

On XPath with Transitive Axes and Data Tests^{*}

Diego Figueira
University of Edinburgh
Edinburgh, UK

ABSTRACT

We study the satisfiability problem for XPath with data equality tests. XPath is a node selecting language for XML documents whose satisfiability problem is known to be undecidable, even for very simple fragments. However, we show that the satisfiability for XPath with the rightward, leftward and downward reflexive-transitive axes (namely *following-sibling-or-self*, *preceding-sibling-or-self*, *descendant-or-self*) is decidable. Our algorithm yields a complexity of 3EXPSPACE , and we also identify an expressive-equivalent normal form for the logic for which the satisfiability problem is in 2EXPSPACE . These results are in contrast with the undecidability of the satisfiability problem as soon as we replace the reflexive-transitive axes with just transitive (non-reflexive) ones.

Categories and Subject Descriptors

I.7.2 [Document Preparation]: Markup Languages; H.2.3 [Database Management]: Languages; H.2.3 [Languages]: Query Languages

General Terms

Algorithms, Languages

Keywords

XML, XPath, unranked unordered tree, reflexive transitive axes, data-tree, infinite alphabet, data values

^{*}We acknowledge the financial support of the Future and Emerging Technologies (FET) programme within the Seventh Framework Programme for Research of the European Commission, under the FET-Open grant agreement FOX, number FP7-ICT-233599.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

PODS'13, June 22–27, 2013, New York, New York, USA.
Copyright 2013 ACM 978-1-4503-2066-5/13/06 ...\$15.00.

1. INTRODUCTION

The simplest way of abstracting an XML document is by seeing it as a tree over a finite alphabet of *tags* or *labels*. However, this abstraction ignores all actual data stored in the document attributes. This is why there has been an increasing interest in data trees: trees that also carry data from an infinite domain. Here, we consider an XML modeled as an unranked ordered finite tree whose every node contains a label, and a vector of data values, one for each attribute. Labels belong to some finite alphabet, and data values to some infinite domain. We call these models *multi-attribute data trees* (see Figure 1). We study logics on these models, that can express data properties, namely equality of attributes' data values.

Here, we show decidability of the satisfiability problem for XPath where navigation can be done going downwards, rightwards or leftwards in the XML document, that is, where navigation is done using the reflexive-transitive XPath axes *descendant-or-self*, *following-sibling-or-self*, and *preceding-sibling-or-self*.

Formalisms for trees with data values

Several formalisms have been studied lately in relation to static analysis on trees with data values.

First-order logic.

One such formalisms is $\text{FO}^2(<_h, \text{succ}_h, <_v, \text{succ}_v, \sim)$, first order logic with two variables, and binary relations to navigate the tree: the descendant $<_v$, child succ_v , next sibling succ_h and following sibling $<_h$ (*i.e.*, the transitive closure of succ_h); and an equivalence relation \sim to express that two nodes of the trees have the same data value. Although the decidability status for the satisfiability problem of this logic is unknown, it is known to be as hard as the reachability problem for BVASS (Branching Vector Addition System with States) [4]. If the signature has only the child and next sibling relation— $\text{FO}^2(\text{succ}_h, \text{succ}_v, \sim)$ —the logic is decidable in 3NEXPTIME as shown in [4].

Automata.

There have also been works on automata models for trees with data. Tree automata with registers to store and compare data values were studied in [20] as an extension to a similar model on words [19, 22]. A decidable alternating version of these automata called ATRA was studied in [18], and it was extended in [9, 12] to show decidability of the satisfiability problem for forward-XPath. The work [3] introduces a simple yet powerful automata model called *Class Automata*

on data trees that can capture $\text{FO}^2(\langle_h, \text{succ}_h, \langle_v, \text{succ}_v, \sim\rangle)$, XPath, ATRA, and other models. Although its emptiness problem is undecidable, classes of data trees for which it is decidable are studied in [1]. Other formalisms include tree automata with set and linear constraints on cardinalities of sets of data values [6, 23].

XPath.

Here we concentrate on XPath, which is incomparable in terms of expressiveness with all the previously mentioned formalisms (except for Class Automata).

XPath is arguably the most widely used XML query language. It is implemented in XSLT and XQuery and it is used as a constituent part of several specification and update languages. XPath is fundamentally a general purpose language for addressing, searching, and matching pieces of an XML document. It is an open standard and constitutes a World Wide Web Consortium (W3C) Recommendation [5].

Query containment and query equivalence are important static analysis problems, which are useful to, for example, query optimization tasks. In logics closed under boolean operators—as the one treated here—, these problems reduce to checking for *satisfiability*: Is there a document on which a given query has a non-empty result? By answering this question we can decide at compile time whether the query contains a contradiction and thus the computation of the query (or subquery) on the document can be avoided. Or, by answering the query equivalence problem, one can test if a query can be safely replaced by another one which is more optimized in some sense (*e.g.*, in the use of some resource). Moreover, the satisfiability problem is crucial for applications on security [7], type checking transformations [21], and consistency of XML specifications.

Core-XPath (term coined in [17]) is the fragment of XPath 1.0 that captures all the navigational behavior of XPath. It has been well studied and its satisfiability problem is known to be decidable even in the presence of DTDs. The extension of this language with the possibility to make equality and inequality tests between attributes of elements in the XML document is named **Core-Data-XPath** in [4].

In an nutshell, the important formulas of **Core-Data-XPath** (henceforth *XPath*) are of the form

$$\langle \alpha @_i = \beta @_j \rangle,$$

where α, β are *path expressions*, that navigate the tree using *axes*: descendant, child, ancestor, next-sibling, etc. and can make tests in intermediary nodes. Such a formula is true at a node x of a multi-attribute data tree if there are two nodes y, z in the tree that can be reached with the relations denoted by α, β respectively, so that the i th attribute of y carries the *same* datum as the j th attribute of z .

Unfortunately, the satisfiability problem for XPath is undecidable [16]. How can we regain decidability for satisfiability of XPath then? We can restrict the models, or restrict the logic. The first possibility is to restrict the classes of documents in which one is interested, which is the approach taken in [1]. Another, more studied, approach is to restrict the syntax, which is the one taken here. One way to regain decidability is to syntactically restrict the amount of nodes that the XPath properties can talk about. In this vein, there have been studies on fragments without negation or without transitive axes [2, 16]. These fragments enjoy a small model property and are decidable. However, they

cannot state global properties, involving all the nodes in an XML document. Ideally, we seek fragments with the following desirable features

- closed under boolean operators,
- having as much freedom as possible to navigate the tree in many directions: up, down, left, right,
- having the possibility to reach any node of the tree, with transitive axes, like descendant, following sibling (the transitive closure of the next sibling axis), etc.

However, decidability results are scarce, and most fragments with the characteristics just described are undecidable. There are, however, some exceptions ([10]).

- The *downward* fragment of XPath, containing the child and descendant axes, is decidable, EXPTIME-complete [8, 13].
- The *forward* fragment of XPath, extending the downward fragment with the next sibling and the following sibling axes, is decidable with non-primitive recursive complexity [9, 12].
- The *vertical* fragment of XPath, extending the downward fragment with the parent and ancestor axes, is decidable with non-primitive recursive complexity [15].
- A last example is the present work: XPath with the reflexive transitive closure of the child, next-sibling and previous-sibling relations is decidable.

All the non-primitive recursive (NPR) upper bounds of the forward and vertical fragments are also matched with NPR lower bounds. That is, there is no primitive recursive function that bounds the time or space needed by any algorithm that computes the satisfiability for any of these two logics. Moreover, it is known that any fragment of XPath containing a transitive rightward, leftward or upward axis has a satisfiability problem which is either undecidable or decidable with a NPR lower bound [14].¹ Further, as soon as we have both the rightward and leftward transitive axes, the satisfiability becomes undecidable [14]. (Indeed, the downward fragment of XPath seemed to be the only one with elementary complexity up to now.)

The aforementioned hardness results make use of non-reflexive transitive relations. Surprisingly, the reductions do not seem to work when the relations are also reflexive. What is then the decidability status of the fragments of XPath with reflexive-transitive relations? This was a question raised in [14].

A partial answer to this question was given in [11]. There, it was shown that XPath restricted to data words is decidable even when we have both a reflexive-transitive future and past relations. (One can think of data words as XML documents of height 1, with only one attribute per node.) This result may seem surprising taking into account that if one of these relations is non-reflexive it is no longer decidable; and if we have only one non-reflexive transitive relation it is decidable with non-primitive recursive complexity. In [11] it was shown that the satisfiability problem is in 2EXPSPACE (or EXPSPACE if we adopt a certain normal form of

¹These are the axes that are called preceding-sibling, following-sibling and ancestor in the XPath jargon.

the formulas). This was a first step in our study of the computational behavior of XPath with reflexive-transitive axes. The present work corresponds to the second part, in which we study XPath on XML documents (*i.e.*, trees) instead of words.

Contribution

We show decidability of the satisfiability for XPath with data equality tests between attributes, where navigation can be done going downwards, rightwards or leftwards in the XML document. The navigation can only be done by reflexive-transitive relations. These correspond to the XPath axes: preceding-sibling-or-self, following-sibling-or-self, and descendant-or-self axes.² Here we denote these axes with $\ast\leftarrow, \rightarrow\ast$ and $\downarrow\ast$ respectively. As already mentioned, the fact that the relations are reflexive-transitive (as opposed to just transitive) is an essential feature to achieve decidability. Given the known complexity results on XPath, this fragment seems to be in balance between navigation and complexity. This work then argues in favor of studying XPath-like logics for trees with data with reflexive-transitive relations, since they may behave computationally much better than the non-reflexive counterpart, as evidenced here.

The extension of the prior work [11] on data words to work with trees with a descendant axis is highly non-trivial, requiring an altogether different formalism and algorithm strategy. Whereas in [11] the main object of study is a transition system—which comes naturally when working with words—this does not adapt well to working with trees. Instead, here we work with an algebra operating on abstractions of *forests* of multi-attribute data trees. Over this algebra, we prove some monotonicity properties, which are necessarily more involved than those used in [11] to account for the interplay between horizontal and vertical navigation of the logic.

Our algorithm yields a 3EXPSpace upper bound for the satisfiability problem of this XPath fragment. We also show that this can be lowered to 2EXPSpace if we work with an expressive-equivalent normal form, called *direct normal form*. Since XPath with just one reflexive-transitive relation is already EXPSpace-hard (even when the formula is in direct normal form) by [11], we cannot aim for much better complexities.

2. PRELIMINARIES

Basic notation.

