ there exists a transition $a(L) \rightarrow q \in \Delta$ such that $q_{i_1} \dots q_{i_n} \in L$,
- ▶ or $c(t_1, \dots, t_n)$, $c \in \Sigma_{AC}$, such that
 - ▶ there exists $i_1, \dots, i_n \leq p$ with $t_j \in L(\mathcal{A}, q_{i_j})$ for all $j \leq n$,
 - ▶ there exists a transition $c(\phi) \rightarrow q \in \Delta$ such that $\#(q_{i_1} \dots q_{i_n}) \models \phi(x_1, \dots, x_p)$.

The language of \mathcal{A} is $L(\mathcal{A}) = \bigcup_{q \in Q_f} L(\mathcal{A}, q)$.

PMA: Properties and Decision Results

Proposition :

The class of PMA languages is closed under all Boolean operations.

Proposition :

$DPMA \equiv NPMA$.

Proposition : \in

Membership is decidable for PMA.

Proposition : \emptyset

Emptiness is decidable for PMA.

Consequences of the analogous results for PHA ($PMA \subseteq PHA$).

PMA: Logic

Theorem :

The class of languages of $\mathcal{M}(\Sigma)$ definable by PMSO formulae is the class of PMA languages.

Corollary

PMSO over $\mathcal{M}(\Sigma)$ is decidable.

Plan

Unranked Trees and Reasoning Tasks over XML Documents

Automata for Unranked Ordered Trees

Automata for Unranked Unordered Trees

Automata for Unranked Mixed Trees

Regular Languages modulo Associativity and Commutativity

- Variants of HA

- Regular Languages modulo Associativity

- Regular Languages modulo Associativity and Commutativity

Verification of XML Updates

Imperative Program

[Bouajjani Touili CAV 02]

void X() {	X	→	Y · X	(r ₁)
while(true) {	Y	→	t	(r ₂)
if Y() {	Y	→	f	(r ₃)
thread_create(&t1,Z)	t · X	→	X Z	(r ₄)
} else { return }	f	→	0	(r ₅)
}				
}				

Reachability analysis:

- ▶ The set of reachable terms is **regular**, but
 - ▶ we want **· Associative**,
 - ▶ and **|| Associative** and **Commutative**,
 - ▶ and regular term languages are not closed **modulo A** and **AC**.
- consider unranked trees as representative.

Extensions of HA

Definition : CF-HA

A CF-HA is a tuple (Σ, Q, Q_f, Δ) , where Q, Q_f are as for HA and the transitions of Δ have the form $a(L) \rightarrow q$ with $a \in \Sigma, q \in Q$, and $L \subseteq Q^*$ is a context-free language.

Definition : CS-HA

A CS-HA is a tuple (Σ, Q, Q_f, Δ) where Q, Q_f are as for HA and the transitions of Δ have the form $a(L) \rightarrow q$ with $a \in \Sigma, q \in Q$, and $L \subseteq Q^*$ is a context-sensitive language.

CF-HA: Example

$\Sigma = \{a, b, f\}$, language L of trees of $\mathcal{O}(\Sigma)$:

- ▶ whose internal nodes are labeled by f ,
- ▶ with the same number of leaves a than leaves b under every node.

CF-HA: Example

$\Sigma = \{a, b, f\}$, language L of trees of $\mathcal{O}(\Sigma)$:

- ▶ whose internal nodes are labeled by f ,
- ▶ with the same number of leaves a than leaves b under every node.

language of the CF-HA (Q, Q_f, Δ) with $Q = \{q, q_a, q_b\}$ and

$$\Delta = \{a \rightarrow q_a, \quad b \rightarrow q_b, \quad f(L) \rightarrow q\}$$

L is the language generated by the context-free grammar

$$\begin{aligned} N &:= \varepsilon \mid \text{all permutations of } NN_aN_b \mid q \\ N_a &:= q_a \quad N_b := q_b \end{aligned}$$

Rem.: L is not a HA language.

Regular Languages modulo A (A-TA)

Signature $\Sigma = \Sigma_{\emptyset} \uplus \{a\}$.

The symbols of a is binary and follows the associativity axiom:

$$a(x_1, a(x_2, x_3)) = a(a(x_1, x_2), x_3) \quad (\text{A})$$

Given a TA \mathcal{B} over Σ , we note

$$A(L(\mathcal{B})) := \{t \in \mathcal{T}(\Sigma) \mid t \xrightarrow[A]{*} s \in L(\mathcal{B})\}$$

(A-TA language)

Proposition :

- ▶ the class of regular tree languages is strictly included in the class of A-TA languages.
- ▶ The class of A-TA languages is not closed under intersection.

Correspondences $\mathcal{T}(\Sigma) \leftrightarrow \mathcal{O}(\Sigma)$

$\Sigma = \Sigma_\emptyset \uplus \{a\}$ where a is the only associative symbol.

$$\text{flat} : \mathcal{T}(\Sigma) \rightarrow \mathcal{O}(\Sigma)$$

$$\text{hflat} : \mathcal{T}(\Sigma)^* \rightarrow \mathcal{H}(\Sigma)$$

$$\text{flat}^{-1} : \mathcal{O}(\Sigma) \rightarrow \mathcal{T}(\Sigma)$$

Definitions ($g \in \Sigma_n \setminus \Sigma_A$):

$$\text{flat}(g(t_1, \dots, t_n)) = g(\text{flat}(t_1) \dots \text{flat}(t_n))$$

$$\text{flat}(a(t_1, t_2)) = a(\text{hflat}(t_1 t_2))$$

$$\text{hflat}(g(s_1, \dots, s_n) t_2 \dots t_m) = \text{flat}(g(s_1, \dots, s_n)) \text{hflat}(t_2 \dots t_m)$$

$$\text{hflat}(a(s_1, s_2) t_2 \dots t_m) = \text{hflat}(s_1 s_2 t_2 \dots t_m)$$

$$\text{flat}^{-1}(g(t_1 \dots t_n)) = g(\text{flat}^{-1}(t_1), \dots, \text{flat}^{-1}(t_n))$$

$$\begin{aligned} \text{flat}^{-1}(a(t_1 \dots t_m)) &= a(\text{flat}^{-1}(t_1), a(\text{flat}^{-1}(t_2), \dots, \\ &\quad a(\text{flat}^{-1}(t_{m-1}), \text{flat}^{-1}(t_m)))) \\ &\quad (m \geq 2) \end{aligned}$$

A-TA \leftrightarrow CF-HA

Proposition :

A-TA \equiv CF-HA via flattening.

