

SAT Solvers for Queries over Tree Automata with Constraints

Pierre-Cyrille Héam, [Vincent Hugot](#), Olga Kouchnarenko
{pcheam,okouchnarenko}@lifc.univ-fcomte.fr,
vhugot@edu.univ-fcomte.fr

Université de Franche-Comté
LIFC-INRIA/CASSIS

April 8, 2010

Plan of the talk

1 Motivating XML example

Plan of the talk

- ① **Motivating XML example**
- ② **Introduction of notions:**
 - ① Tree Automata
 - ② TAGEDs
 - ③ SAT problem

Plan of the talk

- ① **Motivating XML example**
- ② **Introduction of notions:**
 - ① Tree Automata
 - ② TAGEDs
 - ③ SAT problem
- ③ **Main contribution:**
SAT encoding for TAGED Uniform Membership Problem

Plan of the talk

① Motivating XML example

② Introduction of notions:

- ① Tree Automata
- ② TAGEDs
- ③ SAT problem

③ Main contribution:

SAT encoding for TAGED Uniform Membership Problem

④ Some experimental results:

- ① Natural optimisations
- ② The prototype
- ③ Conversion to CNF

Plan of the talk

- ① **Motivating XML example**
- ② **Introduction of notions:**
 - ① Tree Automata
 - ② TAGEDs
 - ③ SAT problem
- ③ **Main contribution:**

SAT encoding for TAGED Uniform Membership Problem
- ④ **Some experimental results:**
 - ① Natural optimisations
 - ② The prototype
 - ③ Conversion to CNF
- ⑤ **Conclusion.**

A small example

Laboratory toy example

```
<university>
  <team>
    <member> Scotty </member>
    <member> Spock </member>
    <member> Uhura </member>
    <laboratory> Enterprise </laboratory>
  </team>
  <team>
    <member> McCoy </member>
    <member> Spock </member>
    <laboratory> Enterprise </laboratory>
  </team>
</university>
```

Objective: check that all teams belong to the same laboratory and no researcher is affected to two different teams.

A small example

Laboratory toy example

```
<university>
  <team>
    <member> Scotty </member>
    <member> Spock </member>
    <member> Uhura </member>
    <laboratory> Enterprise </laboratory>
  </team>
  <team>
    <member> McCoy </member>
    <member> Spock </member>
    <laboratory> Enterprise </laboratory>
  </team>
</university>
```

Objective: check that **all teams belong to the same laboratory** and no researcher is affected to two different teams.

A small example

Laboratory toy example

```
<university>
  <team>
    <member> Scotty </member>
    <member> Spock </member>
    <member> Uhura </member>
    <laboratory> Enterprise </laboratory>
  </team>
  <team>
    <member> McCoy </member>
    <member> Spock </member>
    <laboratory> Enterprise </laboratory>
  </team>
</university>
```

Objective: check that **all teams belong to the same laboratory**
and **no researcher is affected to two different teams**.

Tree automata

Definition through an example

Tree automaton for True propositional formulæ

$$\mathcal{A} \stackrel{\text{def}}{=} (\Sigma = \{ \wedge, \vee/2, \neg/1, 0, 1/0 \}, Q = \{ q_0, q_1 \}, F = \{ q_1 \}, \Delta)$$

$$\begin{aligned} \Delta = \{ & b \rightarrow q_b, \\ & \wedge (q_b, q_{b'}) \rightarrow q_{b \wedge b'}, \\ & \vee (q_b, q_{b'}) \rightarrow q_{b \vee b'}, \\ & \neg(q_b) \rightarrow q_{\neg b} \\ & \mid b, b' \in 0, 1 \} \end{aligned}$$

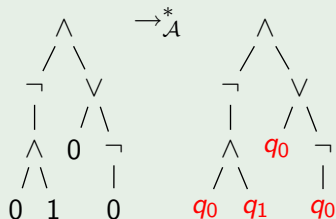
Tree automata

Definition through an example



Tree automata

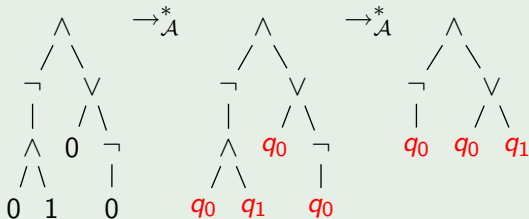
Definition through an example

$$0 \rightarrow q_0, 1 \rightarrow q_1 \in \Delta$$


Tree automata

Definition through an example

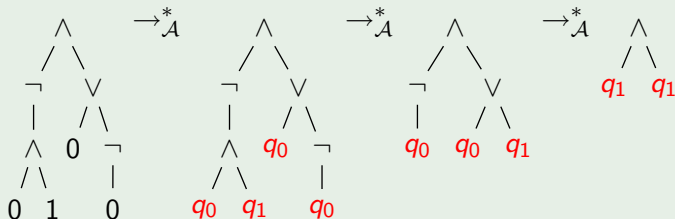
$$\wedge(q_0, q_1) \rightarrow q_0, \neg(q_0) \rightarrow q_1 \in \Delta$$