Let $\mathbb{N}_0 \stackrel{\text{def}}{=} \{0, 1, 2, \dots\}$, $\mathbb{N} \stackrel{\text{def}}{=} \{1, 2, 3, \dots\}$, and let $[n] \stackrel{\text{def}}{=} \{1, \dots, n\}$ for any $n \in \mathbb{N}$. We fix once and for all \mathbb{D} to be any infinite domain of data values; for simplicity in our examples we will consider $\mathbb{D} = \mathbb{N}_0$. In general we use the symbols \mathbb{A} , \mathbb{B} for finite alphabets, and the symbols \mathbb{E} and \mathbb{F} for any kind of alphabet. By \mathbb{E}^* we denote the set of finite sequences over \mathbb{E} , by \mathbb{E}^+ the set of finite sequences with at least one element over \mathbb{E} . We write ‘ ϵ ’ for the empty sequence and ‘ \cdot ’ as the concatenation operator between sequences. By $|S|$ we denote the length of S (if S is a finite sequence), or its cardinality (if S is a set). We use $(a_i)_{i \in \{j, \dots, j+n\}}$ as short for $a_j a_{j+1} \cdots a_{j+n}$.

²Strictly speaking, these axes do not exist in XPath [5]. They must be interpreted as the reflexive version of the preceding-sibling, following-sibling and descendant axes respectively.

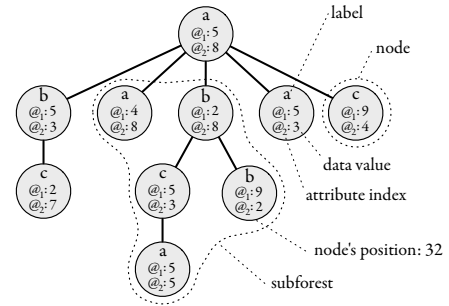


Figure 1: A multi-attribute data tree, where $\mathbb{A} = \{a, b, c\}$ and $k = 2$.

Unranked finite trees with data.

By $Trees(\mathbb{E})$ we denote the set of finite ordered and unranked trees over an alphabet \mathbb{E} . We view each **position** in a tree as an element of \mathbb{N}^* . Formally, we define $POS \subseteq 2^{\mathbb{N}^*}$ as the set of sets of finite tree positions, such that: $X \in POS$ iff (a) $X \subseteq \mathbb{N}^*$, $|X| < \infty$; (b) X is prefix-closed; and (c) if $n \cdot (i+1) \in X$ for $i \in \mathbb{N}$, then $n \cdot i \in X$. A tree is then a mapping from a set of positions to labels of the alphabet $Trees(\mathbb{E}) \stackrel{\text{def}}{=} \{\mathbf{t} : P \rightarrow \mathbb{E} \mid P \in POS\}$. The root’s position is the empty string ϵ . The position of any other node in the tree is the concatenation of the position of its parent and the node’s index in the ordered list of siblings.

Given a tree $\mathbf{t} \in Trees(\mathbb{E})$, $pos(\mathbf{t})$ denotes the domain of \mathbf{t} , which consists of the set of positions of the tree, and $alph(\mathbf{t}) = \mathbb{E}$ denotes the alphabet of the tree. From now on, we informally refer by ‘node’ to a position x together with the value $\mathbf{t}(x)$.

Given two trees $\mathbf{t}_1 \in Trees(\mathbb{E})$, $\mathbf{t}_2 \in Trees(\mathbb{F})$ such that $pos(\mathbf{t}_1) = pos(\mathbf{t}_2) = P$, we define $\mathbf{t}_1 \otimes \mathbf{t}_2 : P \rightarrow (\mathbb{E} \times \mathbb{F})$ as $(\mathbf{t}_1 \otimes \mathbf{t}_2)(x) \stackrel{\text{def}}{=} (\mathbf{t}_1(x), \mathbf{t}_2(x))$.

The set of **multi-attribute data trees** over a finite alphabet \mathbb{A} of labels, k different attributes, and an infinite domain \mathbb{D} is defined as $Trees(\mathbb{A} \times \mathbb{D}^k)$. Note that every tree $\mathbf{t} \in Trees(\mathbb{A} \times \mathbb{D}^k)$ can be decomposed into two trees $\mathbf{a} \in Trees(\mathbb{A})$ and $\mathbf{d} \in Trees(\mathbb{D}^k)$ such that $\mathbf{t} = \mathbf{a} \otimes \mathbf{d}$. Figure 1 shows an example of a multi-attribute data tree. The notation for the set of data values used in a data tree is $data(\mathbf{a} \otimes \mathbf{d}) \stackrel{\text{def}}{=} \{\mathbf{d}(x)(i) \mid x \in pos(\mathbf{d}), i \in [k]\}$. With an abuse of notation we write $data(X)$ to denote all the elements of \mathbb{D} contained in X , for whatever object X may be.

A **forest** is a sequence of trees, and the set of **multi-attribute data forests** over \mathbb{A} and k is $(Trees(\mathbb{A} \times \mathbb{D}^k))^*$. We will normally use the symbol $\bar{\mathbf{t}}$ for a forest of multi-attribute data trees. That is, $\bar{\mathbf{t}} \in (Trees(\mathbb{A} \times \mathbb{D}^k))^*$. (Note that in particular $\bar{\mathbf{t}}$ can be an *empty forest*.) For any $(a, \bar{d}) \in \mathbb{A} \times \mathbb{D}^k$, let us write $(a, \bar{d})\bar{\mathbf{t}}$ for the multi-attribute data tree that results from adding (a, \bar{d}) as a root of $\bar{\mathbf{t}}$. We call this operation **rooting**. We will usually write \mathbf{t} (resp. $\bar{\mathbf{t}}$) to denote multi-attribute data trees (resp. forests) and t (resp. \bar{t}) to denote trees (resp. forests) over a finite alphabet.

XPath.

Next we define transitive XPath, the fragment of XPath where all axes are reflexive and transitive.

Transitive XPath is a two-sorted language, with *path* ex-

$$\begin{aligned} \alpha, \beta &::= o \mid \alpha[\varphi] \mid [\varphi]\alpha \mid \alpha\beta & o \in \{\varepsilon, \downarrow_*, \uparrow^*, \rightarrow^*, * \leftarrow\}, \\ \varphi, \psi &::= a \mid \neg\varphi \mid \varphi \wedge \psi \mid \varphi \vee \psi \mid \langle \alpha \rangle \mid \langle \alpha @_i = \beta @_j \rangle \mid \langle \alpha @_i \neq \beta @_j \rangle & a \in \mathbb{A}, i, j \in [\mathbf{k}]. \end{aligned}$$

$$\begin{aligned} \llbracket \downarrow_* \rrbracket^{\mathbf{t}} &= \{(x, x \cdot i) \mid x \cdot i \in \text{pos}(\mathbf{t})\}^* & \llbracket \uparrow^* \rrbracket^{\mathbf{t}} &= \{(x \cdot i, x) \mid x \cdot i \in \text{pos}(\mathbf{t})\}^* \\ \llbracket \rightarrow^* \rrbracket^{\mathbf{t}} &= \{(x \cdot i, x \cdot (i+1)) \mid x \cdot i, x \cdot (i+1) \in \text{pos}(\mathbf{t})\}^* & \llbracket * \leftarrow \rrbracket^{\mathbf{t}} &= \{(x \cdot (i+1), x \cdot i) \mid x \cdot i, x \cdot (i+1) \in \text{pos}(\mathbf{t})\}^* \\ \llbracket \varepsilon \rrbracket^{\mathbf{t}} &= \{(x, x) \mid x \in \text{pos}(\mathbf{t})\} & \llbracket \alpha \beta \rrbracket^{\mathbf{t}} &= \{(x, z) \mid \text{there exists } y \text{ such that} \\ & & & (x, y) \in \llbracket \alpha \rrbracket^{\mathbf{t}}, (y, z) \in \llbracket \beta \rrbracket^{\mathbf{t}}\} \\ \llbracket [\varphi] \rrbracket^{\mathbf{t}} &= \{(x, x) \mid x \in \text{pos}(\mathbf{t}), x \in \llbracket [\varphi] \rrbracket^{\mathbf{t}}\} & \llbracket \langle \alpha \rangle \rrbracket^{\mathbf{t}} &= \{x \in \text{pos}(\mathbf{t}) \mid \exists y. (x, y) \in \llbracket \alpha \rrbracket^{\mathbf{t}}\} \\ \llbracket a \rrbracket^{\mathbf{t}} &= \{x \in \text{pos}(\mathbf{t}) \mid \mathbf{a}(x) = a\} & \llbracket \varphi \wedge \psi \rrbracket^{\mathbf{t}} &= \llbracket \varphi \rrbracket^{\mathbf{t}} \cap \llbracket \psi \rrbracket^{\mathbf{t}} \\ \llbracket \neg\varphi \rrbracket^{\mathbf{t}} &= \text{pos}(\mathbf{t}) \setminus \llbracket \varphi \rrbracket^{\mathbf{t}} & \llbracket \langle \alpha @_i = \beta @_j \rangle \rrbracket^{\mathbf{t}} &= \{x \in \text{pos}(\mathbf{t}) \mid \exists y, z (x, y) \in \llbracket \alpha \rrbracket^{\mathbf{t}}, \\ & & & (x, z) \in \llbracket \beta \rrbracket^{\mathbf{t}}, \mathbf{d}(y)(i) = \mathbf{d}(z)(j)\} \\ \llbracket \langle \alpha @_i = \beta @_j \rangle \rrbracket^{\mathbf{t}} &= \{x \in \text{pos}(\mathbf{t}) \mid \exists y, z (x, y) \in \llbracket \alpha \rrbracket^{\mathbf{t}}, & & \llbracket \langle \alpha @_i \neq \beta @_j \rangle \rrbracket^{\mathbf{t}} &= \{x \in \text{pos}(\mathbf{t}) \mid \exists y, z (x, y) \in \llbracket \alpha \rrbracket^{\mathbf{t}}, \\ & & & & (x, z) \in \llbracket \beta \rrbracket^{\mathbf{t}}, \mathbf{d}(y)(i) \neq \mathbf{d}(z)(j)\} \end{aligned}$$

Figure 2: The syntax of transitive XPath; and its semantics for a multi-attribute data tree $\mathbf{t} = \mathbf{a} \otimes \mathbf{d}$.

pressions (that we write $\alpha, \beta, \gamma, \delta$) and *node* expressions (that we write φ, ψ, η). Path expressions denote binary relations, resulting from composing the descendant, ancestor, preceding sibling and following sibling relations (which are denoted respectively by \downarrow_* , \uparrow^* , $* \leftarrow$, \rightarrow^* respectively), and node expressions. Node expressions are boolean formulas that test a property of a node like, for example, that it has a certain label, or that it has a descendant labeled a with the same data value in attribute i as the attribute j of an ancestor labeled b , which is expressed by $\langle \downarrow_* [a] @_i = \uparrow^* [b] @_j \rangle$. We write $\text{XPath}(\downarrow_*, \uparrow^*, \rightarrow^*, * \leftarrow, =)$ to denote this logic, and we write $\text{XPath}(\mathcal{O}, =)$ for some $\mathcal{O} \subseteq \{\downarrow_*, \uparrow^*, \rightarrow^*, * \leftarrow\}$, to denote the logic containing only the axes in \mathcal{O} . A *formula* of $\text{XPath}(\downarrow_*, \uparrow^*, \rightarrow^*, * \leftarrow, =)$ is either a node expression or a path expression of the logic. Its syntax and semantics are defined in Figure 2. As another example, we can select the nodes that have a sibling labeled a to the left whose first attribute is the same as the second attribute of some descendant of a right sibling by the formula $\varphi = \langle * \leftarrow [a] @_1 = \rightarrow^* \downarrow_* @_2 \rangle$. Given a tree \mathbf{t} as in Figure 1, we have $\llbracket \varphi \rrbracket^{\mathbf{t}} = \{\epsilon, 2, 3, 4, 5, 311\}$.

