



HAL
open science

Satisfiability of Downward XPath with Data Equality Tests

Diego Figueira

► **To cite this version:**

Diego Figueira. Satisfiability of Downward XPath with Data Equality Tests. Symposium on Principles of Database Systems (PODS), Jun 2009, Providence (RI), United States. hal-02382925

HAL Id: hal-02382925

<https://hal.science/hal-02382925>

Submitted on 27 Nov 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Satisfiability of Downward XPath with Data Equality Tests ^{*}

Diego Figueira
LSV, ENS Cachan,
CNRS, INRIA Saclay, France

ABSTRACT

In this work we investigate the satisfiability problem for the logic XPath(\downarrow^* , \downarrow , =), that includes all downward axes as well as equality and inequality tests. We address this problem in the absence of DTDs and the sibling axis. We prove that this fragment is decidable, and we nail down its complexity, showing the problem to be EXPTIME-complete. The result also holds when path expressions allow closure under the Kleene star operator. To obtain these results, we introduce a new automaton model over data trees that captures XPath(\downarrow^* , \downarrow , =) and has an EXPTIME emptiness problem. Furthermore, we give the exact complexity of several downward-looking fragments.

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, data-tree, infinite alphabet, data values, BIP automaton

1. Introduction

XPath is arguably the most widely used XML query language. It is implemented in XSLT and XQuery and it is

^{*}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.

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 [4], implemented in most languages and XML packages.

Arguably the most important static analysis problem of a query language is that of optimization, which studies the problem of query containment and query equivalence. In logics closed by boolean operators, these problems reduce to *satisfiability* checking: does a given query express some property? I.e., is there a document where this query has a non-empty result? By answering this question we can decide at compile time whether the query contains a contradiction, and thus whether the computation of the query on the document can be avoided, or if one query can be safely replaced by another one. Moreover, this problem becomes crucial for many applications on security, type checking transformations, and consistency of XML specifications.

Core-XPath (introduced in [6]) is the fragment of XPath 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. We consider an extension of this language with the possibility to make equality and inequality tests between attributes of elements in the XML document. This logic is named Core-Data-XPath in [2], and as shown in [5], its satisfiability problem is undecidable. It is then reasonable to study the interaction between different navigational fragments of XPath with equality tests to be able to find decidable and computationally well-behaved fragments. In the present work, we focus on the downward-looking fragments of XPath, where navigation between elements can only be done in the downward direction.

Our main contribution is that the satisfiability problem for XPath(\downarrow^* , \downarrow , =) is decidable. This is the fragment with data equality and inequality tests, with the \downarrow^* axis that can access descendant nodes at any depth and the \downarrow axis to access child elements. We prove a stronger result, showing the decidability of the satisfiability of $\text{regXPath}(\downarrow, =)$, which is the extension of XPath(\downarrow^* , \downarrow , =) with the Kleene star operator to take reflexive-transitive closures of arbitrary path expressions. Moreover, we nail down the precise complexity showing an EXPTIME decision procedure (recall that XPath(\downarrow , \downarrow^*) is already EXPTIME-hard [1]). In order to

do this, we introduce a new class of automata that captures all the expressivity of $\text{regXPath}(\downarrow, =)$. On the other hand, we prove that the fragment $\text{XPath}(\downarrow^*, =)$ without the \downarrow axis is EXPTIME -hard, even for a restricted fragment of $\text{XPath}(\downarrow^*, =)$ without unions of path expressions. This reduction can only be done by using data equality tests, as the corresponding fragment $\text{XPath}(\downarrow^*)$ without unions is shown to be PSPACE -complete. We thus prove that the satisfiability problem for $\text{XPath}(\downarrow, =)$, $\text{XPath}(\downarrow^*, \downarrow, =)$ and $\text{regXPath}(\downarrow, =)$ are all EXPTIME -complete. Additionally, we present a natural fragment of $\text{XPath}(\downarrow^*, =)$ that is PSPACE -complete. We complete the picture showing that satisfiability for $\text{XPath}(\downarrow, =)$ is also PSPACE -complete. Altogether, we establish the precise complexity for all downward fragments of XPath with and without data tests (cf. Figure 4 in Section 6).