Tree automata

Definition through an example

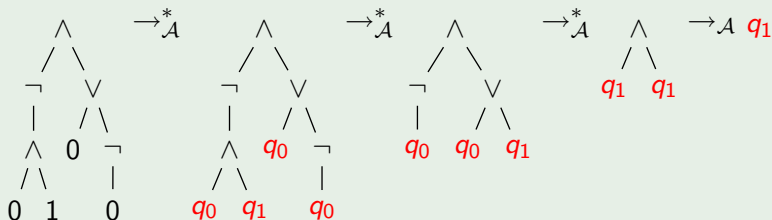
$$\neg(q_0) \rightarrow q_1, \forall(q_0, q_1) \rightarrow q_1 \in \Delta$$



Tree automata

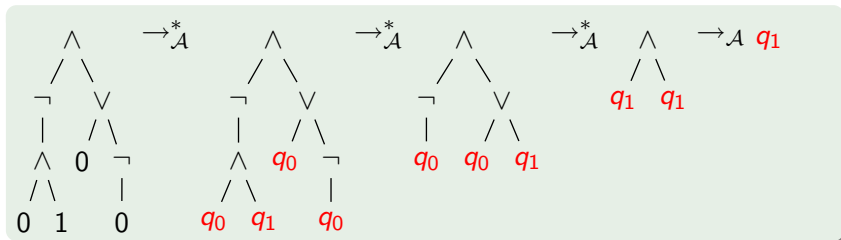
Definition through an example

$$\wedge(q_1, q_1) \rightarrow q_1 \in \Delta$$



Tree automata

Definition through an example

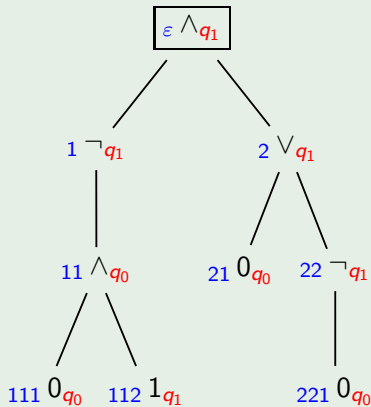


Definition: run of \mathcal{A} on a term $t \in \mathcal{T}(\Sigma)$

A run ρ is a mapping from $\mathcal{Pos}(t)$ to Q compatible with the transition rules.

Tree automata

Definition through an example

 $\rho =$ 

TAGEDs

Tree Automata With Global Equality and Disequality Constraints

Introduced in Emmanuel Filiot's PhD thesis on XML query languages. See [Filiot et al., 2008].

A TAGED is a tuple $\mathcal{A} = (\Sigma, Q, F, \Delta, =_{\mathcal{A}}, \neq_{\mathcal{A}})$, where

- (Σ, Q, F, Δ) is a tree automaton
- $=_{\mathcal{A}}$ is a reflexive symmetric binary relation on a subset of Q
- $\neq_{\mathcal{A}}$ is an irreflexive and symmetric binary relation on Q . Note that in our work, we have dealt with a slightly more general case, where $\neq_{\mathcal{A}}$ is not necessarily irreflexive.

A TAGED \mathcal{A} is said to be *positive* if $\neq_{\mathcal{A}}$ is empty and *negative* if $=_{\mathcal{A}}$ is empty.

Runs must be compatible with equality and disequality constraints.

TAGEDs

Compatibility with global constraints

Let ρ be a run of the TAGED \mathcal{A} on a tree t :

Compatibility with the equality constraint $=_{\mathcal{A}}$

$$\forall \alpha, \beta \in \text{Pos}(t) : \rho(\alpha) =_{\mathcal{A}} \rho(\beta) \implies t|_{\alpha} = t|_{\beta}.$$

Compatibility with the disequality constraint $\neq_{\mathcal{A}}$ (irreflexive)

$$\forall \alpha, \beta \in \text{Pos}(t) : \rho(\alpha) \neq_{\mathcal{A}} \rho(\beta) \implies t|_{\alpha} \neq t|_{\beta}.$$

Compatibility with the disequality constraint $\neq_{\mathcal{A}}$ (non irreflexive)

$$\forall \alpha, \beta \in \text{Pos}(t) : \alpha \neq \beta \wedge \rho(\alpha) \neq_{\mathcal{A}} \rho(\beta) \implies t|_{\alpha} \neq t|_{\beta}.$$

TAGEDs

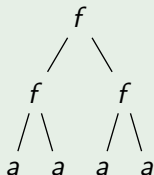
A non-regular language accepted by TAGEDs

TAGED for $\{f(t, t) \mid f \in \Sigma, t \in \mathcal{T}(\Sigma)\}$

$$\mathcal{A} \stackrel{\text{def}}{=} (\Sigma = \{a, f\}, Q = \{q, \hat{q}, q_f\}, F = \{q_f\},$$

$$\Delta, \hat{q} \Rightarrow_{\mathcal{A}} \hat{q}),$$

$$\text{where } \Delta \stackrel{\text{def}}{=} \{f(\hat{q}, \hat{q}) \rightarrow q_f, f(q, q) \rightarrow q, f(q, q) \rightarrow \hat{q}, \\ a \rightarrow q, a \rightarrow \hat{q}, \}$$