We write $\mathbf{t}, x \models \varphi$ (resp. $\mathbf{t}, (x, y) \models \alpha$) for $x \in \text{pos}(\mathbf{t})$ (resp. $x, y \in \text{pos}(\mathbf{t})$) as short for $x \in \llbracket \varphi \rrbracket^{\mathbf{t}}$ (resp. $(x, y) \in \llbracket \alpha \rrbracket^{\mathbf{t}}$). We write $\mathbf{t} \models \varphi$ as short for $\epsilon \in \llbracket \varphi \rrbracket^{\mathbf{t}}$.

In the case of $\text{XPath}(* \leftarrow, \downarrow_*, \rightarrow^*, =)$, we also extend the evaluation to multi-attribute data forests. Let (a, \bar{d}) be an arbitrary fix element of $\mathbb{A} \times \mathbb{D}^{\mathbf{k}}$. Given a forest $\bar{\mathbf{t}}$ and $x, y \in \text{pos}((a, \bar{d})\bar{\mathbf{t}})$, $x, y \neq \epsilon$, we define the satisfaction relation \models , as $\bar{\mathbf{t}}, x \models \varphi$ (resp. $\bar{\mathbf{t}}, (x, y) \models \alpha$) if $(a, \bar{d})\bar{\mathbf{t}}, x \models \varphi$ (resp. $(a, \bar{d})\bar{\mathbf{t}}, (x, y) \models \alpha$). (Note that since $\text{XPath}(* \leftarrow, \downarrow_*, \rightarrow^*, =)$ has no ascending axes, whether $\bar{\mathbf{t}}, x \models \varphi$ or not does not depend on (a, \bar{d}) , we use it as a simple way of defining its semantics.)

The **satisfiability problem** for $\text{XPath}(\mathcal{O}, =)$ (henceforth noted $\text{SAT-XPath}(\mathcal{O}, =)$) is the problem of, given a formula φ of $\text{XPath}(\mathcal{O}, =)$, whether there exists a multi-attribute data tree \mathbf{t} such that $\mathbf{t} \models \varphi$.

3. PROOF SKETCH

The main contribution of this paper is the following.

THEOREM 3.1. *SAT-XPath($* \leftarrow, \downarrow_*, \rightarrow^*, =$) is decidable in $\mathfrak{3EXPSPACE}$.*

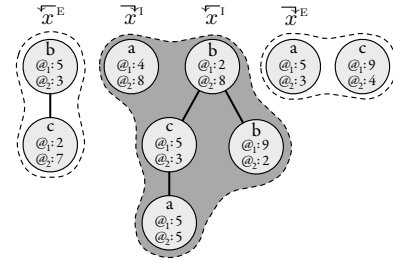


Figure 3: A multi-attribute data forest, with its left and right forests.

We reduce the problem of whether a formula φ of our logic $\text{XPath}(* \leftarrow, \downarrow_*, \rightarrow^*, =)$ is satisfiable, to the problem of whether one can obtain an element with a certain property by repeated applications of operations in some algebra, starting from a basic set of elements. We call it the *derivation problem*. First we introduce the algebra (Section 4), we then solve the derivation problem (Section 5) and finally we show the reduction from the logic into the derivation problem (Section 6).

We first introduce *forest profiles* in Section 4, which constitute the algebra domain. A forest profile is an abstraction of a multi-attribute data forest inside a context, where the context consists of the two (possibly empty) forests that are to the left and to the right. Figure 3 depicts one such possible forest, together with the left and right context forests. A forest profile contains, for each data value d and each path expression, the information of whether d can be reached by the path expression, and *where* it can be reached, either

- inside the main forest, starting from the leftmost root (the node \bar{x}^1 in Figure 3),
- inside the main forest but starting with the rightmost root (the node \bar{x}^1 in Figure 3),
- in the left context forest (starting from the node \bar{x}^E in Figure 3), or
- in the right context forest (starting from the node \bar{x}^E in Figure 3).

In this setting, path expressions are called *patterns* and their navigation is greatly simplified. Patterns can go first to the left, and then down, or first to the right and then down, or only down. They correspond to path expressions like, for example, $\rightarrow^*[a]\rightarrow^*[b]\downarrow_*[c]\downarrow_*[a \vee b]@_1$, or $*\leftarrow[\neg a]\downarrow_*[a]@_2$. Further, node expressions contained in patterns are simple boolean combinations of tests for labels.

A forest profile also keeps track of a set of important data values called the *rigid values*. These are data values that play a determined function in the forest containing the abstracted forest (*i.e.*, in the concatenation of the left, main and right forests). Intuitively, a data value is rigid in a forest if it can be pinpointed by a path expression, in the sense that it is the *only* data value that can be reached with some path expression $\alpha@_i$. At this level of detail, we just mention that some special care must be taken for these rigid data values.

We equip the set of forest profiles with two operations, one that corresponds to concatenating two of the forests being abstracted, and another operation that corresponds to adding a root to the forest, converting it into a tree. This algebra is introduced in Section 4.2. In particular, the root operation is restricted to work only with forest profiles that are from certain set of *consistent* profiles. Consistent profiles will play an important role in the reduction from the logic to the algebra. The idea is that they are those profiles that are not in contradiction with the formula φ to test for satisfiability, that is, that could abstract subforests of a model of φ .

A *root profile*, is a profile that comes from the application of the root operation with a certain label of a certain alphabet \mathbb{A}_{root} of root labels. An *empty profile* is the profile corresponding to the empty forest with an empty context. In Section 5 we define the *derivation problem* for forest profiles as the problem of whether there is a way of obtaining a root profile from the empty profile by repeated applications of the algebra operations.

We show that the derivation problem is decidable in 2EXPSpace in Section 5. We first define a partial ordering on profiles in Section 5.2, this ordering will be of chief importance in our decidability result. We show a series of monotonicity properties that show that the set of derivable profiles is upward-closed. The purpose of the partial ordering is to reduce the derivation problem on the infinite set of forest profiles into a problem on a *finite* set of minimal profiles. The fact that the derivable profiles is upward-closed is indeed a key ingredient for this reduction to work.

However, one problem we need to face is that the ordering has infinite antichains: every two profiles with different set of rigid values are incomparable. We tackle this in Section 5.4, where we show that we can bound the set of rigid values, obtaining an equivalent derivation problem on profiles with a small set of rigid values. Once we obtain this bound, the set of minimal profiles becomes finite, doubly exponential. Next, in Section 5.5 we show that, thanks to the monotonicity properties enjoyed by the algebra, we can work only with minimal elements. Finally, in Section 5.6 we give the concrete saturation-style algorithm that solves the derivation problem using doubly exponential space.

In Section 6 we show that the satisfiability problem for $\text{XPath}(*\leftarrow, \downarrow_*, \rightarrow^*, =)$ can be reduced to the derivation problem in EXPSpace. In Section 6.1 we show a normal form, called *direct unnested normal form*, where direct unnested

path expressions correspond, precisely, to the pattern expressions used in the forest profiles (basically all path expressions are of the form already described). We then show in Section 6.2 that one can reduce, in EXPSpace, the satisfiability problem for formulas in this normal form into the derivation problem, obtaining a 3EXPSpace decidability procedure for $\text{SAT-XPath}(*\leftarrow, \downarrow_*, \rightarrow^*, =)$, obtaining Theorem 3.1.

4. FOREST PROFILES

We define abstractions of forests of multi-attribute data trees. These are called *forest profiles*. They are the main construct in our solution. One must think of a forest profile as the description, for every data value $d \in \mathbb{D}$, of all the possible ways of reaching the data value d via path expressions of $\text{XPath}(*\leftarrow, \downarrow_*, \rightarrow^*, =)$. Some ways of reaching the data value may lie inside the forest being abstracted, and some outside the forest. Take for instance the forest in the middle of Figure 3. For every forest there we identify 4 nodes: the leftmost root, the rightmost root, the node to the left of the leftmost root (if any), and the node to the right of the rightmost root (if any). These are the nodes identified by $\overline{x}^1, \overline{x}^1, \overline{x}^E, \overline{x}^E$ respectively in the figure. The profile of this forest is represented by all the paths that can reach the data value 4, all those that can reach 2, etc. Take as an example the data value 5; this data value can be reached by

- (i) $\rightarrow^*[a]@_1$ from \overline{x}^E ,
- (ii) $*\leftarrow[b]\downarrow_*[a]@_1$ from \overline{x}^1 ,
- (iii) $\rightarrow^*[a]\rightarrow^*[b]\downarrow_*[c]@_1$ from \overline{x}^1 , etc.

Remember that expressions are evaluated in a forest and, for example, an expression starting with \rightarrow^* denotes the possibility to move forward in the sequence of tree roots of the forest. The idea is that we limit ourselves that whenever there are paths departing from \overline{x}^1 or \overline{x}^1 they must be *internal* to the forest (*i.e.*, internal to the gray forest in Figure 3), whenever there are paths from \overline{x}^E or \overline{x}^E they must be *external* to the forest (*i.e.*, either in the forest depicted to the left or to the right of the gray forest in Figure 3).