\subseteq for all TA \mathcal{A} there exists a CF-HA \mathcal{A}' such that
 $L(\mathcal{A}') = \text{flat}(L(\mathcal{A})) = \text{flat}(\mathbf{A}(L(\mathcal{A})))$.

TA \mathcal{A}	CF-HA \mathcal{A}'
$a(q_1, q_2) \rightarrow q$	$N_q := N_{q_1} N_{q_2}, N_q := q$
$g(q_1, \dots, q_k) \rightarrow q$	$g(q_1 \dots q_k) \rightarrow q$

\supseteq for all CF-HA \mathcal{A}' there exists a TA \mathcal{A} such that
 $L(\mathcal{A}) = \text{flat}^{-1}(L(\mathcal{A}'))$ (i.e. $\text{flat}(\mathbf{A}(L(\mathcal{A}))) = L(\mathcal{A}')$).

CF-HA \mathcal{A}'	TA \mathcal{A}
$N := N_1 N_2$	$a(q_{N_1}, q_{N_2}) \rightarrow q_N$
$I := N_1 N_2$	$a(q_{N_1}, q_{N_2}) \rightarrow q$
$g(q_1 \dots q_k) \rightarrow q$	$g(q_1, \dots, q_k) \rightarrow q$

Generalized Tree Automata (GTA)

Definition : GTA

A *generalized Tree Automata* (GTA) over a signature Σ is a tuple $\mathcal{B} = (\Sigma, Q, Q^f, \Delta)$ where Q is a finite set of *states*, $Q^f \subseteq Q$ is the subset of final states and Δ is a set of transition rules of the form $f(q_1, \dots, q_n) \rightarrow q$ or $a(q_1, q_2) \rightarrow a(q'_1, q'_2)$ with $f \in \Sigma_n$ ($n \geq 0$) and $q_1, \dots, q_n, q'_1, \dots, q'_n \in Q$.

For a TA or GTA \mathcal{B} , we define the languages:

$$\begin{aligned} L(\mathcal{B}) &:= \{t \in \mathcal{T}(\Sigma) \mid t \xrightarrow{\Delta^*} q \in Q^f\} \\ L_A(\mathcal{B}) &:= \{t \in \mathcal{T}(\Sigma) \mid t \xrightarrow{\Delta/A^*} q \in Q^f\} \end{aligned}$$

where $\xrightarrow{\Delta/A^*}$ is rewriting modulo A: $\xrightarrow{\Delta/A^*} := \xleftarrow{A^*} \circ \xrightarrow{\Delta} \circ \xleftarrow{A^*}$

Languages of TA and GTA modulo A

Proposition :

For all GTA \mathcal{B} , there exists a TA \mathcal{B}' such that $L(\mathcal{B}') = L(\mathcal{B})$.

Proposition :

For all TA \mathcal{B} , $L_A(\mathcal{B}) = A(L(\mathcal{B}))$.

Csq: the emptiness of $L_A(\mathcal{B})$ for a TA \mathcal{B} is decidable.

Proposition :

For a GTA \mathcal{B} , in general $L_A(\mathcal{B}) \neq A(L(\mathcal{B}))$.

$\Delta_{\mathcal{B}} = \{a \rightarrow q_1, b \rightarrow q_2, f(q_1, q_2) \rightarrow f(q_3, q_3), f(q_3, q_2) \rightarrow q_2\}$.

$L(\mathcal{B}, q_2) = \{b\}$

$L_A(\mathcal{B}, q_2) \ni f(f(a, b), b)$

CS-HA \leftrightarrow GTA modulo A

Proposition :

A-GTA \equiv CS-HA via flattening.

- \subseteq for all GTA \mathcal{A} there exists a CS-HA \mathcal{A}' such that
 $L(\mathcal{A}') = \text{flat}(L_{\mathcal{A}}(\mathcal{A}))$.

GTA \mathcal{A}	CS-HA \mathcal{A}'
$a(q_1, q_2) \rightarrow q$	$N_q := N_{q_1} N_{q_2}$
$a(q_1, q_2) \rightarrow a(q_3, q_4)$	$N_{q_3} N_{q_4} := N_{q_1} N_{q_2}$
	$N_q := q$
$g(q_1, \dots, q_k) \rightarrow q$	$g(q_1 \dots q_k) \rightarrow q$

- \supseteq for all CS-HA \mathcal{A}' there exists a GTA \mathcal{A} such that
 $L_{\mathcal{A}}(\mathcal{A}) = \mathbf{A}(\text{flat}^{-1}(L(\mathcal{A}')))$ (i.e. $\text{flat}(L_{\mathcal{A}}(\mathcal{A})) = L(\mathcal{A}')$)

CS-HA \mathcal{A}'	TA \mathcal{A}
$N := N_1 N_2, N \neq I$	$a(q_{N_1}, q_{N_2}) \rightarrow q_N$
$N_1 N_2 := N_3 N_4$	$a(q_{N_3}, q_{N_4}) \rightarrow a(q_{N_1}, q_{N_2})$
$I := N_1 N_2$ for $L_{a,q}$	$a(q_{N_1}, q_{N_2}) \rightarrow q$
$g(q_1 \dots q_k) \rightarrow q$	$g(q_1, \dots, q_k) \rightarrow q$

GTA & CS-HA: Decision Results

Proposition : \in CS-HA

The membership problem is PSPACE-complete for CS-HA.

Proposition : \in GTA

The membership problem $t \in L_A(\mathcal{B})$ given a GTA \mathcal{B} is PSPACE-complete.

Proposition : \emptyset CS-HA

The emptiness problem is undecidable for CS-HA.

Proposition : \emptyset GTA

The emptiness problem $L_A(\mathcal{B})$ given a GTA \mathcal{B} is undecidable.

Regular Languages modulo AC (AC-TA)

Signature $\Sigma = \Sigma_{\emptyset} \uplus \{a\}$.

The symbol a is binary and follows the axioms of associativity and commutativity:

$$a(x_1, a(x_2, x_3)) = a(a(x_1, x_2), x_3) \quad (\text{A})$$

$$a(x_1, x_2) = a(x_2, x_1) \quad (\text{C})$$

For a TA \mathcal{B} over Σ , we note

$\text{AC}(L(\mathcal{B})) := \{t \in \mathcal{T}(\Sigma) \mid t =_{\text{AC}} s \in L(\mathcal{B})\}$ (language of AC-TA)

Proposition :

- ▶ the class of regular tree languages is strictly included in the class of AC-TA languages.
- ▶ La class of AC-TA languages is closed under Boolean operations.