TAGEDs

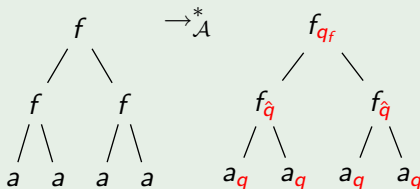
A non-regular language accepted by TAGEDs

TAGED for $\{f(t, t) \mid f \in \Sigma, t \in \mathcal{T}(\Sigma)\}$

$$\mathcal{A} \stackrel{\text{def}}{=} (\Sigma = \{a, f\}, Q = \{q, \hat{q}, q_f\}, F = \{q_f\},$$

$$\Delta, \hat{q} \Rightarrow_{\mathcal{A}} \hat{q}),$$

$$\text{where } \Delta \stackrel{\text{def}}{=} \{f(\hat{q}, \hat{q}) \rightarrow q_f, f(q, q) \rightarrow q, f(q, q) \rightarrow \hat{q}, \\ a \rightarrow q, a \rightarrow \hat{q}, \}$$



TAGEDs

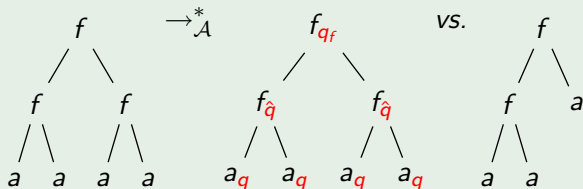
A non-regular language accepted by TAGEDs

TAGED for $\{f(t, t) \mid f \in \Sigma, t \in \mathcal{T}(\Sigma)\}$

$$\mathcal{A} \stackrel{\text{def}}{=} (\Sigma = \{a, f\}, Q = \{q, \hat{q}, q_f\}, F = \{q_f\},$$

$$\Delta, \hat{q} \Rightarrow_{\mathcal{A}} \hat{q}),$$

$$\text{where } \Delta \stackrel{\text{def}}{=} \{f(\hat{q}, \hat{q}) \rightarrow q_f, f(q, q) \rightarrow q, f(q, q) \rightarrow \hat{q}, \\ a \rightarrow q, a \rightarrow \hat{q}, \}$$



TAGEDs

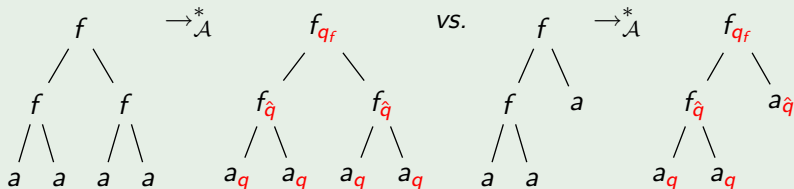
A non-regular language accepted by TAGEDs

TAGED for $\{f(t, t) \mid f \in \Sigma, t \in \mathcal{T}(\Sigma)\}$

$$\mathcal{A} \stackrel{\text{def}}{=} (\Sigma = \{a, f\}, Q = \{q, \hat{q}, q_f\}, F = \{q_f\},$$

$$\Delta, \hat{q} \Rightarrow_{\mathcal{A}} \hat{q}),$$

$$\text{where } \Delta \stackrel{\text{def}}{=} \{f(\hat{q}, \hat{q}) \rightarrow q_f, f(q, q) \rightarrow q, f(q, q) \rightarrow \hat{q}, \\ a \rightarrow q, a \rightarrow \hat{q}, \}$$



TAGED membership

... through SAT solvers

Uniform Membership Problem

INPUT: \mathcal{A} a TAGED and $t \in \mathcal{T}(\Sigma)$ a term.

OUTPUT: Is t accepted by \mathcal{A} ?

TAGED membership

... through SAT solvers

Uniform Membership Problem

INPUT: \mathcal{A} a TAGED and $t \in \mathcal{T}(\Sigma)$ a term.

OUTPUT: Is t accepted by \mathcal{A} ?

Theorem [Filiot2008]

The Uniform Membership Problem for TAGEDs is *NP*-complete.

TAGED membership

... through SAT solvers

Uniform Membership Problem

INPUT: \mathcal{A} a TAGED and $t \in \mathcal{T}(\Sigma)$ a term.

OUTPUT: Is t accepted by \mathcal{A} ?

Theorem [Filiot2008]

The Uniform Membership Problem for TAGEDs is *NP*-complete.

Possible practical approach

XML/... \Rightarrow TAGED membership

TAGED membership

... through SAT solvers

Uniform Membership Problem

INPUT: \mathcal{A} a TAGED and $t \in \mathcal{T}(\Sigma)$ a term.

OUTPUT: Is t accepted by \mathcal{A} ?

Theorem [Filiot2008]

The Uniform Membership Problem for TAGEDs is *NP*-complete.