Let \mathbb{A} be a finite alphabet of **labels**, let $\mathbb{A}_{root} \subseteq \mathbb{A}$ be the set of **root labels**, and let \mathbb{D} be an infinite domain of **data values**. The set $\mathcal{B}(\mathbb{A})$ is the boolean closure of tests for labels from \mathbb{A} . For any $a \in \mathbb{A}$ and $\psi \in \mathcal{B}(\mathbb{A})$, we write $a \models \psi$ if the interpretation assigning *true* to a , and *false* to every other $b \in \mathbb{A}$, satisfies ψ . Let $\mathbf{k} \in \mathbb{N}$ be a fixed natural number, corresponding to the number of **attributes** at each node. We say that $i \in [\mathbf{k}]$ is an **attribute index**. We define the set of **patterns**, as any finite, subword-closed, subset of $(\mathcal{B}(\mathbb{A}))^*$, and we denote it by \mathcal{P} . We generally use the symbols $\alpha, \beta, \gamma, \delta \in \mathcal{P}$ to denote patterns. For every label $a \in \mathbb{A}$ we define the following set of patterns $\sigma_a \subseteq \mathcal{P}$

$$\sigma_a \stackrel{\text{def}}{=} \{\psi_1 \cdots \psi_k \in \mathcal{P} \mid a \models \psi_1 \wedge \cdots \wedge \psi_k\}.$$

Note that $\epsilon \in \sigma_a$. The set of **composed patterns** is

$$\Pi \stackrel{\text{def}}{=} (\mathcal{P} \setminus \{\epsilon\}) \times \mathcal{P} \times [\mathbf{k}].$$

The intended meaning is that the first component operates on the siblings, the second on a downward path, and the third retrieves a data value from an attribute index. We will sometimes use the symbol \bar{a} to represent elements from

Π , or (α, β, i) if we need to explicit the components of the composed pattern.

A **forest profile** \mathfrak{f} is a tuple

$$\mathfrak{f} = (\overleftarrow{\chi}^E, \overrightarrow{\chi}^1, \overleftarrow{\chi}^1, \overrightarrow{\chi}^E, R)$$

where $R \subseteq \mathbb{D}$, and we call it the set of **rigid values** of \mathfrak{f} , and $\overleftarrow{\chi}^E, \overrightarrow{\chi}^1, \overleftarrow{\chi}^1, \overrightarrow{\chi}^E \subseteq \mathbb{D} \times \Pi$, and we call them the set of left/right external/internal **descriptions** respectively. In the example before, one shall interpret (i) as $(5, a, \epsilon, 1) \in \overrightarrow{\chi}^E$, (ii) as $(5, b, a, 1) \in \overleftarrow{\chi}^1$ and (iii) as $(5, a \cdot b, c, 1) \in \overrightarrow{\chi}^1$. We use χ to denote a subset of $\mathbb{D} \times \Pi$; and we write $\bar{\chi}$ (resp. $\bar{\chi}_i$) to denote the 4-uple $(\overleftarrow{\chi}^E, \overrightarrow{\chi}^1, \overleftarrow{\chi}^1, \overrightarrow{\chi}^E)$ (resp. $(\overleftarrow{\chi}_i^E, \overrightarrow{\chi}_i^1, \overleftarrow{\chi}_i^1, \overrightarrow{\chi}_i^E)$). Likewise, we use \mathfrak{f} (resp. \mathfrak{f}_i) to denote $(\bar{\chi}, R)$ (resp. $(\bar{\chi}_i, R_i)$).

We define, for every $\chi \subseteq \mathbb{D} \times \Pi$,

$$\begin{aligned} \chi(d) &\stackrel{\text{def}}{=} \{(\alpha, \beta, i) \in \Pi \mid (d, \alpha, \beta, i) \in \chi\}, \\ \chi(\alpha, \beta, i) &\stackrel{\text{def}}{=} \{d \in \mathbb{D} \mid (d, \alpha, \beta, i) \in \chi\}, \text{ and} \\ \bar{\chi}(d) &\stackrel{\text{def}}{=} (\overleftarrow{\chi}^E(d), \overrightarrow{\chi}^1(d), \overleftarrow{\chi}^1(d), \overrightarrow{\chi}^E(d)). \end{aligned}$$

We define $\text{data}(\mathfrak{f}) \stackrel{\text{def}}{=} R \cup \{d \in \mathbb{D} \mid \bar{\chi}(d) \neq (\emptyset, \emptyset, \emptyset, \emptyset)\}$. We call $\text{data}(\mathfrak{f}) \setminus R$ the set of **flexible values** of \mathfrak{f} . We use the symbol $\bar{\pi}$ to denote $(\overleftarrow{\pi}^E, \overrightarrow{\pi}^1, \overleftarrow{\pi}^1, \overrightarrow{\pi}^E)$ where $\overleftarrow{\pi}^E, \overrightarrow{\pi}^1, \overleftarrow{\pi}^1, \overrightarrow{\pi}^E \subseteq \Pi$. We further say that $\bar{\pi}$ is **the description of $d \in \mathbb{D}$ in \mathfrak{f}** if $\bar{\chi}(d) = \bar{\pi}$.

4.1 Rigid and flexible values

In a forest satisfying some XPath formula, different data values have different roles. We distinguish here two categories of data values: rigid and flexible. Rigid data values are important for the satisfaction of the formula and special care is needed to treat these, whereas flexible values are not crucial, and they can be sometimes removed from the tree. Let us give some more precise intuition. We use the logic XPath to make this intuition clear, but we will then state the definitions in terms of forest profiles.

Given a multi-attribute data forest $\bar{\mathfrak{t}}$ where $\bar{\mathfrak{t}}, i \models \varphi$, suppose there is a data value d such that: there is some position $1 \leq j \leq |\bar{\mathfrak{t}}|$ and some path expression α of φ of the form $\alpha = \rightarrow^* \beta @_k$ or $\alpha = \leftarrow^* \beta @_k$ so that d is the *only* data value that can be reached through α from j . When there is such a d we call it a **rigid value** for j , since the logic can identify it and pinpoint it from the rest of the data values. If d is rigid for at least one position $j \in \{1, \dots, |\bar{\mathfrak{t}}|\}$ we say that d is rigid for $\bar{\mathfrak{t}}$. All the remaining data values of $\bar{\mathfrak{t}}$ (which are the **flexible values**) play the role of assuring that “there are at least two data values reachable through α from position j ” for some α and j . As such, its importance is only relative. In particular, if $\bar{\mathfrak{t}}$ is a forest satisfying φ and containing d as a flexible value, consider $\bar{\mathfrak{t}}'$ as the result of replacing, for some fresh data value d' , every tree \mathfrak{t}'' of $\bar{\mathfrak{t}}$ with the forest $\mathfrak{t}'' \cdot (\mathfrak{t}''[d \mapsto d'])$, where $\mathfrak{t}''[d \mapsto d']$ is the result of replacing the data value d with d' in \mathfrak{t}'' , and leaving all the structures and labels as they were. Indeed, $\bar{\mathfrak{t}}'$ will continue to satisfy φ ; but this is not necessarily true if d was a rigid value. The same notions hold for our algebra on forest profiles. This is a key property that we need to exploit and hence the need to make explicit the set of rigid values of any given profile. We formalize this by defining an ordering on profiles corresponding to the operation just described, so that the forest profile abstracting $\bar{\mathfrak{t}}'$ is bigger than the profile abstracting $\bar{\mathfrak{t}}$. We make explicit (in Lemma 5.1) the aforementioned argument as a monotonicity property of the algebra.

We say that a forest profile $\mathfrak{f} = (\bar{\chi}, R)$ is **valid** if every $d \in \mathbb{D}$ so that $\overrightarrow{\chi}^E(\alpha, \beta, i) = \{d\}$ or $\overleftarrow{\chi}^E(\alpha, \beta, i) = \{d\}$ for some $(\alpha, \beta, i) \in \Pi$, is in R . We define \mathfrak{F} as the set of all valid profiles.

4.2 Algebra

We equip \mathfrak{F} with two operations. The idea is that these operations correspond to the *concatenation* of two forests, and to the addition of a root to a forest (called *rooting*), turning it into a tree.

Preliminaries

The set of **root patterns** of a forest profile \mathfrak{f} , denoted by $[\mathfrak{f}]^{\triangleright}, \llbracket \mathfrak{f} \rrbracket \subseteq \mathcal{P}$ is defined as follows

$$\begin{aligned} [\mathfrak{f}]^{\triangleright} &\stackrel{\text{def}}{=} \{\alpha \mid (d, \alpha, \beta, i) \in \overrightarrow{\chi}^1 \text{ for some } d, \beta, i\}, \\ \llbracket \mathfrak{f} \rrbracket &\stackrel{\text{def}}{=} \{\alpha \mid (d, \alpha, \beta, i) \in \overleftarrow{\chi}^1 \text{ for some } d, \beta, i\}. \end{aligned}$$

Given $P \subseteq \mathcal{P}$ and $\chi \subseteq \mathbb{D} \times \Pi$, we define the **extension of χ by P** , denoted by $P \cdot \chi$, as the set

$$\begin{aligned} P \cdot \chi &\stackrel{\text{def}}{=} \chi \cup \{(d, \alpha' \cdot \alpha, \beta, i) \in \mathbb{D} \times \Pi \mid (d, \alpha, \beta, i) \in \chi, \\ &\quad \alpha' \in P\}. \end{aligned}$$

It is easy to see that the extension operation distributes over union (i.e., $P \cdot (\chi \cup \chi') = P \cdot \chi \cup P \cdot \chi'$ and $(P \cup P') \cdot \chi = P \cdot \chi \cup P' \cdot \chi$).

Fingerprints

We now define the *fingerprint* of a forest profile. It contains a summary information, sufficient to decide whether the tree abstracted by the profile satisfies a formula of XPath—as we show in Section 6.

Let $\mathcal{A} = \{\bar{o}, \bar{o}^{\triangleright}, \bar{o}^{\perp}, \bar{o}\}$. Given a profile $\mathfrak{f} \in \mathfrak{F}$ and $a \in \mathcal{A}$, we define the set $\mathfrak{f} \cdot \chi_a$ as

- $\overleftarrow{\chi}^1 \cup \llbracket \mathfrak{f} \rrbracket \cdot \overleftarrow{\chi}^E$ if $a = \bar{o}$,
- $\overrightarrow{\chi}^1 \cup [\mathfrak{f}]^{\triangleright} \cdot \overrightarrow{\chi}^E$ if $a = \bar{o}^{\triangleright}$,
- $\{(d, \alpha', \beta, i) \in \mathbb{D} \times \Pi \mid \exists \alpha. (d, \alpha, \beta, i) \in \overleftarrow{\chi}^1 \cup \overrightarrow{\chi}^1\}$ if $a = \bar{o}^{\perp}$, or
- $\{(d, \alpha, \beta, i) \in \overrightarrow{\chi}^1 \cup \overleftarrow{\chi}^1 \mid \beta = \epsilon\}$ if $a = \bar{o}$.

