

A Tree Logic...

... and an Application for the Analysis of Cascading Style Sheets

Pierre Genevès

CNRS – Tyrex team
pierre.geneves@inria.fr

Toccatà seminar, LRI – Feb. 22nd, 2013

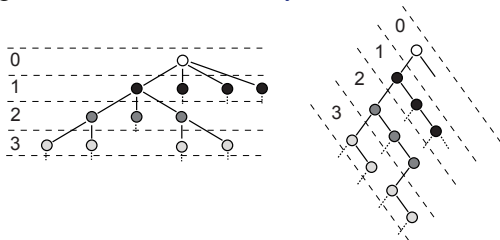
Outline

- 1 Insights on the \mathcal{L}_μ Tree Logic
- 2 Overview of Perspectives and Applications
- 3 Zoom on the Analysis of CSS

Data Model for the Logic

Trees: the logic was originally designed for XML trees

- Specifically: finite binary labeled trees
- They model finite ordered unranked labeled trees wlog
- Bijective encoding of unranked trees as **binary** trees:



Formulas of the \mathcal{L}_μ Logic

- Programs $\alpha \in \{1, 2, \bar{1}, \bar{2}\}$ for navigating binary trees ($\overline{\bar{\alpha}} = \alpha$)



$\mathcal{L}_\mu \ni \varphi, \psi ::=$

	\top
	$p \mid \neg p$
	$n \mid \neg n$
	$\varphi \vee \psi \mid \varphi \wedge \psi$
	$\langle \alpha \rangle \varphi \mid \neg \langle \alpha \rangle \top$
	$\mu X. \varphi$
	$\mu \overline{X}_i. \varphi_i \text{ in } \psi$

formula

true
atomic prop (negated)
nominal (negated)
disjunction (conjunction)
existential (negated)
unary fixpoint (finite recursion)
n -ary fixpoint

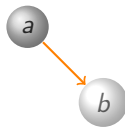
Sample Formula and Satisfying Tree

a



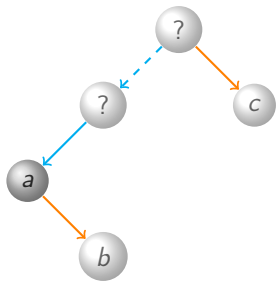
Sample Formula and Satisfying Tree

$a \wedge \langle 2 \rangle b$

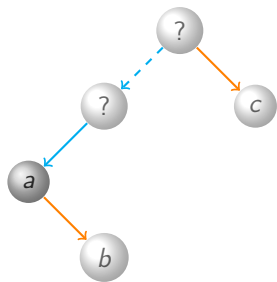


Sample Formula and Satisfying Tree

$a \wedge \langle 2 \rangle b \wedge \mu X. \langle 2 \rangle c \vee \langle \bar{1} \rangle X$

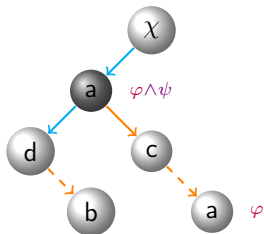


Sample Formula and Satisfying Tree

$$a \wedge \langle 2 \rangle b \wedge \mu X. \langle 2 \rangle c \vee \langle \bar{1} \rangle X$$


- Semantics: models of φ are finite trees for which φ holds at some node
- ✓ Interesting balance between succinctness and expressive power: XPath, CSS selectors, and XML types can be translated into the logic, [linearly](#)

Example: Translation of an XPath Expression into \mathcal{L}_μ



- Formula holds at **selected** nodes
- $\mu Z.\varphi$: finite recursion
- Converse programs are crucial
- More generally, we have a compiler:
 - $t_{\text{xpath}}(e, \chi) : \mathcal{L}_{\text{XPath}} \times \mathcal{L}_\mu \rightarrow \mathcal{L}_\mu$
 - χ is the latest navigation step
 - initially, $\chi = \neg \langle \bar{1} \rangle \top \wedge \neg \langle \bar{2} \rangle \top$ for absolute expressions

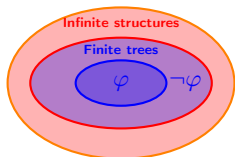
Translated query: $\text{child}::a$ $[\text{child}::b]$

$$\underbrace{a \wedge (\mu Z. \langle \bar{1} \rangle \chi \vee \langle \bar{2} \rangle Z)}_{\varphi} \quad \wedge \quad \underbrace{\langle 1 \rangle \mu Y. b \vee \langle 2 \rangle Y}_{\psi}$$