Correspondence between PA, HA and AC-TA

Proposition :

There is a equivalence, via flattening, between

1. PA
2. CF-HA whose (CF) languages in transitions are closed under permutation
3. AC-TA

pr.:

- ▶ $1 \equiv 2$ by Parikh's theorem (equivalence between solutions of Presburger formulas and Parikh projection of CF languages)
- ▶ $2 \equiv 3$ with above constructions for A-TA and of previous lecture with $flat, flat^{-1}$.

Summary

CS-HA \equiv A-GTA | Boolean closures
 \emptyset undecidable.

CF-HA \equiv A-TA | no closure under \cap and \neg
 \emptyset decidable.

HA \equiv TA | Boolean closures
 \emptyset decidable.

PA \equiv CF-HA/perm \equiv AC-TA
Boolean closures
 \emptyset decidable.

Plan

Unranked Trees and Reasoning Tasks over XML Documents

Automata for Unranked Ordered Trees

Automata for Unranked Unordered Trees

Automata for Unranked Mixed Trees

Regular Languages modulo Associativity and Commutativity

Verification of XML Updates

Verification of XML Transformations

Typechecking XML transformations languages: different models for different languages

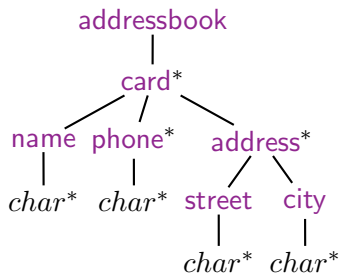
- ▶ [Milo Suciu Vianu 03]: k-PTTs (XSLT fragment without data joins); backward type inference.
- ▶ [Maneth Berlea Perst Seidl 05], [Engelfriet Maneth Seidl 09]: language LT and MTTs.
- ▶ [Martens Neven 04]: top-down tree transducers.
- ▶ [Tozawa 01] backward type inference for XSLT without XPath, with recursive calls (alternating TA).
- ▶ [Frish Hozawa 07] backward type inference for MTTs (optimizations).

XQuery Update Facility

W3C recommendation 2011

- ▶ [Fundulaki Maneth 04] XACU: model based on the W3C XQuery Update Facility draft.
- ▶ [Benedikt Cheney 10] formal model, operational semantics.
- ▶ [Bravo Cheney Fundulaki 08] synthesis of schema, verification tool ACCoN.
- ▶ [Gardner et al 08] local Hoare reasoning about W3C DOM update library (Context Logic).
- ▶ [JR PPDP 10]
 - formal model: parameterized [rewrite rules](#)
 - forward/backward type inference, typechecking, reachability verification [regular tree model checking](#)
 - verification of [access control policies](#)

A DTD and a HA



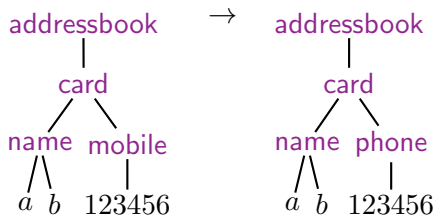
$addressbook(p_c^*) \rightarrow p_b$
 $card(p_n p_h^* p_a^*) \rightarrow p_c$
 $card(p_n) \rightarrow p_{ec}$
 $name(p^*) \rightarrow p_n$
 $phone(p^*) \rightarrow p_h$
 $address(p_{st} p_{ci}) \rightarrow p_a$
 $street(p^*) \rightarrow p_{st}$
 $city(p^*) \rightarrow p_{ci}$
 $a \rightarrow p$
 $b \rightarrow p$
 $c \rightarrow p$
 \vdots

XUF_{reg} Rule: Rename

"replace a label **mobile** by **phone**"

$$\text{mobile}(x) \rightarrow \text{phone}(x)$$

- ▶ the variable x stands for a sequence of trees (hedge).
- ▶ the rule can be applied to any node labeled by **mobile**.



XUF_{reg} Rule: Insert first

"insert a tree of type p_{ec} (card with only a name) as the first children of **addressbook**"

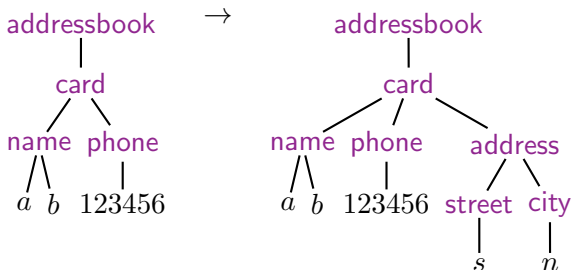
$$\text{addressbook}(x) \rightarrow \text{addressbook}(p_{ec} x)$$

- ▶ the rewrite rule can be applied to any node labeled by **addressbook**.
- ▶ p_{ec} is a state of a given HA \mathcal{A} .
- ▶ it stands for an arbitrary tree in the language recognized by \mathcal{A} in state p_{ec} .
- ▶ this parametrized rule represents an infinity of rules.
see [Gilleron 91], [Löding 02, 07].

XUF_{reg} Rule: Insert last

"insert a tree of type p_a (address) as the last children of $card$ "

$$card(x) \rightarrow card(x p_a)$$



XUF_{reg} Rule: Insert into

"insert a tree of type p_{ec} as a children of **addressbook**"

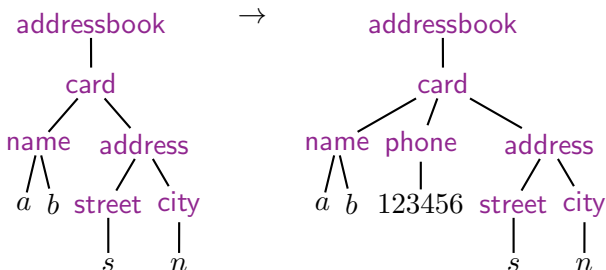
$$\text{addressbook}(x y) \rightarrow \text{addressbook}(x p_{ec} y)$$

- ▶ each of the variables x and y stands for an arbitrary sequence of trees (hedge).