Note that $\mathfrak{f} \cdot \chi_a(\alpha, \beta, i)$ is independent of α when $a = \bar{o}^{\perp}$, but it takes an element of Π as argument for the sake of uniformity of notation. The **fingerprint** of a profile \mathfrak{f} , noted $\xi(\mathfrak{f})$, is an element of

$$\begin{aligned} \mathcal{F} &\stackrel{\text{def}}{=} \Pi \times \mathcal{A} \rightarrow \{0, 1, 2+\} \quad \cup \\ &\quad \Pi \times \mathcal{A} \times \Pi \times \mathcal{A} \rightarrow \{0, 1+\}, \end{aligned}$$

where for $\bar{\alpha}, \bar{\alpha}' \in \Pi$, $a, a' \in \mathcal{A}$, we define $\xi(\mathfrak{f})(\bar{\alpha}, a, \bar{\alpha}', a')$ as 0 or 1+ depending on whether $|\mathfrak{f} \cdot \chi_a(\bar{\alpha}) \cap \mathfrak{f} \cdot \chi_{a'}(\bar{\alpha}')| = 0$ or not; and we define $\xi(\mathfrak{f})(\bar{\alpha}, a)$ as 0, 1, or 2+ depending on $|\mathfrak{f} \cdot \chi_a(\bar{\alpha})|$ being 0, 1 or greater than 1 respectively.

We fix the set of **consistent fingerprints**, as a set of fingerprints $\Gamma \subseteq \mathcal{F}$. The usefulness of this set will become apparent in the reduction from XPath to the derivation problem of forest profiles in Section 6.2, but we can anticipate that this set will represent all the profiles abstracting multi-attribute data trees that do not contradict the formula we are trying to satisfy. For the moment, however, the reader may simply consider Γ as a given arbitrary set of fingerprints.

Concatenation

For every two $f_1, f_2 \in \mathfrak{F}$ so that

- (a) $R_1 = R_2$,
- (b) $\vec{\chi}_1^E = \vec{\chi}_2^1 \cup [f_2] \cdot \vec{\chi}_2^E$, and
- (c) $\overleftarrow{\chi}_2^E = \overleftarrow{\chi}_1^1 \cup \llbracket f_1 \rrbracket \cdot \overleftarrow{\chi}_1^E$;

we define the **concatenation** of f_1 and f_2 , denoted as $f_1 + f_2$ as f_3 , where

$$R_3 = R_1 = R_2 \quad (+1)$$

$$\vec{\chi}_3^E = \vec{\chi}_2^E \quad (+2)$$

$$\overleftarrow{\chi}_3^E = \overleftarrow{\chi}_1^E \quad (+3)$$

$$\vec{\chi}_3^1 = \vec{\chi}_1^1 \cup [f_1] \cdot \vec{\chi}_2^1 \quad (+4)$$

$$\overleftarrow{\chi}_3^1 = \overleftarrow{\chi}_2^1 \cup \llbracket f_2 \rrbracket \cdot \overleftarrow{\chi}_1^1. \quad (+5)$$

Notice that

- the concatenation is associative $((f_1 + f_2) + f_3 = f_1 + (f_2 + f_3))$,
- the extension operation \cdot distributes over the concatenation operation $+$ $((f_1 + f_2) \cdot \chi = [f_1] \cdot ([f_2] \cdot \chi))$,
- if $f_1 + f_2 = f_3$ and $f_1, f_2 \in \mathfrak{F}$, then $f_3 \in \mathfrak{F}$.

Rooting

Given $a \in \mathbb{A}$, and $\vec{d} \in \mathbb{D}^k$, we define $(a, \vec{d})f_1 \subseteq \mathfrak{F}$, where $f_2 \in (a, \vec{d})f_1$ if

- (a) $\xi(f_2) \in \Gamma$,
- (b) $\vec{\chi}_1^E = \overleftarrow{\chi}_1^E = \emptyset$,
- (c) $\vec{\chi}_2^1 = \overleftarrow{\chi}_2^1 = \{(d, \alpha, \beta \cdot \gamma, i) \in \mathbb{D} \times \Pi \mid \exists \alpha'. (d, \alpha', \gamma, i) \in \vec{\chi}_1^1 \cup \overleftarrow{\chi}_1^1, \alpha, \beta \in \sigma_a\} \cup \bigcup_{i \in [k]} (\{d(i)\} \times (\sigma_a \setminus \{\epsilon\}) \times \sigma_a \times \{i\})$

We say that f_2 is a **rooting of f_1 with (a, \vec{d})** .

Notice that since the root pattern of any pair of profiles $f_1, f_2 \in (a, \vec{d})f_3$ is the same, it is idempotent and absorbing $([f_1] \cdot [f_2] \cdot \chi = [f_1] \cdot \chi = [f_2] \cdot \chi, [f_1] \cdot \vec{\chi}_1^1 = \vec{\chi}_1^1)$.

4.3 The derivation problem

We define the **empty profile** as $f_\emptyset \stackrel{def}{=} (\emptyset, \emptyset, \emptyset, \emptyset, \emptyset)$. Note that $f_\emptyset \in \mathfrak{F}$. The set of profiles that can be obtained from empty profiles by applying the rooting and concatenation operations is called the set of **derivable profiles**, and noted \mathfrak{D} . We say that f is a **derivable root profile** if $\vec{\chi}^E = \overleftarrow{\chi}^E = \emptyset$ and $f \in (a, \vec{d})f'$ for some $f' \in \mathfrak{D}$, $a \in \mathbb{A}_{root}$ and $\vec{d} \in \mathbb{D}^k$. Let a **derivation tree** for f be a tree t whose every node is labeled by a forest profile and an element from $\mathbb{A} \times \mathbb{D}^k$, except the leaves that are labeled only by the forest profile f_\emptyset and

- the root is labeled with f ,
- every internal node x of t labeled with a forest profile f' and (a, \vec{d}) is so that $f' \in (a, \vec{d})(f_1 + \dots + f_n)$, where f_1, \dots, f_n are the labels of the children of x .

Similarly, a **derivation forest \vec{t}** for f is a forest of derivation trees $\vec{t} = t_1 \dots t_n$ for some profiles f_1, \dots, f_n so that $f = f_1 + \dots + f_n$. Therefore, a profile f is derivable if, and only if, there is a derivation forest for f .

We can now state the *derivation problem*, that is, whether there exists a derivable root profile, given \mathbb{A} , \mathbb{A}_{root} , \mathcal{P} and Γ .

PROBLEM:	The derivation problem
INPUT:	A finite alphabet \mathbb{A} , $\mathbb{A}_{root} \subseteq \mathbb{A}$, a set of patterns \mathcal{P} , a set of fingerprints $\Gamma \subseteq \mathcal{F}$.
QUESTION:	Is there a derivable root profile?

In the next section we show that this problem is decidable. Later, in Section 6, we show that this problem is reducible from SAT-XPath($\ast \leftarrow, \downarrow \ast, \rightarrow \ast, =$).

5. COMPUTING DERIVABLE PROFILES

In this section we solve the derivation problem, showing that it is decidable in 2EXPSpace. To show this problem we work with some partial ordering on forest profiles (Section 5.2) that has some good monotonicity closure properties with our forest profile algebra (Section 5.3). This allows us to reduce the problem to a restricted derivation problem in which solutions can be found by only inspecting profiles with a bounded number of rigid values (Section 5.4), that are minimal elements of the ordering (Section 5.5). These are bounded and computable, allowing us to produce an algorithm solving the problem (Section 5.6).

5.1 Preliminaries

Given $f_1, f_2 \in \mathfrak{F}$ we define that f_1 and f_2 are **equivalent**, and we note it $f_1 \sim f_2$, if there is some bijection $g : \mathbb{D} \rightarrow \mathbb{D}$ so that f_2 is the result of replacing d by $g(d)$ in f_1 ; in this case we write $g(f_1) = f_2$. For a set $C \subseteq \mathfrak{F}$, we write $f \in C$ if there is $f' \sim f$ so that $f' \in C$. Given a forest profile f and two data values $d \in data(f)$, $d' \notin data(f)$, we define $f[d \mapsto d']$ as the result of replacing d by d' in f . Note that $f[d \mapsto d'] \sim f$. Given two data values d, d' we write $f[d \mapsto d, d']$ to denote f' where $R' = R$, $\vec{\chi}'(d') = \vec{\chi}(d)$ and $\vec{\chi}'(e) = \vec{\chi}(e)$ for every other $e \neq d'$. Note that if $d \in data(f) \setminus R$ and $d' \notin data(f)$, we have that if $f \in \mathfrak{F}$ then $f[d \mapsto d, d'] \in \mathfrak{F}$.

We say that a data value $d \in \mathbb{D}$ is an **external data value** of f if $\vec{\chi}^E(d) \cup \overleftarrow{\chi}^E(d) \neq \emptyset$. If further $\vec{\chi}^1(d) \cap \overleftarrow{\chi}^1(d) = \emptyset$, we say that d is a **strict external data value** of f . If $d \in data(f)$ is not a strict external data value, it is then an **internal data value**, and if it is not an external data value, it is then a **strict internal data value**.

5.2 Ordering on profiles

We define a partial order \preceq on forest profiles, that follows from our discussion of Section 4 on the role of flexible and rigid data values. It is the order in which we can make a profile bigger by adding a fresh data value to it, with the same description as that of a flexible data value already contained in it.

Given $f_1, f_2 \in \mathfrak{F}$, we define $f_1 \preceq f_2$ if either $f_1 = f_2$, or there is a flexible datum d of f_1 so that $f_1[d \mapsto d, d'] \preceq f_2$ for some $d' \notin data(f_1)$. Note that \preceq is recursive, reflexive and transitive, and it is hence a partial order.

Note that if $f_1 \preceq f_2$ then $\llbracket f_1 \rrbracket = \llbracket f_2 \rrbracket$ and $[f_1] \cdot = [f_2] \cdot$. Note also that if $f \preceq f'$ then $\xi(f) = \xi(f')$.

We write $f \succsim f'$ if $f \preceq f''$ for some $f'' \sim f'$. We say that a set of forest profiles $G \subseteq \mathfrak{F}$ is **upward closed** (resp. **downward closed**) with respect to \succsim , if for every $f \in G$ and $f' \succsim f$ (resp. $f \succsim f'$), we have $f' \in G$. We write

$$\uparrow G \stackrel{def}{=} \{f \in \mathfrak{F} \mid f \succsim f' \text{ for some } f' \in G\}$$

$$\downarrow G \stackrel{def}{=} \{f \in \mathfrak{F} \mid f' \succsim f \text{ for some } f' \in G\}$$

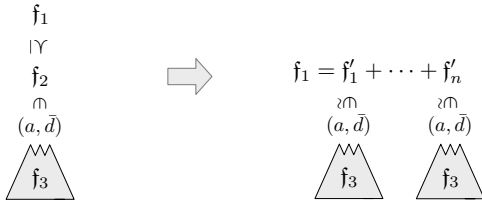


Figure 4: Statement of Lemma 5.1.