Possible practical approach

XML/... \Rightarrow TAGED membership \Rightarrow SAT problem

TAGED membership

... through SAT solvers

Uniform Membership Problem

INPUT: \mathcal{A} a TAGED and $t \in \mathcal{T}(\Sigma)$ a term.

OUTPUT: Is t accepted by \mathcal{A} ?

Theorem [Filiot2008]

The Uniform Membership Problem for TAGEDs is *NP-complete*.

Possible practical approach

XML/... \Rightarrow TAGED membership \Rightarrow *SAT problem*

TAGED membership

... through SAT solvers

Uniform Membership Problem

INPUT: \mathcal{A} a TAGED and $t \in \mathcal{T}(\Sigma)$ a term.

OUTPUT: Is t accepted by \mathcal{A} ?

Theorem [Filiot2008]

The Uniform Membership Problem for TAGEDs is *NP*-complete.

Possible practical approach

XML/... \Rightarrow TAGED membership \Rightarrow SAT problem \Rightarrow *answer*

The SAT Problem

... and applications

Definition: The SAT problem

Given a propositional formula, for instance

$$\varphi = X \vee (\neg X \wedge \neg Y) \text{ or } \psi = X \wedge (\neg X \wedge \neg Y),$$

is there a valuation such that the formula evaluates to true?

NP-complete

The SAT problem is the first known *NP*-complete decision problem (Cook, 1971).

In practice

There are very efficient heuristics implanted in modern SAT solvers.

Plan

- 1 The SAT encoding
- 2 Implementation and Experiments

The SAT encoding

Choice of variables and Θ_{\rightarrow}

Variables: X_q^α

A run is a mapping from $\mathcal{Pos}(t)$ to Q . So we take variables of the form X_q^α , meaning:

$$\forall \alpha \in \mathcal{Pos}(t), q \in Q : X_q^\alpha \iff \rho(\alpha) = q$$

The SAT encoding

Choice of variables and Θ_{\rightarrow}

Variables: X_q^α

A run is a mapping from $\text{Pos}(t)$ to Q . So we take variables of the form X_q^α , meaning:

$$\forall \alpha \in \text{Pos}(t), q \in Q : X_q^\alpha \iff \rho(\alpha) = q$$

Partial function constraint Θ_{\rightarrow} : “ ρ is a function”

$$\Theta_{\rightarrow} \stackrel{\text{def}}{=} \bigwedge_{\substack{\alpha \in \text{Pos}(t) \\ q \in Q}} \left[X_q^\alpha \implies \bigwedge_{\substack{p \in Q \\ p \neq q}} \neg X_p^\alpha \right]$$

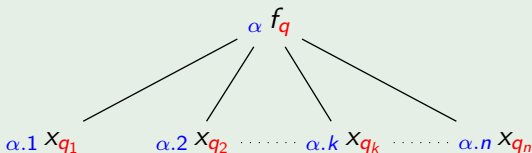
The SAT encoding

Rule application and compatibility: $\Psi^\alpha(r)$ and $\Phi^\varepsilon(t)$

Rule application constraint $\Psi^\alpha(r)$

We define, for any $\alpha \in \mathcal{Pos}(t)$, and any transition rule $f(q_1, \dots, q_n) \rightarrow q \in \Delta$,

$$\Psi^\alpha(f(q_1, \dots, q_n) \rightarrow q) \stackrel{\text{def}}{=} X_q^\alpha \wedge \bigwedge_{k=1}^n X_{q_k}^{\alpha.k}.$$



The SAT encoding

Rule application and compatibility: $\Psi^\alpha(r)$ and $\Phi^\varepsilon(t)$

Rule application constraint $\Psi^\alpha(r)$

We define, for any $\alpha \in \mathcal{Pos}(t)$, and any transition rule $f(q_1, \dots, q_n) \rightarrow q \in \Delta$,

$$\Psi^\alpha(f(q_1, \dots, q_n) \rightarrow q) \stackrel{\text{def}}{=} X_q^\alpha \wedge \bigwedge_{k=1}^n X_{q_k}^{\alpha.k}.$$

Rules compatibility constraint $\Phi^\varepsilon(t)$

$$\Phi^\varepsilon(t) \stackrel{\text{def}}{=} \bigwedge_{\alpha \in \mathcal{Pos}(t)} \left[\bigvee_{r \in \Delta_{t(\alpha)}} \Psi^\alpha(r) \right]$$

where $\Delta_f = \{ f(\dots) \rightarrow \dots \in \Delta \}$.

The SAT encoding

Connecting to TAGEDs: Θ_{\Leftarrow}

Accepting run for tree automata

We have coded runs for tree automata: one more constraint:

$\bigvee_{q \in F} X_q^\varepsilon$ makes sure we end up in a final state.

The SAT encoding

Connecting to TAGEDs: Θ_{\Leftarrow}

Accepting run for tree automata

We have coded runs for tree automata: one more constraint:

$\bigvee_{q \in F} X_q^\varepsilon$ makes sure we end up in a final state.

New variables: T_u^q

We need new variables to link subterms and states: T_u^q denotes “the subterm u evaluates to q ”, for any $u \trianglelefteq t$ and $q \in Q$.