\mathcal{L}_μ Closure under Negation

Cycle-freeness: A key property

- If both a program and its converse occur between a $\mu X.$ binder and X , formula has a cycle, e.g.: $\mu X. \langle \alpha \rangle X \vee \langle \bar{\alpha} \rangle X$
- Otherwise the formula is **cycle-free**
- in practice, most (all?) formulas are **cycle-free** (e.g. XPath translations are always cycle-free)



- **Cycle-freeness** of \mathcal{L}_μ implies **closure under negation**
 - The negation of finite recursion is finite recursion (see paper)
 - $\neg\varphi$ is easily (linearly) expressible in \mathcal{L}_μ for all $\varphi \in \mathcal{L}_\mu$
- Crucial for BC: implication (subtyping, containment tests...)
- Crucial for implementation

Deciding \mathcal{L}_μ Satisfiability

Is a formula $\psi \in \mathcal{L}_\mu$ satisfiable?

- Given ψ , determine whether there exists a finite tree that satisfies ψ
- Validity: test $\neg\psi$

Principles: Automatic Theorem Proving

- Search for a proof tree
- Build the proof bottom up:
“if ψ holds then it is necessarily somewhere up”

Search Space Optimization

Idea: Truth Status is Inductive

- The truth status of ψ can be expressed as a function of its subformulas
- For boolean connectives, it can be deduced (truth tables)
- Only base subformulas really matter: $\text{Lean}(\psi)$



A Tree Node: Truth Assignment of $\text{Lean}(\psi)$ Formulas

- With some additional constraints, e.g. $\neg \langle \bar{1} \rangle \top \vee \neg \langle \bar{2} \rangle \top$

Satisfiability-Testing Algorithm: Principles

Bottom-up construction of proof tree

- A set of nodes is repeatedly updated (fixpoint computation)

Satisfiability-Testing Algorithm: Principles



Bottom-up construction of proof tree

- Step 1: all possible leaves are added

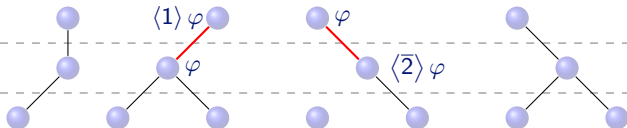
Satisfiability-Testing Algorithm: Principles



Bottom-up construction of proof tree

- Step $i > 1$: all possible parents of previous nodes are added

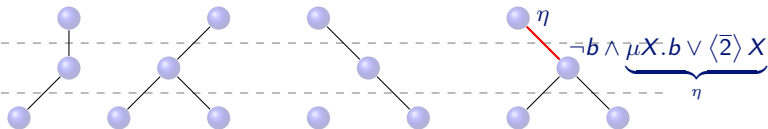
Satisfiability-Testing Algorithm: Principles



Compatibility relation between nodes

- Nodes from previous step are proof support:
 $\langle \alpha \rangle \varphi$ is added if φ holds in some node added at previous step

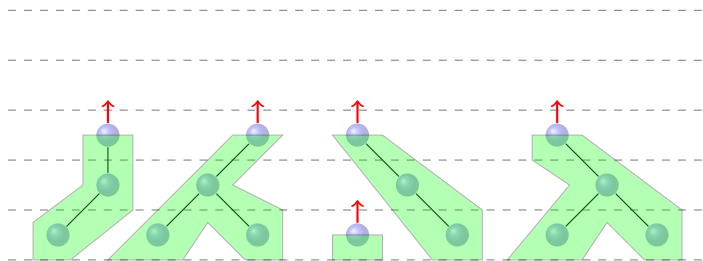
Satisfiability-Testing Algorithm: Principles



Compatibility relation between nodes

- Nodes from previous step are proof support:
 $\langle \alpha \rangle \varphi$ is added if φ holds in some node added at previous step

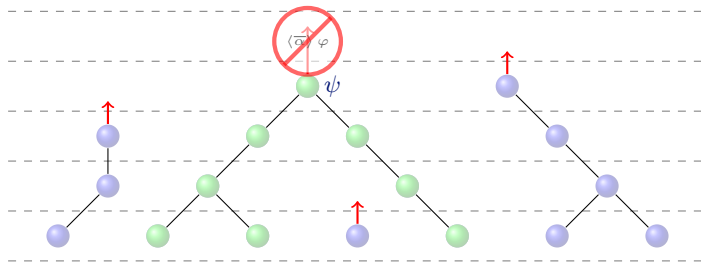
Satisfiability-Testing Algorithm: Principles



Progressive bottom-up reasoning (partial satisfiability)