XUF_{reg} Rule: Insert before

"insert a tree of type p_h (phone) as preceding sibling of **address**"

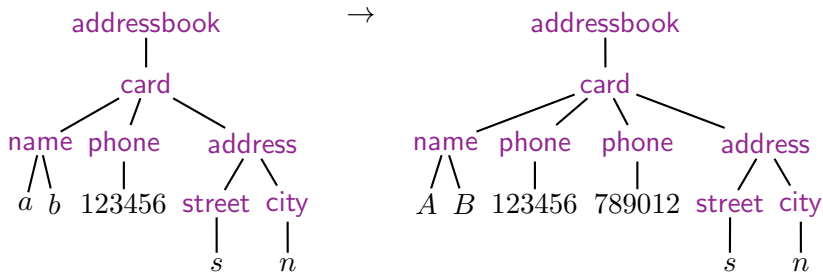
$$\text{address}(x) \rightarrow p_h \text{ address}(x)$$



XUF_{reg} Rule: Insert after

"insert a tree of type p_h (phone) as following sibling of $phone$ "

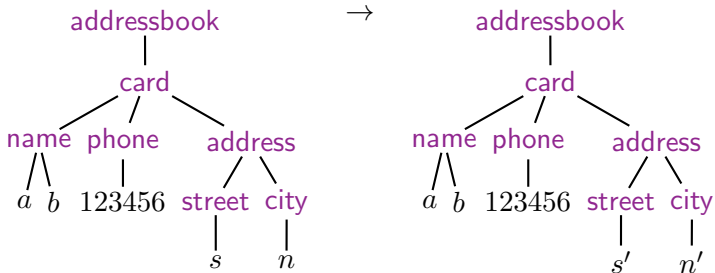
$phone(x) \rightarrow phone(x) p_h$



XUF_{reg} Rule: Replace

"replace a subtree (headed by) **address** by an arbitrary tree of type p_a "

$$\text{address}(x) \rightarrow p_a$$



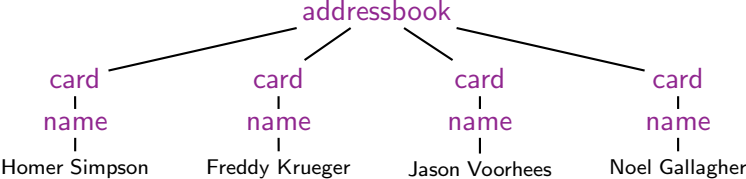
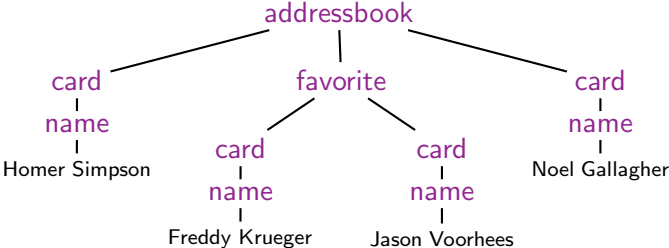
XUF_{reg} Rule: Delete

"delete a whole subtree headed by **card**"

$$\text{card}(x) \rightarrow ()$$

- ▶ $()$ is the empty sequence of trees.

Delete single node (not a XQuery Update Primitive)



Extended XUF_{reg} Rule: Delete single node

"delete a single node labeled by *favorite*"

$$\text{favorite}(x) \rightarrow x$$

- ▶ the trees in the sequence of children x are moved up to the position of the deleted node.
- ▶ *collapsing rule*

Extended XUF_{reg} Rule: Multiple Replace

"replace a subtree (headed by) card by an sequence of n trees of respective types p_1, \dots, p_n "

$$\text{card}(x) \rightarrow p_1 \dots p_n$$

- ▶ this parametrized rule represents an infinity of rules.
- ▶ the right hand sides of these rules are hedges (not trees).

Summary

XUF _{reg}			XUF		
$a(x)$	\rightarrow	$b(x)$	REN		
$a(x)$	\rightarrow	$a(px)$	INS _{first}	$a(x)$	\rightarrow $pa(x)$ INS _{before}
$a(x)$	\rightarrow	$a(xp)$	INS _{last}	$a(x)$	\rightarrow $a(x)p$ INS _{after}
$a(xy)$	\rightarrow	$a(xpy)$	INS _{into}		
$a(x)$	\rightarrow	p	RPL ₁	$a(x)$	\rightarrow $p_1 \dots p_n$ RPL
$a(x)$	\rightarrow	$()$	DEL	$a(x)$	\rightarrow x DEL _s

Forward Type Inference for XUF_{reg}

Theorem :

For \mathcal{R} in XUF_{reg} , L_{in} HA language, $\mathcal{R}^*(L_{\text{in}})$ is a HA language.

pr.:

- ▶ let \mathcal{A}_{in} be the HA for L_{in} and \mathcal{A} be the HA parameter of \mathcal{R} ,
- ▶ the HA \mathcal{A}_{out} is obtained from $\mathcal{A}_{\text{in}} \uplus \mathcal{A}$, by adding transitions to the horizontal NFAs:

Forward Type Inference for XUF_{reg}

Theorem :

For \mathcal{R} in XUF_{reg} , L_{in} HA language, $\mathcal{R}^*(L_{\text{in}})$ is a HA language.

pr.:

- ▶ let \mathcal{A}_{in} be the HA for L_{in} and \mathcal{A} be the HA parameter of \mathcal{R} ,
- ▶ the HA \mathcal{A}_{out} is obtained from $\mathcal{A}_{\text{in}} \uplus \mathcal{A}$, by adding transitions to the horizontal NFAs:

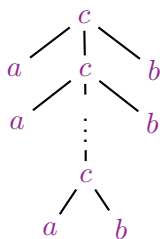
case INS_{first} : $a(x) \rightarrow a(px) \in \mathcal{R}$, $a \left(\begin{array}{c} \rightarrow \textcircled{i_{a,q}} \text{---} \textcircled{\text{f}_{a,q}} \\ \end{array} \right) \rightarrow q$

Forward Type Inference for XUF

Lemma

For some \mathcal{R} in XUF and L_{in} HA language, $\mathcal{R}^*(L_{\text{in}}) \notin \text{HA}$.

example: $\mathcal{R} = \{c(x) \rightarrow x\}$ (one DEL_S) and $L_{\text{in}} =$



$$\mathcal{R}^*(L_{\text{in}}) \cap c(\{a, b\}^*) = \{c(a^n b^n) \mid n \geq 0\}$$

Forward Type Inference for XUF

Theorem :

For \mathcal{R} in XUF, L_{in} a CF-HA language, $\mathcal{R}^*(L_{\text{in}})$ is a CF-HA language, with a PTIME, polynomial size construction.

pr.: similar as XUF_{reg} for REN , $\text{INS}_{\text{first}}$, INS_{last} , INS_{into} .

- ▶ **REN:** if $a(x) \rightarrow b(x) \in \mathcal{R}$, then add $I'_{b,q} := I'_{a,q}$
- ▶ **INS_{first}:** if $a(x) \rightarrow b(px) \in \mathcal{R}$, then add $I'_{b,q} := p I'_{a,q}$

add **collapsing transitions** for $\text{INS}_{\text{before}}$, $\text{INS}_{\text{after}}$, RPL , DEL , DEL_S .