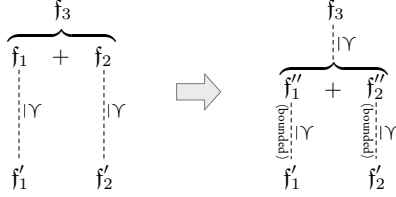


Figure 5: Statement of Lemma 5.3.

for the upward and downward closure of G with respect to \lesssim . We say that G is $\uparrow\downarrow$ -closed, if it is both upward and downward closed, that is, $G = \uparrow\downarrow G$.

5.3 Monotonicity properties

In order to devise an algorithm that tests the existence of a derivable root profile, we will need some monotonicity lemmas evidencing the relationship between \preceq and the rooting and concatenation operations on profiles. The ultimate goal of these lemmas is to restrict the derivation problem to profiles that are *minimal* with respect to \lesssim .

The next Lemma 5.1 states that for any two profiles $f_1 \succeq f_2$, f_1 can be seen as a concatenation of profiles that share the same descriptions of internal values as f_1 , under certain restrictions, as it is shown next. This is a crucial property that follows from our discussion in Section 4.1.

LEMMA 5.1 (FIGURE 4). *For every $f_1 \succeq f_2 \in (a, \bar{d})f_3$, there is $n \in \mathbb{N}$, and $f'_i \preceq (a, \bar{d})f_3$ for every $i \in [n]$ so that*

$$f_1 = f'_1 + \dots + f'_n.$$

The lemma above implies that the set of derivable profiles is upward closed.

LEMMA 5.2. $\mathfrak{D} = \uparrow\mathfrak{D}$.

We finally state two other monotonicity properties that will be required to reduce the derivation problem into a similar problem that works only with minimal profiles in Section 5.5.

We say that a profile f' is a **bounded extension** of a profile f if $f \preceq f'$ and $|\text{data}(f')| \leq |\text{data}(f)| + 3|\Pi|^4$. The following lemma tells us that for any $G \subseteq \mathfrak{F}$ and any profiles $f_1, f_2 \in \uparrow G$, there are bounded extensions f'_1, f'_2 of profiles of G so that $f'_1 + f'_2 \lesssim f_1 + f_2$, as in Figure 5.

LEMMA 5.3 (FIGURE 5). *If $f_1 + f_2 = f_3$ and $f'_1 \preceq f_1$, $f'_2 \preceq f_2$, then $f'_1 + f'_2 \preceq f_3$, for some $f''_1, f''_2 \in \mathfrak{F}$ so that f''_i is a bounded extension of f'_i , for all $i \in \{1, 2\}$.*

A similar lemma holds for the rooting operation.

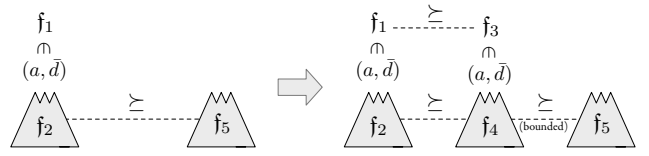


Figure 6: Statement of Lemma 5.4.

LEMMA 5.4 (FIGURE 6). *For every $f_1 \in (a, \bar{d})f_2$ and $f_2 \succeq f_5$, there is $f_4 \succeq f_5$ and $f_3 \preceq (a, \bar{d})f_4$ so that $|\text{data}(f_4)| \leq |\text{data}(f_5)| + |\Pi|^4 + |R_1|$ and $f_3 \preceq f_1$, $f_4 \preceq f_2$.*

5.4 Bounding the rigid values

In this section we show that we can reduce the derivation problem into a similar problem where all the profiles have boundedly many rigid values. This will be combined with the result of the next sections, stating that the derivation problem restricted to profiles with boundedly many rigid values is decidable in 2EXPSpace, to solve the derivation problem.

LEMMA 5.5. *If there is a derivable root profile, then there is a derivation tree for a root profile so that all the profiles in the forest have no more than $2|\Pi|$ rigid values.*

Let \mathfrak{F}_b be the set of all $f \in \mathfrak{F}$ that have no more than $2|\Pi|$ rigid values. Let \mathfrak{D}_b be the set of derivable profiles restricted to \mathfrak{F}_b .

REMARK 5.6. *By Lemma 5.5 and in light of the definition of bounded extension, it follows that Lemma 5.4, when applied to profiles of \mathfrak{F}_b , yields a profile f_4 that is a bounded extension of f_5 .*

By the Lemma just shown, we have the following

LEMMA 5.7. *There is a derivable root profile in \mathfrak{D} if and only if there is a derivable root profile in \mathfrak{D}_b .*

We have then reduced the derivation problem into a simpler problem, the **bounded derivation problem**: testing whether there is a derivable root profile in \mathfrak{D}_b .

REMARK 5.8. *We have that \mathfrak{D}_b is upward closed since \mathfrak{D} is upward closed. That is, $\mathfrak{D}_b = \uparrow\mathfrak{D}_b$.*

Note that \mathfrak{F}_b has boundedly many \lesssim -minimal elements. In the next section we show how to restrict the problem to a problem that uses only these \lesssim -minimal profiles. We will show how this yields a 2EXPSpace algorithm in Section 5.6.

5.5 Restricting to minimal elements

Thanks to the result from the previous section stating that \mathfrak{D}_b is upward closed, we can now show that we can work only with the minimal elements of \mathfrak{F}_b . The main necessary property concerns all those profiles $f' \in \mathfrak{F}_b$ that are ' \lesssim -related' to a profile $f'' \in \mathfrak{D}_b$, in the sense that $f' \succsim f \preceq f'' \in \mathfrak{D}_b$ for some f . (Note that this set of profiles is precisely $\uparrow\downarrow\mathfrak{D}_b$.) The property states that the forest profiles algebra preserves the \lesssim -relatedness.

Given $G \subseteq \mathfrak{F}_b$, let

$$R_{up}^{(a, \bar{d})}(G) \stackrel{\text{def}}{=} \{f \in \mathfrak{F}_b \mid f \in (a, \bar{d})f', f' \in G\}$$

for $(a, \bar{d}) \in \mathbb{A} \times \mathbb{D}^k$,

$$\begin{aligned} \mathbf{R}_{up}(G) &\stackrel{def}{=} \bigcup_{(a,\bar{a}) \in \mathbb{A} \times \mathbb{D}^k} \mathbf{R}_{up}^{(a,\bar{a})}(G), \\ \mathbf{R}_+(G) &\stackrel{def}{=} \{f \in \mathfrak{F}_b \mid f = f_1 + f_2 \text{ where } f_1, f_2 \in G\}, \\ \mathbf{R}(G) &\stackrel{def}{=} \mathbf{R}_{up}(G) \cup \mathbf{R}_+(G). \end{aligned}$$

LEMMA 5.9. $\mathbf{R}(\uparrow\downarrow\mathcal{D}_b) \subseteq \downarrow\mathcal{D}_b$.

5.6 The algorithm

In this section we show how to compute, in 2EXPSpace, whether there exists a derivable root profile in \mathcal{D}_b , solving thus the derivation problem.

For $G \subseteq \mathfrak{F}_b$, we define $G^{\sim} \stackrel{def}{=} \{f \mid f \sim f' \text{ for some } f' \in G\}$. We define G/\sim as the set containing one representative profile of G for each \sim -equivalence class. We define $\text{MIN}(G)$ as the set of \preceq -minimal elements of G ,

$$\text{MIN}(G) \stackrel{def}{=} \{f \in G \mid \text{for all } f' \in G \text{ so that } f' \preceq f \text{ we have } f \sim f'\}.$$

For any $f \in \mathfrak{F}$, we write $|f|$ —the **size** of f —, as the size needed to write f . Note that for all $f \in \text{MIN}(\mathfrak{F}_b)$, $|f|$ is at most exponential in $|\mathcal{P}|$. For any $G \subseteq \mathfrak{F}$, we write $|G|$ to denote $\sum_{f \in G} |f|$.

Let us define \mathbf{C}_i for every $i \in \mathbb{N}_0$ as

$$\begin{aligned} \mathbf{C}_0 &\stackrel{def}{=} \{f_\emptyset\}, \\ \mathbf{C}_{i+1} &\stackrel{def}{=} \mathbf{C}_i \cup \text{MIN}(\downarrow\mathbf{R}(\uparrow\downarrow\mathbf{C}_i))/\sim. \end{aligned}$$

Let $k_0 \in \mathbb{N}_0$ be the first index so that $\mathbf{C}_{k_0}^{\sim} = \mathbf{C}_{k_0+1}^{\sim}$.

REMARK 5.10. For every $i \in \mathbb{N}_0$, $\mathbf{C}_i \subseteq \text{MIN}(\mathfrak{F}_b)$.

As a consequence of the property of the preceding section, we have that this algorithm computes $\text{MIN}(\downarrow\mathcal{D}_b)$.

LEMMA 5.11. $\mathbf{C}_{k_0}^{\sim} = \text{MIN}(\downarrow\mathcal{D}_b)$.

We further have that this computation is in 2EXPSpace since $|\text{MIN}(\mathfrak{F}_b)/\sim|$ is doubly exponential in $|\mathcal{P}|$, hence we obtain the following.

PROPOSITION 5.12. *The derivation problem is decidable in 2EXPSpace.*

6. FROM XPATH TO FOREST PROFILES

In this section we reduce the satisfiability problem for XPath($\ast\leftarrow, \downarrow_\ast, \rightarrow^\ast, =$) into the derivation problem for forest profiles.

In Section 6.1 we define a normal form for XPath($\ast\leftarrow, \downarrow_\ast, \rightarrow^\ast, =$), called *direct unnested normal form*, and in Section 6.2 we show the reduction from the satisfiability problem of direct unnested XPath($\ast\leftarrow, \downarrow_\ast, \rightarrow^\ast, =$) formulas into the derivation problem for forest profiles.

6.1 Normal forms

We will assume a certain normal form of the formula $\varphi \in \text{XPath}(\ast\leftarrow, \downarrow_\ast, \rightarrow^\ast, =)$ to test for satisfiability. This will simplify the reduction into the derivation problem for forest profiles.

The normal form has two main properties. Firstly, it contains only path expressions that are *direct*, in the sense that the navigation consists in going left and then down, or going right and then down. And secondly, path expressions do not contain data tests as node expressions, in other words the formula is *unnested*. Next, we explain in detail these properties.

Preliminaries

Let $\alpha = a_1 \cdots a_n$ with $n > 0$ be a XPath($\ast\leftarrow, \downarrow_\ast, \rightarrow^\ast, =$) path expression, where for every i , $a_i = [\psi]$ for some node expression ψ , or $a_i \in \{\varepsilon, \ast\leftarrow, \downarrow_\ast, \rightarrow^\ast\}$. We say that α is in **alternating path normal form** if either $\alpha = \varepsilon$, or n is even and for all $1 \leq i \leq n$

- if i is even, $a_i = [\psi]$ for some node expression ψ ,
- if i is odd, $a_i \in \{\ast\leftarrow, \downarrow_\ast, \rightarrow^\ast\}$.