Related work

In [1] there is a study of the satisfiability problem for many XPath logics, mostly fragments without negation or without data equality tests. Also, the fragment $\text{XPath}(\downarrow, =)$ is proved to be in NEXPTIME . We improve this result by providing an optimal PSPACE upper bound. It is also known that $\text{XPath}(\downarrow)$ is already PSPACE -hard, and in this work we match the upper bound showing PSPACE -completeness. Furthermore, in [1] $\text{XPath}(\downarrow, \downarrow^*)$ is proved EXPTIME -complete. In this work we prove that this complexity is preserved in the presence of data values and even under closure with Kleene star. We also consider a fragment that is not mentioned in [1]: $\text{XPath}(\downarrow^*, =)$ and show that $\text{XPath}(\downarrow^*)$ is PSPACE -complete while $\text{XPath}(\downarrow^*, =)$ is EXPTIME -complete. In this case, data tests make a real difference in complexity.

First-order logic with two variables and data equality tests is investigated in [2]. Although in the absence of data values FO^2 is expressive-equivalent to Core-XPath (cf. [8]), FO^2 with data equality tests becomes incomparable with respect to all the data aware fragments treated here. [2] also shows the decidability of a fragment of $\text{XPath}(\uparrow, \downarrow, \leftarrow, \rightarrow, =)$ with sibling and upward axes but restricted to local elements accessible by a ‘one step’ relation, and to data formulæ of the kind $\varepsilon = p$ (or \neq). However, most of the fragments we treat here disallow upward and sibling axes but allow the descendant \downarrow^* axis and arbitrary $p = p'$ data test expressions.

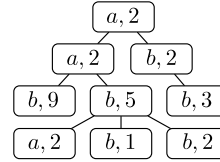
In [7] a fragment of $\text{XPath}(\downarrow, \downarrow^*, \rightarrow, \rightarrow^*, =)$ is treated, denominated ‘forward XPath ’. In the cited work, the full set of downward and rightward axes are allowed, while the fragments treated here only allow the downward axis. As in [2], the language is restricted to data test formulæ of the form $\varepsilon = p$ contrary to the ones studied here, and hence no decidability results can be inferred. It is shown that its satisfiability problem is decidable, but with a non-primitive recursive algorithm, while in our work all the fragments considered are in EXPTIME . The question of whether the forward fragment with arbitrary tests is decidable is still open.

2. Statement of the problem and main result

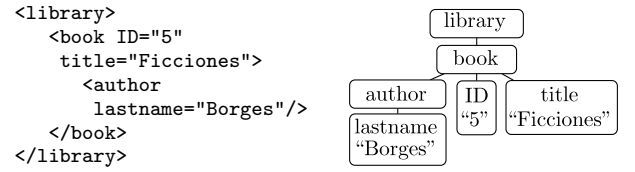
2.1 Data trees

The structure of an XML document can be seen as an unranked tree with attributes and data values in its nodes. We work with an abstraction that we call *data tree*, that is, an unranked finite tree where every node contains a *symbol* from a finite alphabet Σ and a *data value* from some infinite domain Δ . Below we show an example of a data tree with $\Sigma = \{a, b\}$ and $\Delta = \mathbb{N}$.

Example 1.



It is important to mention that this model has only one data value on each node, whilst an XML document element may typically have several (0 or more) *attributes*, each with a data value associated. We address this issue by coding each attribute element by a child as shown next.