- ▶ **RPL:** if $a(x) \rightarrow p_1 \dots p_n \in \mathcal{R}$ and $a(L) \rightarrow q$ transition, $L \neq \emptyset$, then add the collapsing transition $p_1 \dots p_n \rightarrow q$.

Backward Type Inference for XUF

Theorem :

For \mathcal{R} in XUF, L_{out} a HA language, $(\mathcal{R}^{-1})^*(L_{\text{out}})$ is a HA language, with a EXPTIME, exponential size construction.

pr.: tree automata completion; construction of a finite sequence of HA

$$\mathcal{A}_{\text{out}} = \mathcal{A}_0, \mathcal{A}_1, \dots, \mathcal{A}_n = (\mathcal{R}^{-1})^*(\mathcal{A}_{\text{out}})$$

by addition of new transitions.

Applications

Theorem :

For \mathcal{R} in XUF, L_{in} a CF-HA language, $\mathcal{R}^*(L_{in})$ is a CF-HA language, with a PTIME, polynomial size construction.

Consequences:

- ▶ **reachability** for XUF
 - ▶ decidable in PTIME.
- ▶ **typechecking** XUF
given two HA for L_{in} , L_{out} and \mathcal{R} in XUF, $\mathcal{R}^*(L_{in}) \subseteq L_{out}$?
 - ▶ EXPTIME-complete
 - ▶ PTIME if L_{out} is given by a deterministic and complete HA.
- ▶ **type synthesis**:
 - ▶ given \mathcal{R} in XUF and an input type L_{in} , a CF-HA for L_{out} is constructed in PTIME.
 - ▶ given \mathcal{R} in XUF and an output type L_{out} , a HA for L_{in} is constructed in EXPTIME.

Rule Based Access Control Policies

Definition : [Fundulaki Maneth][Bravo et al, ACCOn]

An access control policy (ACP) is given by two finite sets of rules

- ▶ \mathcal{R}_+ : authorized operations of XUF
- ▶ \mathcal{R}_- : forbidden operations of XUF

example

- ▶ $\mathcal{R}_+ = \left\{ \begin{array}{l} \text{addressbook}(x) \rightarrow \text{addressbook}(p_{ec} x), \\ \text{card}(x) \rightarrow () \end{array} \right\}$
 - ▶ user can insert card with name, delete card.
- ▶ $\mathcal{R}_- = \{ \text{name}(x) \rightarrow p_n \}$
 - ▶ user cannot change a name.

Inconsistency

Inconsistency

An ACP $\langle \mathcal{R}_+, \mathcal{R}_- \rangle$ is **inconsistent** if one rule of \mathcal{R}_- can be simulated through a sequence of rules of \mathcal{R}_+

example: changing name in a card is simulated by deleting and then inserting.

Theorem : [Fundulaki Maneth 04] [Moore 09]

Inconsistency is undecidable for XUF_{reg} .

Local Inconsistency

Definition : Local Inconsistency

An ACP $\langle \mathcal{R}_+, \mathcal{R}_- \rangle$ is **locally inconsistent** for t if there exists u such that $t \xrightarrow{\mathcal{R}_-} u$ and $t \xrightarrow{\mathcal{R}_+^*} u$.

Theorem : [JR PPDP 10]

Local inconsistency is decidable in PTIME for XUF.

pr.:

- ▶ compute a HA recognizing $L_- = \{u \mid t \xrightarrow{\mathcal{R}_-} u\}$
- ▶ compute a CF-HA recognizing $L_+ = \mathcal{R}_+^*(\{t\})$
- ▶ check that $L_- \cap L_+ \neq \emptyset$.

Summary

- ▶ model for Xquery Update facility primitives (and extensions) as parameterized rewrite rules for unranked trees
- ▶ results of forward/backward type inference (rewrite closure)
- ▶ reachability decision
- ▶ decision of local inconsistency of ACPs

Extensions

- ▶ rules controlled with downward regular XPath expressions.
- ▶ better model: generalize \mathcal{R}^* to \mathcal{R}^e (e regular expression over rule names).
- ▶ unranked unordered trees (Presburger Automata).

Part IV

Tree Automata Defined as Sets of Horn Clauses

Plan

Definition

Saturation Results

Tree Automata with Equality Constraints

Tree Automata with One Memory

Definition of tree automata as set of first order (universal) clauses.
Languages = Herbrand models.

- + uniform formalism for the definition of several classes of automata (alternating, 2-ways, with constraints...)
- + enables the use of techniques and tools from automatic deduction in order to solve the classical decision problems
- complexity
- not easy to analyse the history
(and construct a witness, a run...)

Clauses: Syntax

- ▶ **terms** in $\mathcal{T}(\Sigma, \mathcal{X})$ over signature Σ (Σ_n : symbol of arity n)
- ▶ finite set \mathcal{P} of **predicate** symbols P, Q, \dots (notation \mathcal{P}_n)
basically we will only consider predicates of arity 0 or 1.
- ▶ **literals** positive: $P(t)$, denoted $+P(t)$
negative: $\neg P(t)$, denoted $-P(t)$
- ▶ **clause**: disjunction of literals $\pm P_1(t_1) \vee \dots \vee \pm P_k(t_k)$
empty clause ($k = 0$), denoted \perp .
- ▶ **Horn clause**: at most one positive literal
 $-P_1(t_1) \vee \dots \vee -P_k(t_k) \vee +P(t)$, denoted
 $P_1(t_1), \dots, P_k(t_k) \Rightarrow P(t)$.
- ▶ **goal** = negative clause
 $-P_1(t_1) \vee \dots \vee -P_k(t_k)$, denoted $P_1(t_1), \dots, P_k(t_k) \Rightarrow \perp$.

Herbrand Models

- ▶ a Herbrand structure \mathcal{H} has domain $\mathcal{T}(\Sigma)$ and functions $f^{\mathcal{H}}(t_1, \dots, t_n) := f(t_1, \dots, t_n)$ (ground term with the symbol f at the root).
- ▶ \mathcal{H} is completely defined by the set of ground atoms $P(t)$ such that $\mathcal{H} \models P(t)$.

Theorem :

A set of clauses S is satisfiable iff it admits a Herbrand model.

Theorem :

Every satisfiable set S of Horn clauses admits a smallest (wrt inclusion) Herbrand model \mathcal{H}_S .

Smallest Herbrand Models

Theorem :

Every satisfiable set S of Horn clauses admits a smallest (wrt inclusion) Herbrand model \mathcal{H}_S .