- $\langle \bar{\alpha} \rangle \varphi$ are left unproved until a parent is connected

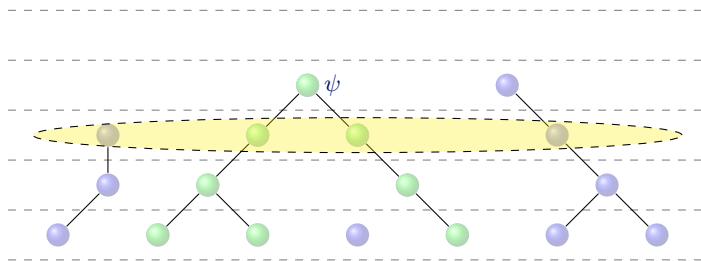
Satisfiability-Testing Algorithm: Principles



Termination

- If ψ is present in some **root** node, then ψ is satisfiable
- Otherwise, the algorithm terminates when no more nodes can be added

Satisfiability-Testing Algorithm: Principles



Implementation techniques

- Crucial optimization: symbolic representation

Correctness & Complexity

Theorem

The satisfiability problem for a formula $\psi \in \mathcal{L}_\mu$ is decidable in time $2^{O(n)}$ where $n = |\text{Lean}(\psi)|$.

System fully implemented

- decision procedure
- compilers (XPath, DTD, XML Schema, CSS selectors, ...)

Overview of Some Experiments

DTD	Symbols	Binary type variables
SMIL 1.0	19	11
XHTML 1.0 Strict	77	325

Table: Types used in experiments.

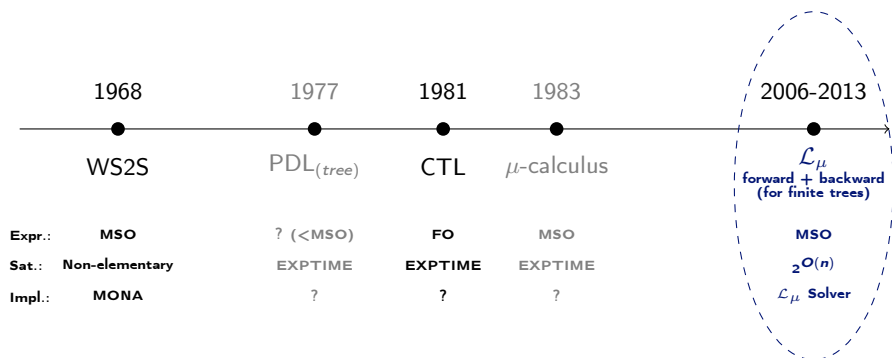
XPath decision problem	XML type	Time (ms)
$e_1 \subseteq e_2$ and $e_2 \not\subseteq e_1$	none	353
$e_4 \subseteq e_3$ and $e_4 \subseteq e_3$	none	45
$e_6 \subseteq e_5$ and $e_5 \not\subseteq e_6$	none	41
e_7 is satisfiable	SMIL 1.0	157
e_8 is satisfiable	XHTML 1.0	2630
$e_9 \subseteq (e_{10} \cup e_{11} \cup e_{12})$	XHTML 1.0	2872

Table: Some decision problems and corresponding results.

For the last test, size of the Lean is 550. The search space is $2^{550} \approx 10^{165}$... more than the square number of atoms in the universe 10^{80}

Tree Logics: an Overview

- On the theoretical side: \mathcal{L}_μ offers an interesting expressivity, succinctness, optimal complexity bound



On the practical side:

- except (hyperexponential) MONA, this is the only one implementation of a satisfiability solver for such an expressive logic
- It can be useful for graphs too: the sublogic without backward modalities enjoys the finite tree model property

Going Further: Challenges

Several directions

- Growing logical **expressive power**? (currently MSO)
- Decreasing combined **complexity**? (impossible without dropping features: containment for regular tree grammars is hard for EXPTIME)
- Augmenting **succinctness** of the logic → good potential

Succinctness is crucial

- A blow-up in the logical translations affects the combined complexity
- Augmenting succinctness is a way to address more problems in EXPTIME

Further Perspectives: $\text{card}(\text{phi})=n$

$\text{card}(\text{phi})=n$

- Even if this remains **regular**, this is not a priori **succinct**
- For instance, L_{2a2b} : set of strings over $\Sigma = \{a, b, c\}$ containing at least 2 occurrences of a and at least two occurrences of b