The SAT encoding

Connecting to TAGEDs: Θ_{\Leftarrow}

Accepting run for tree automata

We have coded runs for tree automata: one more constraint:

$\bigvee_{q \in F} X_q^\varepsilon$ makes sure we end up in a final state.

New variables: T_u^q

We need new variables to link subterms and states: T_u^q denotes “the subterm u evaluates to q ”, for any $u \trianglelefteq t$ and $q \in Q$.

Structural glue: Θ_{\Leftarrow}

$$\Theta_{\Leftarrow} \stackrel{\text{def}}{=} \bigwedge_{\substack{\alpha \in \text{Pos}(t) \\ q \in Q}} [X_q^\alpha \implies T_{t|_\alpha}^q].$$

The SAT encoding

Compatibility with \Rightarrow_A and \neq_A

Compatibility with \Rightarrow_A : Θ_{\Rightarrow_A}

$$\Theta_{\Rightarrow_A} \stackrel{\text{def}}{=} \bigwedge_{\substack{\alpha \in \text{Pos}(t) \\ q \in Q}} \left[X_q^\alpha \implies \bigwedge_{\substack{p \in Q \\ p \Rightarrow_A q}} \bigwedge_{\substack{u \triangleleft t \\ u \neq t|_\alpha}} \neg T_u^p \right]$$

Compatibility with \neq_A ($p \neq q$): Θ_{\neq_A}

$$\Theta_{\neq_A} \stackrel{\text{def}}{=} \bigwedge_{\substack{\alpha \in \text{Pos}(t) \\ q \in Q}} \left[X_q^\alpha \implies \bigwedge_{\substack{p \in Q \\ p \neq_A q \\ p \neq q}} \neg T_{t|_\alpha}^p \right]$$

The SAT encoding

Compatibility with \Rightarrow_A and $\not\Rightarrow_A$

Compatibility with $\not\Rightarrow_A$ with structural glue Θ_{\Leftrightarrow}

$$\Theta_{\not\Rightarrow_A} \wedge \Theta_{\Leftrightarrow} \stackrel{\text{def}}{=} \bigwedge_{\substack{\alpha \in \text{Pos}(t) \\ q \in Q}} \left[X_q^\alpha \Rightarrow T_{t|_\alpha}^q \wedge \bigwedge_{\substack{p \in Q \\ p \not\Rightarrow_A q \\ p \neq q}} \neg T_{t|_\alpha}^p \right]$$

The idea: a counterexample

Let $t \in \mathcal{T}(\Sigma)$, $\alpha, \beta \in \text{Pos}(t)$ and $u = t|_\alpha = t|_\beta$. Suppose that ρ is a run such that $\rho(\alpha) = p$ and $\rho(\beta) = q$ with $p \not\Rightarrow_A q$. Then we have X_p^α, X_q^β ,

The SAT encoding

Compatibility with \Rightarrow_A and \neq_A

Compatibility with \neq_A with structural glue Θ_{\Leftrightarrow}

$$\Theta_{\neq_A} \wedge \Theta_{\Leftrightarrow} \stackrel{\text{def}}{=} \bigwedge_{\substack{\alpha \in \text{Pos}(t) \\ q \in Q}} \left[X_q^\alpha \Rightarrow T_{t|_\alpha}^q \wedge \bigwedge_{\substack{p \in Q \\ p \neq_A q \\ p \neq q}} \neg T_{t|_\alpha}^p \right]$$

The idea: a counterexample

Let $t \in \mathcal{T}(\Sigma)$, $\alpha, \beta \in \text{Pos}(t)$ and $u = t|_\alpha = t|_\beta$. Suppose that ρ is a run such that $\rho(\alpha) = p$ and $\rho(\beta) = q$ with $p \neq_A q$. Then we have $X_p^\alpha, X_q^\beta, T_{t|_\alpha}^p, \neg T_{t|_\alpha}^q$,

The SAT encoding

Compatibility with \Rightarrow_A and \neq_A

Compatibility with \neq_A with structural glue Θ_{\Leftarrow}

$$\Theta_{\neq_A} \wedge \Theta_{\Leftarrow} \stackrel{\text{def}}{=} \bigwedge_{\substack{\alpha \in \text{Pos}(t) \\ q \in Q}} \left[X_q^\alpha \Rightarrow T_{t|_\alpha}^q \wedge \bigwedge_{\substack{p \in Q \\ p \neq_A q \\ p \neq q}} \neg T_{t|_\alpha}^p \right]$$

The idea: a counterexample

Let $t \in \mathcal{T}(\Sigma)$, $\alpha, \beta \in \text{Pos}(t)$ and $u = t|_\alpha = t|_\beta$. Suppose that ρ is a run such that $\rho(\alpha) = p$ and $\rho(\beta) = q$ with $p \neq_A q$. Then we have $X_p^\alpha, X_q^\beta, T_{t|_\alpha}^p, \neg T_{t|_\alpha}^q, T_{t|_\beta}^q, \neg T_{t|_\beta}^p$.