A set of Horn clauses S defines the following operator T_S over sets of ground atoms

$$T_S(L) = \left\{ \begin{array}{l} P(t\sigma) \mid t\sigma \text{ ground, } P_1(t_1), \dots, P_n(t_n) \Rightarrow P(t) \in S, \\ P_1(t_1\sigma), \dots, P_n(t_n\sigma) \in L \\ \cup \{ \perp \} \text{ if } P_1(t_1), \dots, P_n(t_n) \Rightarrow \perp \in S, \\ P_1(t_1\sigma), \dots, P_n(t_n\sigma) \in L \end{array} \right\}$$

The smallest fixpoint of T_S is $\bigcup_{n \geq 1} T_S^n(\emptyset)$,

- ▶ if it contains \perp , then S is not satisfiable,
- ▶ otherwise, it is the smallest Herbrand model of S .

Languages and Automata

The **language** of a satisfiable set S of Horn clauses for a predicate Q is:

$$L(S, Q) = \{t \mid Q(t) \in \mathcal{H}_S\}$$

Let $\mathcal{A} = (\Sigma, \{q_1, \dots, q_k\}, F, \Delta)$ be a bottom-up tree automaton.

Let $\mathcal{P} = \{Q_1, \dots, Q_k\}$ be a set of unary predicates.

We associate to \mathcal{A} the (satisfiable) set of Horn clauses

$$S_{\mathcal{A}} := \left\{ \begin{array}{l} Q_1(x_1), \dots, Q_n(x_n) \Rightarrow Q(f(x_1, \dots, x_n)) \\ \mid f(q_1, \dots, q_n) \rightarrow q \in \Delta \end{array} \right\}$$

Lemma :

For all state q , $L(\mathcal{A}, q) = L(S_{\mathcal{A}}, Q)$.

Clauses/Classes of Automata

clauses of standard automata (x_1, \dots, x_n pairwise distinct)

$$Q_1(x_1), \dots, Q_n(x_n) \Rightarrow Q(f(x_1, \dots, x_n)) \quad (\text{reg})$$

ε -transitions

$$Q_1(x) \Rightarrow Q(x) \quad (\varepsilon)$$

alternating clauses

$$Q_1(x), \dots, Q_n(x) \Rightarrow Q(x) \quad (\text{alt})$$

2-ways (bidirectional) clauses (x_1, \dots, x_n pairwise distinct)

$$Q(f(x_1, \dots, x_n)) \Rightarrow Q_i(x_i) \quad (\text{bidi})$$

Plan

Definition

Saturation Results

Tree Automata with Equality Constraints

Tree Automata with One Memory

Decision Problems, Satisfiability

Let S be a satisfiable set of Horn clauses and let Q be a predicate.

- ▶ **membership**: the ground term $t \in L(S, Q)$ iff $S \cup \{Q(t) \Rightarrow \perp\}$ is unsatisfiable
- ▶ **emptiness**: $L(S, Q) \neq \emptyset$ iff $S \cup \{Q(x) \Rightarrow \perp\}$ is unsatisfiable.
- ▶ **membership of instance**: there exists σ such that $t\sigma \in L(S, Q)$ iff $S \cup \{Q(t) \Rightarrow \perp\}$ is unsatisfiable.
- ▶ **emptiness of intersection**: $L(S, Q_1) \cap \dots \cap L(S, Q_p) \neq \emptyset$ iff $S \cup \{Q_1(x), \dots, Q_p(x) \Rightarrow \perp\}$ is unsatisfiable.

\Rightarrow we are interested in automated deduction techniques for deciding the satisfiability when S represents an automaton.

Resolution

Clauses:

$$\frac{C \vee +Q(s) \quad -Q(t) \vee D}{C\sigma \vee D\sigma}$$

where σ is the most general unifier (*mgu*) of s and t .

Horn clauses:

$$\frac{P_1(s_1), \dots, P_m(s_m) \Rightarrow Q_1(s) \quad Q_1(t_1), Q_2(t_2), \dots, Q_n(t_n) \Rightarrow Q(t)}{P_1(s_1\sigma), \dots, P_m(s_m\sigma), Q_2(t_2\sigma), \dots, Q_n(t_n\sigma) \Rightarrow Q(t\sigma)}$$

where σ is the *mgu* of s and t_1 .

Theorem : correction, completeness

A set S of Horn clauses is unsatisfiable iff one can derive \perp by resolution starting from S .

Termination of Resolution

The application of the resolution rule to automata clauses (reg) does not terminate.

$$\frac{P_1(x_1), P_2(x_2) \Rightarrow Q_1(g(x_1, x_2)) \quad Q_1(y_1), Q_2(y_2) \Rightarrow Q(f(y_1, y_2))}{P_1(x_1), P_2(x_2), Q_2(y_2) \Rightarrow Q(f(g(x_1, x_2), y_2))}$$

Complete Strategies for Resolution

$$\frac{\begin{array}{c} (C) \\ P_1(s_1), \dots, P_m(s_m) \Rightarrow Q_1(s) \end{array} \quad \begin{array}{c} (D) \\ Q_1(t_1), \dots, Q_n(t_n) \Rightarrow Q(t) \end{array}}{P_1(s_1\sigma), \dots, P_m(s_m\sigma), Q_1(t_1\sigma), \dots, Q_n(t_n\sigma) \Rightarrow Q(t\sigma)}$$

ordered resolution for \succ :

- ▶ $Q_1(s)$ maximal for \succ in C ,
- ▶ $Q_1(t_1)$ maximal for \succ in D .

ordered resolution with selection :

selection function : clause \mapsto subset of negative literals.

- ▶ no literal is selected in C ,
- ▶ $Q_1(s)$ is maximal for \succ in C ,
- ▶ $Q_1(t_1)$ is selected in D or
- ▶ no literal is selected in D and $Q_1(t)$ is maximal in D .

Completeness of Ordered Resolution

Theorem :

Ordered resolution with selection is complete for Horn clauses.

Starting from any unsatisfiable set S , we shall derive \perp .

Choice of an Ordered Strategy with Selection

- ▶ ordering \succ s.t. $P(s) \succ Q(t)$ iff $s > t$ for the subterm ordering $>$.
- ▶ selection function sel :
 - ▶ negative literals $\neg Q(t)$ where t is not a variable.

Lemma :

Every tree automaton (finite set clauses of type (reg)) is saturated under resolution ordered by \succ .