In other words, the path alternates between axes and tests for node expressions. We say that a formula is in alternating path normal form if all its path expressions are in alternating path normal form. Note that one can turn any formula $\varphi \in \text{XPath}(\ast\leftarrow, \downarrow_\ast, \rightarrow^\ast, =)$ into an equivalent formula φ' in alternating path normal form in polynomial time, using the equivalences

$$\begin{aligned} \langle [\psi]\alpha @_i \odot \beta @_j \rangle &\equiv \psi \wedge \langle \alpha @_i \odot \beta @_j \rangle \text{ for } \odot \in \{=, \neq\}, \\ \langle \alpha @_i \odot [\psi]\beta @_j \rangle &\equiv \psi \wedge \langle \alpha @_i \odot \beta @_j \rangle \text{ for } \odot \in \{=, \neq\}, \quad (\Delta) \\ \alpha[\psi_1][\psi_2]\beta &\equiv \alpha[\psi_1 \wedge \psi_2]\beta, \text{ and,} \\ \text{if } \alpha\beta \neq \varepsilon, \quad \alpha\beta &\equiv \alpha[\top]\beta \text{ and } \alpha\varepsilon\beta \equiv \alpha\beta. \end{aligned}$$

For simplicity and without any loss of generality we can further assume that all our formulas do not contain formulas of the type $\langle \alpha \rangle$, since it is equivalent to $\langle \alpha @_1 = \alpha @_1 \rangle$. We will henceforth assume that all the formulas we work with are in this form.

We say that a path expression in alternating path normal form is a **rightward path expression**, if it starts with \rightarrow^\ast and all the axes in it are \rightarrow^\ast (similarly with **leftward**, **downward** and $\ast\leftarrow, \downarrow_\ast$). Notice that, for example, a leftward expression may contain node tests using rightward or downward axes. For example, $\ast\leftarrow[\downarrow_\ast[a]]\ast\leftarrow[b]$ is a leftward expression while $\ast\leftarrow[a]\downarrow_\ast[\ast\leftarrow[a]]$ is not.

Direct normal form

The object of the direct normal form is to avoid having unnecessary mixed directions in path formulas, that use perhaps \rightarrow^\ast and $\ast\leftarrow$ in the same expression, or that contain a $\ast\leftarrow$ (or \rightarrow^\ast) axis after a \downarrow_\ast axis. That is, we avoid having formulas like

$$\langle \rightarrow^\ast[a]\ast\leftarrow @_1 = \downarrow_\ast[b]\rightarrow^\ast @_2 \rangle$$

in favor of equivalent formulas with a more *direct* navigation, like

$$\begin{aligned} \langle \rightarrow^\ast[\langle \rightarrow^\ast[a] \rangle @_1 = \downarrow_\ast[\langle \ast\leftarrow[b] \rangle @_2] \rangle \vee \\ \langle \langle \rightarrow^\ast[a] \rangle \ast\leftarrow @_1 = \downarrow_\ast[\langle \ast\leftarrow[b] \rangle @_2] \rangle. \end{aligned} \quad (\ddagger)$$

In the formula above we factor the loops that may be in the navigation of the path expression to obtain a simple navigation that goes in only one horizontal direction.

We say that a formula $\varphi \in \text{XPath}(\ast\leftarrow, \downarrow_\ast, \rightarrow^\ast, =)$ is in **direct normal form**, if every path expression is ε , or of the form $\alpha.\beta$, where $\alpha.\beta \neq \varepsilon$ (i.e., it is not the empty string), α is leftward, rightward or empty, and β is downward or empty. Note that, strictly speaking, the formula (\ddagger) is not in direct normal form since its second disjunct is not in alternating path normal form, but the equivalent alternating path expression—using (Δ) —is in direct normal form.

LEMMA 6.1 (DIRECT NORMAL FORM). *There exists an exponential time translation that for every node expression*

$\varphi \in \text{XPath}(*\leftarrow, \downarrow_*, \rightarrow^*, =)$ returns an equivalent node expression ψ in direct normal form.

Unnested normal form

The second normal form consists in having formulas without nesting of data tests. That is, we avoid treating formulas like, for example

$$\langle \downarrow_*[\underbrace{*\leftarrow[a]@_1 \rightarrow^*[b]@_1}_{\text{nested data test}}]@_1 = \rightarrow^*[c]@_2 \rangle.$$

If a formula is such that all its path expressions α contain only (boolean combinations of) tests for labels we call it a **non-recursive** formula.

We say that φ is in unnested normal form if $\varphi = \varphi_1 \wedge \varphi_2$ where $\varphi_1 \in \mathcal{B}(\mathbb{A})$ and φ_2 is a conjunction of tests of the form “if a node has some of the labels $\{a_1, \dots, a_n\}$ then it satisfies ψ ” for some non-recursive formula ψ and labels $a_1, \dots, a_n \in \mathbb{A}$. Formally, φ_2 contains a conjunction of tests of the form

$$\neg \langle \downarrow_*[\tau \wedge \neg\psi] \rangle$$

for τ a disjunction of labels and ψ a non-recursive formula. Given $\varphi = \varphi_1 \wedge \varphi_2$ in unnested normal form, we write $\gamma_\varphi(a)$ for $a \in \mathbb{A}$ to denote the function where $\gamma_\varphi(a)$ is the conjunction of all the formulas ψ such that φ_2 contains $\neg \langle \rightarrow^*[\tau \wedge \neg\psi] \rangle$ as a subformula, for some disjunctive formula τ containing the label a .

Then, we obtain the following.

LEMMA 6.2 (UNNESTED NORMAL FORM). *There exists an exponential time translation that for every formula $\eta \in \text{XPath}(*\leftarrow, \downarrow_*, \rightarrow^*, =)$ returns a formula φ in unnested normal form such that η is satisfiable iff φ is satisfiable. Further, the translation of a formula in direct normal form is in direct normal form.*

COROLLARY 6.3. *About the translation of Lemma 6.2:*

1. *The set of path subformulas resulting from the translation has cardinality polynomial in η .*
2. *Every path subformula resulting from the translation can be written using polynomial space.*

6.2 Reduction to the derivation problem

In this section we show how we can reduce the satisfiability problem of direct unnested $\text{XPath}(*\leftarrow, \downarrow_*, \rightarrow^*, =)$ formulas into the derivation problem for forest profiles.

Let us fix $\phi = \phi_1 \wedge \phi_2$ in direct unnested normal form, where \mathbb{A} as the finite alphabet, \mathbf{k} as the number of attributes, \mathbb{D} as any infinite domain, and \mathbb{A}_{root} is the set of all $a \in \mathbb{A}$ that make ϕ_1 true.

Given a pattern $\alpha = \psi_1 \cdots \psi_k \in \mathcal{P}$, and an axis $o \in \{*\leftarrow, \downarrow_*, \rightarrow^*\}$, we can convert α into a path expression as follows:

$$\begin{aligned} P_o(\epsilon) &\stackrel{def}{=} \epsilon && \text{if } k = 0, \\ P_o(\psi_1 \cdots \psi_k) &\stackrel{def}{=} o[\psi_1] o \cdots o[\psi_k] && \text{if } k > 0. \end{aligned}$$

Note that P_o is injective.

Let us define \mathcal{P}_ϕ as the set of patterns consisting of

- the constant \top and the empty string ϵ ,
- ψ , for every $\psi \in \mathcal{B}(\mathbb{A})$ that is a subformula of ϕ ,

- every $\alpha \in (\mathcal{B}(\mathbb{A}))^*$ so that $P_{\rightarrow^*}(\alpha)$, $P_{*\leftarrow}(\alpha)$, or $P_{\downarrow_*}(\alpha)$ is a substring of ϕ .

It follows that \mathcal{P}_ϕ is finite and subword-closed.

For any direct non-recursive formula ψ that is a boolean combination of subformulas of ϕ and forest profile \mathbf{f} , we define $\mathbf{f} \vdash \psi$ as follows. If $\psi \in \mathbb{A}$, then $\mathbf{f} \vdash \psi$ if and only if there is some $d \in \mathbb{D}$ and $i \in [\mathbf{k}]$ so that $(\psi, \epsilon, i) \in \overline{\chi}^1(d)$. For all the boolean cases \vdash is homomorphic. Suppose now that $\psi = \langle \alpha \cdot \beta @_i \neq \gamma \cdot \delta @_j \rangle$ where α is leftward, ϵ or empty, γ is rightward, ϵ or empty, and β, δ are downward or empty. We define $\mathbf{f} \vdash \psi$ if there are some $d, d' \in \mathbb{D}$ so that $d \neq d'$ and

- if $\alpha = \epsilon$ or $\alpha = \epsilon$, $(\top, P_{\downarrow_*}^{-1}(\beta), i) \in \overline{\chi}^1(d)$,
- if $\alpha \neq \epsilon$, $\alpha \neq \epsilon$, $(P_{*\leftarrow}^{-1}(\alpha), P_{\downarrow_*}^{-1}(\beta), i) \in (\llbracket \mathbf{f} \rrbracket \cdot \overline{\chi}^E \cup \overline{\chi}^1)(d)$,
- if $\gamma = \epsilon$ or $\gamma = \epsilon$, $(\top, P_{\downarrow_*}^{-1}(\delta), j) \in \overline{\chi}^1(d')$,
- if $\gamma \neq \epsilon$, $\gamma \neq \epsilon$, $(P_{\rightarrow^*}^{-1}(\gamma), P_{\downarrow_*}^{-1}(\delta), j) \in (\llbracket \mathbf{f} \rrbracket \cdot \overline{\chi}^E \cup \overline{\chi}^1)(d')$.

Note that if $\alpha = \epsilon$ then $\beta = \epsilon$ (resp. with γ and δ). If both α and γ are rightwards or leftwards it is defined in an analogous way. The case for $=$ is also analogous, where $d = d'$. The idea is that $\mathbf{f} \vdash \psi$ makes only sense when the derivation forest for \mathbf{f} is a tree, and the multi-attribute data tree \mathbf{t} associated to the derivation tree is so that $\mathbf{t} \models \psi$.

For example, testing ψ is the same as testing if there is some pattern $(\psi, -, -)$ in $\overline{\chi}^1$ or $\overline{\chi}^1$. In a similar way, checking a formula like

$$\langle \rightarrow^*[a] \downarrow_*[b] @_1 = \downarrow_*[c] @_2 \rangle$$

reduces to checking if there is a data value $d \in \mathbb{D}$ that can be reached with $(\top, c, 2)$ in the main forest (i.e., in $\overline{\chi}^1$ or $\overline{\chi}^1$), and either

- d can be reached by $(a, b, 1)$ in the main forest, that is, $(a, b, 1) \in \overline{\chi}^1$ (or equivalently $\overline{\chi}^1$), or
- d can be reached in the right forest by $(a, b, 1)$, where a could be tested in the main forest (i.e., $a \in \llbracket \mathbf{f} \rrbracket$), that is, $(a, b, 1) \in \llbracket \mathbf{f} \rrbracket \cdot \overline{\chi}^E$.