The SAT encoding

Compatibility with \Rightarrow_A and \neq_A

Compatibility with \neq_A with structural glue Θ_{\Leftarrow}

$$\Theta_{\neq_A} \wedge \Theta_{\Leftarrow} \stackrel{\text{def}}{=} \bigwedge_{\substack{\alpha \in \text{Pos}(t) \\ q \in Q}} \left[X_q^\alpha \Rightarrow T_{t|_\alpha}^q \wedge \bigwedge_{\substack{p \in Q \\ p \neq_A q \\ p \neq q}} \neg T_{t|_\alpha}^p \right]$$

The idea: a counterexample

Let $t \in \mathcal{T}(\Sigma)$, $\alpha, \beta \in \text{Pos}(t)$ and $u = t|_\alpha = t|_\beta$. Suppose that ρ is a run such that $\rho(\alpha) = p$ and $\rho(\beta) = q$ with $p \neq_A q$. Then we have $X_p^\alpha, X_q^\beta, T_{t|_\alpha}^p, \neg T_{t|_\alpha}^q, T_{t|_\beta}^q, \neg T_{t|_\beta}^p$. Since $u = t|_\alpha = t|_\beta$ we have $T_u^q, \neg T_u^q, T_u^p, \neg T_u^p$, hence the formula is not satisfiable.

The SAT encoding

Compatibility with \neq_A , non-irreflexive case: Ω_{\neq_A}

Compatibility with $\neq_A (p \neq q)$: Θ_{\neq_A}

$$\Theta_{\neq_A} \stackrel{\text{def}}{=} \bigwedge_{\substack{\alpha \in \text{Pos}(t) \\ q \in Q}} \left[X_q^\alpha \implies \bigwedge_{\substack{p \in Q \\ p \neq_A q \\ p \neq q}} \neg T_{t|_\alpha}^p \right]$$

New variables: S_u^α

We need new variables to link subterms and positions: S_u^α encodes the statement “the subterm u is rooted in α ”.

Compatibility with \neq_A (non-irreflexive; $q \neq_A q$): Ω_{\neq_A}

$$\Omega_{\neq_A} \stackrel{\text{def}}{=} \bigwedge_{\alpha \in \text{Pos}(t)} S_{t|_\alpha}^\alpha \wedge \bigwedge_{\substack{\alpha \neq \beta \in \text{Pos}(t) \\ q \neq_A q}} \left[X_q^\alpha \wedge X_q^\beta \implies \neg S_{t|_\beta}^\alpha \right]$$

The SAT encoding

Completeness and soundness

Definition (SAT encoding of TAGED membership problem $\Delta_{\mathcal{A}}(t)$)

Let $\mathcal{A} = (\Sigma, \Delta, Q, F, =_{\mathcal{A}}, \neq_{\mathcal{A}})$ be a TAGED and $t \in \mathcal{T}(\Sigma)$; then we define

$$\Delta_{\mathcal{A}}(t) \stackrel{\text{def}}{=} \Theta_{\rightarrow} \wedge \Phi^{\varepsilon}(t) \wedge \bigvee_{q \in F} X_q^{\varepsilon} \wedge \Theta_{=_{\mathcal{A}}} \wedge \Theta_{\neq_{\mathcal{A}}} \wedge \Omega_{\neq_{\mathcal{A}}}.$$

Theorem (TAGED membership, correctness and soundness)

There exists a successful run ρ of the TAGED \mathcal{A} on a term t iff $\Delta_{\mathcal{A}}(t)$ is satisfiable. Moreover, if $I \models \Delta_{\mathcal{A}}(t)$, then for any $\alpha \in \text{Pos}(t)$ we have $\rho(\alpha) = q \iff I \models X_q^{\alpha}$.

Plan

- 1 The SAT encoding
- 2 Implementation and Experiments

Possible optimisations

... from simple observations

The formulæ can be trimmed down: consider

Structural glue: Θ_{\Leftrightarrow}

$$\Theta_{\Leftrightarrow} \stackrel{\text{def}}{=} \bigwedge_{\substack{\alpha \in \text{Pos}(t) \\ q \in Q}} [X_q^\alpha \implies T_{t|_\alpha}^q].$$

Not all couples (α, q) are necessary because not all states are obtainable at any given position.

Definition (Possibly obtainable states at position α)

$$\sigma(\alpha) \stackrel{\text{def}}{=} \{q \in Q / \exists t(\alpha)(\dots) \rightarrow q \in \Delta\}$$

Given a position α , we only need to deal with $q \in \sigma(\alpha)$.

Possible optimisations

... from simple observations

The formulæ can be trimmed down: consider

Structural glue: Θ_{\Leftarrow}

$$\Theta_{\Leftarrow} \stackrel{\text{def}}{=} \bigwedge_{\substack{\alpha \in \text{Pos}(t) \\ q \in \sigma(\alpha)}} [X_q^\alpha \implies T_{t|_\alpha}^q].$$

Not all couples (α, q) are necessary because not all states are obtainable at any given position.