Resolution ordered by \succ and with selection by sel cannot be applied between automata clauses (reg) like in

$$\frac{P_1(x_1), \dots, P_m(x_m) \Rightarrow Q_1(g(\bar{x})) \quad Q_1(y_1), \dots, Q_n(y_n) \Rightarrow Q(f(\bar{y}))}{P_1(x_1), \dots, P_m(x_m), Q_2(y_2), \dots, Q_n(y_n) \Rightarrow Q(f(g(\bar{x}), y_2, \dots, y_n))}$$

because no literal are selected in the clauses and $Q_1(y_1)$ is not maximal in $\{Q_1(y_1), \dots, Q_n(y_n), Q(f(\bar{y}))\}$ ($Q(f(\bar{y})) \succ Q_1(y_1)$).

Transformation of Alternating Automata (reg + alt = reg)

Alternating automata = finite set of clauses

$$Q_1(x_1), \dots, Q_n(x_n) \Rightarrow Q(f(x_1, \dots, x_n)) \quad (\text{reg})$$

and

$$Q_1(x), \dots, Q_n(x) \Rightarrow Q(x) \quad (\text{alt})$$

Proposition :

Given an alternating tree automaton \mathcal{A} over Σ , we can construct in exponential time a deterministic bottom-up tree automaton \mathcal{A}' recognizing the same language.

Construction by application of ordered resolution with selection, following an appropriated strategy.

Transformation of Alternating Automata

- ▶ start from a set \mathcal{A} of clauses of the form (reg) and (alt).
- ▶ saturate with resolution ordered (by \succ) with selection (by *sel*).
- ▶ all the clauses produced belong to a type containing an exponential number of clauses

$$Q_1(x_{i_1}), \dots, Q_k(x_{i_k}), \underline{Q'_1(f(\bar{x}))}, \dots, \underline{Q'_m(f(\bar{x}))} \Rightarrow Q(f(\bar{x})) \quad (\text{f})$$

- ▶ hence saturation terminates with a set \mathcal{A}''
- ▶ the application of resolution to $\mathcal{A}'' \cup \{Q(t) \Rightarrow \perp\}$ (for a ground term t) only involves clauses of the form (reg).
- ▶ hence, for all Q , $L(\mathcal{A}'', Q) = L(\mathcal{A}''|_{\text{reg}}, Q)$.

Transformation of Bidirectional Alternating Automata (reg + alt + bidi = reg)

Bidirectional (2-way) alternating automata = finite set of clauses

$$Q_1(x_1), \dots, Q_n(x_n) \Rightarrow Q(f(x_1, \dots, x_n)) \quad (\text{reg})$$

and

$$Q_1(x), \dots, Q_n(x) \Rightarrow Q(x) \quad (\text{alt})$$

and

$$Q(f(x_1, \dots, x_n)) \Rightarrow Q_i(x_i) \quad (\text{bidi})$$

Proposition :

Given a bidirectional alternating tree automaton \mathcal{A} over Σ , we can construct in exponential time a bottom-up deterministic tree automaton \mathcal{A}' recognizing the same language.

same principle as for alternating tree automata, with

- ▶ other ordering and selection function for defining the resolution strategy,
- ▶ and a new rule called ε -splitting.

Decision of Instance Membership

tree automaton

$$Q_1(x_1), \dots, Q_n(x_n) \Rightarrow Q(f(x_1, \dots, x_n)) \quad (\text{reg})$$

+ query

$$Q(t) \Rightarrow \perp$$

Resolution ordered by \succ with selection by *sel* and ε -splitting terminates. **invariant:** the resolution only produces clauses of the following 2 types:

$$P_1(s_1), \dots, P_m(s_m), q_1, \dots, q_k \Rightarrow [q] \quad (\text{gs})$$

where $m, k \geq 0$, and s_1, \dots, s_m are subterms of t .

$$P_1(y_{i_1}), \dots, P_k(y_{i_k}), \underline{P'_1(f(y_1, \dots, y_n))}, \dots, \underline{P'_m(f(\bar{y}))} \Rightarrow [q] \quad (\text{gf})$$

where $k, m \geq 0$, $k + m > 0$, $i_1, \dots, i_k \leq n$, and y_1, \dots, y_n distinct.

Plan

Definition

Saturation Results

Tree Automata with Equality Constraints

Tree Automata with One Memory

Testing Equalities between Brother Subterms

In standard tree automata clauses, the variables x_1, \dots, x_n are pairwise distinct

$$Q_1(x_1), \dots, Q_n(x_n) \Rightarrow Q(f(x_1, \dots, x_n)) \quad (\text{reg})$$

With variable sharing in

$$Q_1(x_1), \dots, Q_n(x_n) \Rightarrow Q(f(x_1, \dots, x_n)) \quad (\text{brother})$$

we force equalities between brother subterms.

example:

$$\begin{aligned} & \Rightarrow Q(a), \\ Q(x_1), Q(x_2) & \Rightarrow Q(f(x_1, x_2)), \\ Q(x), Q(x) & \Rightarrow Q_f(f(x, x)) \end{aligned}$$

Testing Equality between Brother Subterms

Tree automata with equality tests between brother subterms are strictly more expressive than bottom-up tree automata.

Theorem :

The emptiness problem is EXPTIME-complete for tree automata with equality tests between brother subterms.

Testing Equalities and Disequalities between Brother Subterms

[Bogaert Tison STACS 1992]: tests of $=$ and \neq between brother subterms (see chapter 4 TATA book).

- ▶ determinizable in exponential time
- ▶ all Boolean closures
- ▶ emptiness decidable in PTIME for deterministic
- ▶ emptiness EXPTIME-complete for non-deterministic

[Reuss Seidl 2010]: clausal presentation with \neq .

Testing Arbitrary Equalities

One can perform more general tests with clauses with equalities

$$Q_1(x_1), \dots, Q_n(x_n), u_1 = v_1, \dots, u_k = v_k \Rightarrow Q(f(x_1, \dots, x_n))$$

(test)

where $k \geq 0$, $u_1, v_1, \dots, u_k, v_k \in \mathcal{T}(\Sigma, \{x_1, \dots, x_n\})$.

Without restrictions, emptiness is **undecidable**.

Arbitrary Equality Tests: Decidable Class

[JRV JLAP 08]

We distinguish some predicates call **test predicates**, and assume a partial ordering \succ over predicates such that $Q \succ Q_0$ for all Q , test predicate and Q_0 non-test.

$$Q_1(x_1), \dots, Q_n(x_n), u_1 = v_1, \dots, u_k = v_k \Rightarrow Q(f(x_1, \dots, x_n)) \quad (\text{test})$$

where Q is a test predicate, and for all $i \leq n$, $Q \succ Q_i$

$$Q_1(x_1), \dots, Q_n(x_n) \Rightarrow Q(f(x_1, \dots, x_n)) \quad (\text{reg}'')$$

where either all Q, Q_1, \dots, Q_n are not test predicates, or Q is a test predicate and at most one $Q_i = Q$ and the others are not test.