Further Perspectives: $\text{card}(\text{phi})=n$

$\text{card}(\text{phi})=n$

- Even if this remains **regular**, this is not a priori **succinct**
- For instance, L_{2a2b} : set of strings over $\Sigma = \{a, b, c\}$ containing at least 2 occurrences of a and at least two occurrences of b

$$\begin{aligned} & (a|b|c)^* a(a|b|c)^* a(a|b|c)^* b(a|b|c)^* b(a|b|c)^* | \\ & (a|b|c)^* a(a|b|c)^* b(a|b|c)^* a(a|b|c)^* b(a|b|c)^* | \\ & (a|b|c)^* a(a|b|c)^* b(a|b|c)^* b(a|b|c)^* a(a|b|c)^* | \\ & (a|b|c)^* b(a|b|c)^* b(a|b|c)^* a(a|b|c)^* a(a|b|c)^* | \\ & (a|b|c)^* b(a|b|c)^* a(a|b|c)^* b(a|b|c)^* a(a|b|c)^* | \\ & (a|b|c)^* b(a|b|c)^* a(a|b|c)^* a(a|b|c)^* b(a|b|c)^* \end{aligned}$$

Further Perspectives: $\text{card}(\phi)=n$

- If we add \cap to the regular expression operators:

$$((a|b|c)^* a(a|b|c)^* a(a|b|c)^*) \cap ((a|b|c)^* b(a|b|c)^* b(a|b|c)^*)$$

- In logical terms, conjunction offers a dramatic reduction in expression size
- If we now consider the ability to describe numerical constraints on the frequency of occurrences, we get another exponential reduction in size:

$$((a|b|c)^* a(a|b|c)^*)^2 \cap ((a|b|c)^* b(a|b|c)^*)^2$$

- Crucial when the complexity of the decision procedure depends on the formula size

Further Perspectives: $\text{card}(\phi)=n$

Querying all the articles with 4 or more authors

- Navigational XPath expression:

```
article[author/following-sibling::author/following-sibling::author/following-sibling::author]
```

or, using the counting operator in XPath:

```
article[count(author)>=4]
```

- The counting operator is **exponentially** more succinct
- Again, we would like efficient static analyzers that directly operate on the succinct form! (i.e. not pay the price of the blow-up)

Facts

Nominals + Backward modalities + $\text{card}(\phi)=n$

- undecidable over graphs [Bonatti-AI'04]
- decidable over finite trees

Ongoing research...

- What is the precise complexity for $\text{card}(\phi)=n$ for finite trees?
 - ... or more generally of rich logical **combinators** that may duplicate formulas of arbitrary length (but in a particular manner)?
- Hint: look at the factorization power of the Lean

Further Perspectives: Follow the Arrows

- So far: logical description of **structural constraints** stemming from queries and schemas
- Can we also logically capture a notion of **computation** performed by programs (i.e. functions)?
- For example, can the logic capture the type algebra on which CDuce sits?

τ	::=	b	basic type
		$\tau \times \tau$	product type
		$\tau \rightarrow \tau$	function type
		$\tau \vee \tau$	union type
		$\neg \tau$	complement type
		0	empty type
		v	recursion variable
		$\mu v. \tau$	recursive type

- Yes. We interpret the type algebra in a purely logical manner...

Further Perspectives: Follow the Arrows

Representing functions

$f = \{(d_1, d'_1), (d_2, d'_2), \dots\}$ modelizes a function such that:

- $f d_i$ may evaluate (nondeterministically) to d'_i
- $f x$ where $x \notin \{d_i\}$ never terminates (and is well-typed)
- if $d'_i = ERR$ then $f d_i$ is a type error

Lemma (Frisch et al.): considering only finite such sets of pairs is sufficient for defining semantic subtyping.

Further Perspectives: Follow the Arrows

Types as Logical Formulas (detailed encoding in [ICFP'11])

- Interpretation of $\tau_1 \rightarrow \tau_2$: all finite f s such that $f : \tau_1 \rightarrow \tau_2$
- $\text{form}(\tau_1 \rightarrow \tau_2) = (\rightarrow) \wedge [1] \mu X. ([2] X \wedge \langle 1 \rangle (\neg \text{form}(\tau_1) \vee \langle 2 \rangle \text{form}(\tau_2)))$
- with the shorthand $[\alpha] \varphi = \neg \langle \alpha \rangle \top \vee \langle \alpha \rangle \varphi$
- Intuitively: “a (\rightarrow) node whose first child, if it exists, satisfies X ”
- where X = “a node whose next sibling, if it exists, satisfies X , and which has a first child which either does not satisfy $\text{form}(\tau_1)$ or has a next sibling which satisfies $\text{form}(\tau_2)$.”