Definition (Possibly obtainable states at position α)

$$\sigma(\alpha) \stackrel{\text{def}}{=} \{q \in Q / \exists t(\alpha)(\dots) \rightarrow q \in \Delta\}$$

Given a position α , we only need to deal with $q \in \sigma(\alpha)$.

The prototype

Intro/ Input format

Prototype (0Cam1). Takes a TAGED and a term as input (syntax close to the Tree Automata library *Timbuk*).

```

(** TAGED Automaton for {f(x,x)} *)
  Taged fxxA
  Alphabet f a b
  States q qq qf
  Final qf
  Rules
    f qq qq : qf
    f q q   : q
    f q q   : qq
    a:q a:qq
    b:q b:qq
  Equal
    qq qq
  Different
    qq qf

```

f(f(a,a), f(a,a))
// in a_fxx

The prototype

Generated formula

Automaton: $fixA$

Alphabet: $\{f, a, b\}$

States: $\{q, qq, qf\}$

Final States: $\{qf\}$

Transition Rules: $\{$

$f(qq, qq) \rightarrow qf,$

$f(q, q) \rightarrow q,$

$f(q, q) \rightarrow qq,$

$a \rightarrow q,$

$a \rightarrow qq,$

$b \rightarrow q,$

$b \rightarrow qq$

$\}$

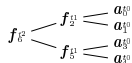
Global State Equality: $\{qq = qq\}$

Global State Disequality: $\{qq \neq qf\}$

End Automaton.

Term as expression: $f\{f[a, a], f[a, a]\}$

Term as tree:



Membership formula = $[[([X_q^3 \vee X_{qq}^3]) \wedge ([X_{qf}^5 \wedge X_{qq}^3 \wedge X_q^4] \vee [X_q^5 \wedge X_q^3 \wedge X_q^4] \vee [X_{qq}^5 \wedge X_q^3 \wedge X_q^4]) \wedge ([X_{qf}^2 \wedge X_{qq}^0 \wedge X_{qq}^1] \vee [X_q^2 \wedge X_q^0 \wedge X_q^1] \vee [X_{qq}^2 \wedge X_q^0 \wedge X_q^1]) \wedge (X_q^1 \vee X_{qq}^1) \wedge ([X_{qf}^6 \wedge X_{qq}^2 \wedge X_{qq}^5] \vee [X_q^6 \wedge X_q^2 \wedge X_q^5] \vee [X_{qq}^6 \wedge X_q^2 \wedge X_q^5]) \wedge (X_q^0 \vee X_{qq}^0) \wedge (X_q^4 \vee X_{qq}^4)] \wedge [\neg X_q^6 \wedge \neg X_{qq}^6] \wedge \{ \{ X_{qq}^3 \Rightarrow T_0^{qq1} \} \wedge \{ X_{qq}^5 \Rightarrow T_1^{qq1} \} \wedge \{ X_{qq}^2 \Rightarrow T_1^{qq1} \} \wedge \{ X_{qq}^1 \Rightarrow T_0^{qq1} \} \wedge \{ X_{qq}^6 \Rightarrow T_2^{qq1} \} \wedge \{ X_{qq}^0 \Rightarrow T_0^{qq1} \} \wedge \{ X_q^4 \Rightarrow T_0^{qq1} \} \wedge \{ X_{qq}^5 \Rightarrow \neg T_1^{qq1} \} \wedge \{ X_{qq}^2 \Rightarrow \neg T_1^{qq1} \} \wedge \{ X_{qq}^1 \wedge \{ X_{qq}^6 \Rightarrow \neg T_2^{qq1} \} \wedge \{ X_q^3 \Rightarrow [\neg T_2^{qq1} \wedge \neg T_1^{qq1}] \} \wedge \{ X_{qq}^5 \Rightarrow [\neg T_2^{qq1} \wedge \neg T_0^{qq1}] \} \wedge \{ X_{qq}^2 \Rightarrow [\neg T_2^{qq1} \wedge \neg T_0^{qq1}] \} \wedge \{ X_q^1 \Rightarrow [\neg T_2^{qq1} \wedge \neg T_1^{qq1}] \} \wedge \{ X_{qq}^6 \Rightarrow [\neg T_1^{qq1} \wedge \neg T_0^{qq1}] \} \wedge \{ X_{qq}^0 \Rightarrow [\neg T_2^{qq1} \wedge \neg T_1^{qq1}] \} \wedge \{ X_q^4 \Rightarrow [\neg T_2^{qq1} \wedge \neg T_1^{qq1}] \}]]]$

The prototype

CNF Conversion: the BAT

Definition: Conjunctive Normal Form

A formula is in Conjunctive Normal Form (CNF) if it is a conjunction of disjunctions of literals and contains only \neg , \wedge or \vee .