Arbitrary equality tests: decidable class

example: stuttering lists

$$\begin{array}{l} \Rightarrow Q_0(0) \qquad \qquad Q_0(x) \Rightarrow Q_0(s(x)) \\ \Rightarrow Q_1(\text{nil}) \quad Q_0(x), Q_1(y) \Rightarrow Q_1(\text{cons}(x, y)) \\ \qquad \qquad \qquad Q_0(x), Q_2(y) \Rightarrow Q_2(\text{cons}(x, y)) \end{array}$$

$$Q_0(x), Q_1(y), y = \text{cons}(x, y') \Rightarrow Q_2(\text{cons}(x, y))$$

Arbitrary Equality Tests: Decidable Class

Theorem :

The satisfiability of a set of clauses of type (test) and (reg') and a goal clause $Q(t) \Rightarrow \perp$ is decidable.

pr.: Saturation by ordered paramodulation with selection and ϵ -splitting.

Extension to langages (and equality tests) modulo equational theories, by adding clauses

$$\Rightarrow \ell = x$$

of a restricted form.

Example:

$$\text{car}(\text{cons}(x, y)) = x, \text{cdr}(\text{cons}(x, y)) = y, \text{cons}(\text{car}(y), \text{cdr}(y)) = y$$

Plan

Definition

Saturation Results

Tree Automata with Equality Constraints

Tree Automata with One Memory

Pushdown Tree Automata with One Tree Memory

[Guessarian 83, Schimpf Gallier 85, Coquidé et al 94]

- ▶ Σ : input signature; on input: terms of $\mathcal{T}(\Sigma)$
- ▶ Γ : stack alphabet; the auxiliary memory is a stack of Γ^*

$$Q_1(x_1, s_1), \dots, Q_n(x_n, s_n) \Rightarrow Q(f(\bar{x}), h(y_1, \dots, y_m)) \quad (\text{read})$$

where $f \in \Sigma$, $h \in \Gamma$, $s_1, \dots, s_n \in \mathcal{T}(\Gamma, \{y_1, \dots, y_m\})$

$$Q_1(x, s) \Rightarrow Q(x, h(y_1, \dots, y_m)) \quad (\text{pda-}\varepsilon)$$

where $h \in \Gamma$, $s \in \mathcal{T}(\Gamma, \{y_1, \dots, y_m\})$.

Same expressiveness as [context-free tree grammars](#)

$$X(x_1, \dots, x_n) \rightarrow r$$

where $X \in \mathbb{N}$, non-terminal of arity n , x_1, \dots, x_n distinct variables, $r \in \mathcal{T}(\Sigma \cup \mathbb{N}, \{x_1, \dots, x_n\})$

Tree Automata with One Memory

[Comon Cortier 05]

$$Q_1(x_1, y_1), Q_2(x_2, y_2) \Rightarrow Q(f(x_1, x_2), h(y_1, y_2)) \quad (\text{push})$$

$$Q_1(x_1, h(y_{11}, y_{12})), Q_2(x_2, y_2) \Rightarrow Q(f(x_1, x_2), y_{11}) \quad (\text{pop}_{11})$$

$$Q_1(x_1, \perp), Q_2(x_2, y_2) \Rightarrow Q(f(x_1, x_2), \perp)$$

$$\dots (\text{pop}_{12}), (\text{pop}_{21}), (\text{pop}_{22})$$

$$\Rightarrow Q(a, \perp) \quad (\text{int}_0)$$

$$Q_1(x_1, y_1), Q_2(x_2, y_2) \Rightarrow Q(f(x_1, x_2), y_1) \quad (\text{int}_1)$$

$$Q_1(x_1, y_1), Q_2(x_2, y_2) \Rightarrow Q(f(x_1, x_2), y_2) \quad (\text{int}_2)$$

where $h \in \Gamma_2$, $a \in \Sigma_0$, $f \in \Sigma_2$.

Tree Automata with One Memory: Languages and Closure

(input) Language

$$L(\mathcal{A}, Q) = \{t \mid Q(t, s) \in \mathcal{H}_S\}$$

Memory Language

$$M(\mathcal{A}, Q) = \{s \mid Q(t, s) \in \mathcal{H}_S\}$$

Languages of tree automata with one memory are closed under \cup but not under \cap or \neg .

[CJP FOSSACS 07] addition of a **visibly** condition:

- ▶ $\Sigma = \Sigma_{push} \uplus \Sigma_{pop_{11}} \uplus \dots \uplus \Sigma_{int_0} \uplus \Sigma_{int_1} \uplus \Sigma_{int_2}$
- ▶ the input symbol determines the possible operations on memory
- ▶ full Boolean closure

Tree Automata with One Memory: Languages and Closure

Theorem :

The emptiness problem is decidable in polynomial time for tree automata with one memory.

Lemma :

For all \mathcal{A}, Q , $L(\mathcal{A}, Q) = \emptyset$ iff $M(\mathcal{A}, Q) = \emptyset$.

- ▶ For all \mathcal{A}, Q , $M(\mathcal{A}, q)$ is a language of **bidirectional alternating automata** (clauses $\text{reg} + \text{alt} + \text{bidi} = \text{reg}$).
- ▶ definition by "*projection*" of clauses on second component of states.
- ▶ actually $M(\mathcal{A}, q)$ in a smaller class \mathcal{H}_3 [Nielson Seidl SAS 02] decidable in cubic time.

Tree Automata with One Memory and Constraints

extension with constraints testing the content of the memory.

$$Q_1(x_1, y_1), Q_2(x_2, y_2), y_1 = y_2 \Rightarrow Q(f(x_1, x_2), y_1) \quad (\text{int}_1^-)$$

$$Q_1(x_1, y_1), Q_2(x_2, y_2), y_1 = y_2 \Rightarrow Q(f(x_1, x_2), y_2) \quad (\text{int}_2^-)$$

also \neq tests in the model of [CJP FOSSACS 07].

- ▶ emptiness is decidable
- ▶ Boolean closure with visibly condition and restriction of constraints to structural equality.

TA1M for Verification

Tree automata with one memory and constraints can recognize the following data-structures

- ▶ balanced binary trees
- ▶ powerlists
(description and verification of data parallel algorithms)
- ▶ red-black trees (binary search trees)
 1. every node is black or red
 2. the root is black
 3. all the leaves are black
 4. the 2 children of a red node are black
 5. all paths have the same number of black nodes