Note that checking $\mathbf{f} \vdash \psi$ takes polynomial time in the size of \mathbf{f} and ψ . Also, whether $\mathbf{f} \vdash \bigwedge_{a \in \mathbb{A}} (a \Rightarrow \gamma_\varphi(a))$ holds or not depends only on $\xi(\mathbf{f})$.

LEMMA 6.4. *Given a direct non-recursive formula ψ that is a boolean combination of subformulas of ϕ , and two forest profiles $\mathbf{f}, \mathbf{f}' \in \mathfrak{F}$ so that $\xi(\mathbf{f}) = \xi(\mathbf{f}')$ then $\mathbf{f} \vdash \psi$ if, and only if, $\mathbf{f}' \vdash \psi$.*

We can then write $\xi \models \psi$ for $\xi \in \mathcal{F}$ instead of $\mathbf{f} \models \psi$ for any \mathbf{f} so that $\xi(\mathbf{f}) = \xi$. We define the set of consistent profiles Γ_ϕ as all $\xi \in \mathcal{F}$ so that $\xi \vdash \bigwedge_{a \in \mathbb{A}} (a \Rightarrow \gamma_\varphi(a))$. The following lemma follows straight from the above definition of \vdash .

LEMMA 6.5. $\mathbf{f} \vdash \bigwedge_{a \in \mathbb{A}} (a \Rightarrow \gamma_\varphi(a))$ iff $\xi(\mathbf{f}) \in \Gamma_\phi$.

Abstractions.

Given multi-attribute data forests $\bar{\mathbf{t}}_l, \bar{\mathbf{t}}, \bar{\mathbf{t}}_r$, we define

$$\text{abs}(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}, \bar{\mathbf{t}}_r)$$

as the forest profile that abstracts the forest $\bar{\mathbf{t}}$ in the context of the forests $\bar{\mathbf{t}}_l$ to the left and $\bar{\mathbf{t}}_r$ to the right. We

have already discussed the idea of this abstraction in Section 4. For example, for the forest of Figure 3, assuming $\mathcal{P} = \{\top, b \cdot c, b, c, \epsilon\}$, we would obtain an abstraction where

$$\overleftarrow{\chi}^E = \{(5, b, b, 1), (5, b, \epsilon, 1), (3, b, \epsilon, 2), (2, b, c, 1), \dots\}.$$

We have that abs is basically an algebra morphism between multi-attribute data forests with rooting and concatenation and forest profiles with profile rooting and profile concatenation. Further, the profile $abs(\epsilon, \mathbf{t}, \epsilon)$ is a derivable root profile whenever $\mathbf{t} \models \phi$; and every derivable root profile is the abstraction of some tree \mathbf{t} so that $\mathbf{t} \models \phi$. As a corollary from these properties, we have the following.

COROLLARY 6.6. *There is a derivable root forest profile if, and only if, ϕ is satisfiable.*

By the above Corollary 6.6 and Proposition 5.12, we can check in 2EXPSpace if there is a derivable root profile. This is 2EXPSpace in the size of \mathcal{P}_ϕ . Although bringing a formula φ into direct unnested normal form may result in a doubly exponential formula, by Corollary 6.3 it can be stored in exponential space, and \mathcal{P}_ϕ is then singly exponential. Hence, the procedure is 3EXPSpace in the original formula φ . Thus, the decision procedure is in 3EXPSpace and Theorem 3.1 follows.

Note that if the input formula is in direct normal form then we save one exponential in the reduction and we hence obtain a 2EXPSpace decision procedure.

THEOREM 6.7. *The satisfiability problem for formulas of $XPath(\leftarrow, \downarrow_*, \rightarrow^*, =)$ in direct normal form is decidable in 2EXPSpace.*

7. DISCUSSION

We have shown that XPath with downward, rightward and leftward reflexive-transitive axes is decidable. To show this, we devised an algebra with good monotonicity properties. This seems to be the right kind of approach to work with transitive relations, and it generalizes and simplifies, in some aspects, the work of [11].

Upward axes

One natural question that stems from the result presented here is whether it can be extended to work with an upward axis as well. However, we claim (without a proof) that already SAT- $XPath(\uparrow^*, \rightarrow^*, =)$ has a non-primitive recursive lower bound. Indeed, this can be proved by reusing the results on lower bounds of [14]. The cited work shows that XPath with one non-reflexive transitive axis is enough to prove non-primitive recursiveness provided that the axis is functional (*i.e.*, the transitive closure of an axis like $\rightarrow, \leftarrow, \uparrow$ but unlike \downarrow). Here, however, we feature reflexive-transitive axes instead of only transitive. Therefore, in principle we cannot use this result. However, one can somehow code \uparrow^+ with $\rightarrow^*[a]\uparrow^*[-a]$ for some label a . We leave the proof of this claim for the journal version of the present work.

By the previous claim, although it could be that full transitive XPath is decidable, it would have a non-primitive recursive lower bound. We can then answer negatively to the conjecture proposed in [11, Conjecture 2], stating that $XPath(\leftarrow, \downarrow_*, \uparrow^*, \rightarrow^*, =)$ be decidable in elementary time.

Future work

- The present work can be seen as a step forward in answering [11, Conjecture 1], suggesting that the extension of $XPath(\leftarrow, \downarrow_*, \rightarrow^*, =)$ with the *child* axis is decidable with elementary complexity. Our approach may perhaps be extended to handle the child relation.
- We suspect that $XPath(\leftarrow, \downarrow_*, \rightarrow^*, =)$ is in fact hard for 2EXPSpace, even when the formulas are in direct normal form, and hence that SAT-direct- $XPath(\leftarrow, \downarrow_*, \rightarrow^*, =)$ is 2EXPSpace-complete.
- We would also like to investigate further the approach taken in this paper to attempt to generalize it to work with the class of reflexive-transitive closures of regular languages.

8. REFERENCES

- [1] Vince Bárány, Mikołaj Bojańczyk, Diego Figueira, and Paweł Parys. Decidable classes of documents for XPath. In *IARCS Annual Conference on Foundations of Software Technology and Theoretical Computer Science (FSTTCS'12)*, Leibniz International Proceedings in Informatics (LIPIcs), Hyderabad, India, 2012. Leibniz-Zentrum für Informatik.
- [2] Michael Benedikt, Wenfei Fan, and Floris Geerts. XPath satisfiability in the presence of DTDs. *Journal of the ACM*, 55(2):1–79, 2008.
- [3] Mikołaj Bojańczyk and Sławomir Lasota. An extension of data automata that captures XPath. In *Annual IEEE Symposium on Logic in Computer Science (LICS '10)*, 2010.
- [4] Mikołaj Bojańczyk, Anca Muscholl, Thomas Schwentick, and Luc Segoufin. Two-variable logic on data trees and XML reasoning. *Journal of the ACM*, 56(3):1–48, 2009.
- [5] James Clark and Steve DeRose. XML path language (XPath). Website, 1999. W3C Recommendation. <http://www.w3.org/TR/xpath>.
- [6] Claire David, Leonid Libkin, and Tony Tan. Efficient reasoning about data trees via integer linear programming. *ACM Transactions on Database Systems*, 37(3):19, 2012.
- [7] Wenfei Fan, Chee Yong Chan, and Minos N. Garofalakis. Secure XML querying with security views. In *ACM SIGACT-SIGMOD-SIGART International Conference on Management of Data (SIGMOD'04)*, pages 587–598. ACM Press, 2004.
- [8] Diego Figueira. Satisfiability of downward XPath with data equality tests. In *ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS'09)*, pages 197–206. ACM Press, 2009.
- [9] Diego Figueira. Forward XPath and extended register automata on data-trees. In *International Conference on Database Theory (ICDT'10)*. ACM Press, 2010.
- [10] Diego Figueira. *Reasoning on Words and Trees with Data*. Phd thesis, Laboratoire Spécification et Vérification, ENS Cachan, France, December 2010.
- [11] Diego Figueira. A decidable two-way logic on data words. In *Annual IEEE Symposium on Logic in Computer Science (LICS'11)*, pages 365–374, Toronto, Canada, 2011. IEEE Computer Society Press.

- [12] Diego Figueira. Alternating register automata on finite data words and trees. *Logical Methods in Computer Science*, 8(1), 2012.
- [13] Diego Figueira. Decidability of downward XPath. *ACM Trans. Comput. Log.*, 13(4), 2012.
- [14] Diego Figueira and Luc Segoufin. Future-looking logics on data words and trees. In *Int. Symp. on Mathematical Foundations of Comp. Sci. (MFCS'09)*, volume 5734 of *LNCS*, pages 331–343. Springer, 2009.
- [15] Diego Figueira and Luc Segoufin. Bottom-up automata on data trees and vertical XPath. In *International Symposium on Theoretical Aspects of Computer Science (STACS'11)*, Leibniz International Proceedings in Informatics (LIPIcs). Leibniz-Zentrum für Informatik, 2011.
- [16] Floris Geerts and Wenfei Fan. Satisfiability of XPath queries with sibling axes. In *International Symposium on Database Programming Languages (DBPL'05)*, volume 3774 of *Lecture Notes in Computer Science*, pages 122–137. Springer, 2005.
- [17] Georg Gottlob, Christoph Koch, and Reinhard Pichler. Efficient algorithms for processing XPath queries. *ACM Transactions on Database Systems*, 30(2):444–491, 2005.
- [18] Marcin Jurdziński and Ranko Lazić. Alternating automata on data trees and xpath satisfiability. *ACM Trans. Comput. Log.*, 12(3):19, 2011.
- [19] Michael Kaminski and Nissim Francez. Finite-memory automata. *Theoretical Computer Science*, 134(2):329–363, 1994.
- [20] Michael Kaminski and Tony Tan. Tree automata over infinite alphabets. In *Pillars of Computer Science*, volume 4800 of *Lecture Notes in Computer Science*, pages 386–423. Springer, 2008.
- [21] Wim Martens and Frank Neven. Frontiers of tractability for typechecking simple xml transformations. *J. Comput. Syst. Sci.*, 73(3):362–390, 2007.
- [22] Frank Neven, Thomas Schwentick, and Victor Vianu. Finite state machines for strings over infinite alphabets. *ACM Trans. Comput. Log.*, 5(3):403–435, 2004.
- [23] Tony Tan. An automata model for trees with ordered data values. In *Annual IEEE Symposium on Logic in Computer Science (LICS'12)*, pages 586–595. IEEE Computer Society Press, 2012.