Further Perspectives: Parametric Polymorphism

- We can go even further and support parametric polymorphism
 - We add type variables α to the type algebra
 - Intuition of subtyping in the presence of type variables:
 $\tau_1(\bar{\alpha}) \leq \tau_2(\bar{\alpha})$ whenever, independently of the variables $\bar{\alpha}$, any value of type τ_1 has type τ_2 as well.
- Neat formal definition of subtyping by Castagna and Xu (ICFP'11)
- Complete logical encoding in [ICFP'11] (Gesbert, Genevès and Layaïda)
- We can solve subtyping with the satisfiability solver

Interesting facts

- **The complexity bound is not affected:** $2^{\mathcal{O}(|\tau_1|+|\tau_2|)}$ for checking $\tau_1 \leq \tau_2$
- The \mathcal{L}_μ logic is expressive and robust by (intricate) extension

Further Perspectives: Type Synthesis

- **Objective:** static type checking for programming languages that do not require type annotations
- **Method:** (i) type inference, (ii) containment check (unsatisfiability check)
- If the containment check fails between the inferred type and e.g. the expected output type, an error is reported
- **Novelty:** Take advantage of the logic succinctness to represent inferred type portions (ongoing research...)
- A possible application: enhancing static type checking for XQuery
- Current XQuery standardized type system is unsound so far
 - if a program involves an upward navigation such as `parent::*`, the type Any (true in logic) is inferred
 - false negatives may be reported

Some Already Investigated Applications

- Containment for XML queries [PLDI'07, ICDE'10]
 - equivalence test for monadic queries: $\forall t, \forall n \in t, q_1(t, n) \stackrel{?}{=} q_2(t, n)$
- Modeling interleaving and counting [IJCAI'11]
- Dead code analysis for XQuery [ICSE'10, ICSE'11]
- Impact of schema evolution [ICFP'09, TOIT'11]
 - Schema S evolves into S' : impact on a query written against S ?
- Deciding subtyping for rich type algebras [ICFP'11]
 - Intersection, negation, **function**, and **polymorphic** types
- Containment for SPARQL queries (polyadic, graphs) under constraints [AAAI'12, IJCAR'12]
- CSS Analysis [WWW'12]

Try it online*: <http://wam.inrialpes.fr/websolver>



XML Reasoning Solver Project

Home **Demo** Documentation Publications Team

Enter your formula below:

```
bool() = {true|false};
list() = let $l = (_a * $1) | {nil} in $l;
odd() = let $o = (_a * _a * $o) | (_a * {nil}) in $o;
even() = let $e = (_a * _a * $e) | {nil} in $e;
nsubtype ( odd() -> {true} & (even() -> {false}), list() -> bool() )
```

See [user manual](#) or pick an example

- [XPath Satisfiability #1](#)
- [XPath Satisfiability #2](#)
- [XPath Containment](#)
- [XPath Equivalence](#)
- [Mu-formula with values](#)
- [Mu-formula with recursion](#)
- [XHTML Type Evolution](#)
- [MathML Query Evolution](#)
- [Polymorphism with arrow types #1](#)
- [Polymorphism with arrow types #2](#)
- [Regular expression intersection](#)
- [Regular expression equivalence](#)

► Advanced Options

This online demo is a 100% Java implementation of the solver that runs inside a Tomcat servlet. It is based on a thread-safe re-implementation of a BDD package (JavaBDD). However, the performance of this package is very slow compared to what can be achieved with an off-line solver implementation with native BDDs. Ask us if you are interested in the high-speed off-line version of the solver.

* or offline if performance is critical: the offline version is much faster (native BDD library, further optimizations like compression of symbols)

Long-Term Goal

Long-term view

Heterogeneity is here to stay: JSON (JS serialization) + XML + RDF (knowledge)

A unified verification toolbox

- for type-checking web programs: XQuery, XPath, “X...”, Jaql etc.
- for reasoning at the layout level: CSS
- for supporting heterogenous and rich data values: XML, RDF, JSON ...
- possibly constrained by some schema languages (XML Schema, RDFS, Schematron, etc.)

Long-Term Goal

Long-term view

Heterogeneity is here to stay: JSON (JS serialization) + XML + RDF (knowledge)

A unified verification toolbox

- for type-checking web programs: XQuery, XPath, “X...”, Jaql etc.
- for reasoning at the layout level: CSS
- for supporting heterogenous and rich data values: XML, RDF, JSON ...
- possibly constrained by some schema languages (XML Schema, RDFS, Schematron, etc.)