DIMACS CNF

$$\varphi = (X \vee \neg Z) \wedge (Y \vee Z \vee \neg X).$$

c DIMACS CNF for φ

p cnf 3 2

1 -3 0

2 3 -1 0

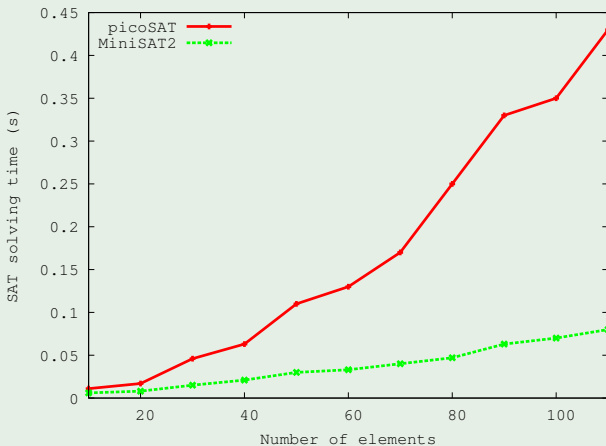
Problem!

Our formula is not in CNF! We must convert it. The BAT solves it.

The prototype

Results: Laboratory example

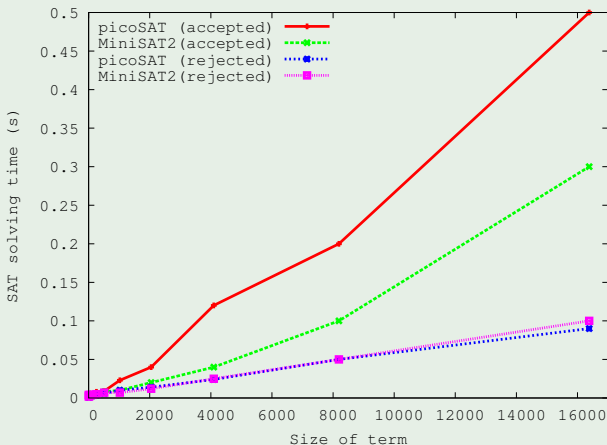
CNF solving time, Laboratory example



The prototype

Results: $\{f(t, t) \mid f \in \Sigma, t \in \mathcal{T}(\Sigma)\}$ example

CNF solving time, $\{f(t, t) \mid f \in \Sigma, t \in \mathcal{T}(\Sigma)\}$



Conclusion, step by step

XML \Rightarrow TAGED $\Rightarrow \Delta_{\mathcal{A}}(t) \Rightarrow$ CNF \Rightarrow SAT solver \Rightarrow *answer*

Conclusion, step by step

XML \Rightarrow TAGED $\Rightarrow \Delta_{\mathcal{A}}(t) \Rightarrow$ CNF \Rightarrow SAT solver \Rightarrow *answer*

Theoretical contribution

SAT encoding for the TAGED uniform membership problem. Natural optimisations.

Complexity

The formula is quadratic in the size of the input, and so is the generation time (worst case).

Implementation

Unoptimised OCaml prototype generating optimised $\Delta_{\mathcal{A}}(t)$ from input TAGED \mathcal{A} and term t .

Conclusion, step by step

XML \Rightarrow TAGED $\Rightarrow \Delta_{\mathcal{A}}(t) \Rightarrow$ CNF \Rightarrow SAT solver \Rightarrow *answer*

Conversion to CNF

We use an existing tool, BAT [Manolios and Vroon, 2009], to convert our formula to DIMACS CNF.

Caveats

For now, CNF conversion is the bottleneck of our experiments:

- BAT is 4.5 times slower than formula generation on big formulæ.
- BAT crashed (stack overflow) on big formulæ

So we could not produce tests big enough to really push the SAT solvers to their limits.

Conclusion, step by step

XML \Rightarrow TAGED $\Rightarrow \Delta_{\mathcal{A}}(t) \Rightarrow$ CNF \Rightarrow SAT solver \Rightarrow *answer*

Tested SAT solvers

Solvers picoSAT and MiniSAT2 both display good performances, MiniSAT2 seeming faster in general.

Results

Efficient (sub-second) SAT solving even in largest tested cases^a

^aOrder of magnitude: Term of 20 000 nodes, formula of 70 000 variables, 120 000 clauses and 250 000 literals (in CNF)

Conclusion, step by step



XML \Rightarrow TAGED $\Rightarrow \Delta_{\mathcal{A}}(t)$ \Rightarrow CNF \Rightarrow SAT solver \Rightarrow *answer*

Conclusion

Overall, the current experimental limitations lie in the computationally easier part of the problem, while its inherent difficulty (*NP*-completeness) seems well overcome by the heuristics of tested SAT solvers.

Some references

[Comon et al., 2007, Filiot et al., 2008, Manolios and Vroon, 2009]

-  Comon, H., Dauchet, M., Gilleron, R., Löding, C., Jacquemard, F., Lugiez, D., Tison, S., and Tommasi, M. (2007).
Tree automata techniques and applications.
release October, 12th 2007.
-  Filiot, E., Talbot, J.-M., and Tison, S. (2008).
Tree Automata with Global Constraints.
In 12th International Conference on Developments in Language Theory (DLT), pages 314–326, Kyoto Japan.
-  Manolios, B. and Vroon, D. (2009).
Faster SAT Solving with Better CNF Generation.
Design, Automation and Test in Europe, DATE.