In the above example the data values of the nodes tagged with non-attribute elements (library, book, author) may have any data value. For every fragment of XPath with the child axis (\downarrow) we can enforce that attributes are leaves and we can translate any XPath expression on XML to an equivalent one on data trees. For the fragments with the single descendant axis \downarrow^* , this is not true anymore, but a more careful analysis shows that all complexity results still hold in this case. Indeed, any $\text{XPath}(\downarrow^*, =)$ formula on data trees can be seen as an $\text{XPath}(\downarrow^*, =)$ formula on XMLs that makes use of only one fixed attribute, and we can thus transfer the lower bound (the upper bound follows from the more expressive fragment $\text{XPath}(\downarrow^*, \downarrow, =)$). Summing up, all our forthcoming results also hold on arbitrary XML documents with multiple attributes per element. This implies, for example, that the results hold for the satisfiability problem for data trees that may have some nodes with no data values.

We use the standard representation of unranked trees by a nonempty, prefix-closed set T of elements from \mathbb{N}^* such that whenever $x(i+1) \in T$ then $xi \in T$; together with a labeling function $\sigma : T \rightarrow \Sigma$ and a data value function $\delta : T \rightarrow \Delta$. A *data tree model* \mathcal{T} is then a tuple $\langle T, \sigma, \delta \rangle$, and we call $\text{Pos}(\mathcal{T}) = T$ the set of *positions* of \mathcal{T} . We denote by $\mathcal{T}|_x$ the subtree of \mathcal{T} with root x , and $\delta(\mathcal{T}) = \{\delta(x) \mid x \in \text{Pos}(\mathcal{T})\}$. In the Example 1 presented before, $T = \{\varepsilon, 1, 2, 11, 12, 121, 122, 123, 22\}$.

2.2 The logic XPath

We work with a simplification of XPath , stripped of its syntactic sugar. Actually, we consider fragments of XPath that correspond to the navigational part of XPath 1.0 with data equality and inequality. Let us give the formal definition of this logic. XPath is a two-sorted language, with

path expressions (α, β, \dots) and node expressions (φ, ψ, \dots) . The fragment $\text{XPath}(\mathcal{O}, =)$, with $\mathcal{O} \subseteq \{\downarrow, \downarrow^*\}$ is defined by mutual recursion as follows:

$$\alpha ::= o \mid \alpha[\varphi] \mid [\varphi]\alpha \mid \alpha\beta \mid \alpha \cup \beta \quad o \in \mathcal{O} \cup \{\varepsilon\}$$

$$\varphi ::= a \mid \neg\varphi \mid \varphi \wedge \psi \mid \langle \alpha \rangle \mid \alpha \otimes \beta \quad \otimes \in \{=, \neq\}, a \in \Sigma$$

A *formula* of $\text{XPath}(\mathcal{O}, =)$ is either a node expression or a path expression of the logic. $\text{XPath}(\mathcal{O})$ is the fragment $\text{XPath}(\mathcal{O}, =)$ without the node expressions of the form $\alpha \otimes \beta$.

There have been efforts to extend this navigational core of XPath in order to have the full expressivity of MSO, e.g. by adding a least fix-point operator (cf. [9, Sect. 4.2]), but these logics generally lack clarity and simplicity. However, a form of recursion can be added by means of the Kleene star, which allows to take the transitive closure of any path expression. Although in general this is not enough to already have MSO –as shown in [10]–, it does give an intuitive language with counting ability. By $\text{regXPath}(\downarrow, =)$ we refer to the language where path expressions are extended

$$\alpha ::= o \mid \alpha[\varphi] \mid [\varphi]\alpha \mid \alpha\beta \mid \alpha \cup \beta \mid \alpha^* \quad o \in \{\downarrow, \varepsilon\}$$

by allowing the Kleene star on *any* path expression. In terms of expressivity, we see that $\text{XPath}(\downarrow^*, =) \subset \text{XPath}(\downarrow^*, \downarrow, =) \subset \text{regXPath}(\downarrow, =) = \text{regXPath}(\downarrow^*, \downarrow, =)$.

Let $\mathcal{T} = \langle T, \sigma, \delta \rangle$, we define the semantics of XPath:

$$\llbracket \downarrow \rrbracket^{\mathcal{T}} = \{(x, xi) \mid xi \in T\}$$

$$\llbracket \alpha^* \rrbracket^{\mathcal{T}} = \text{the reflexive transitive closure of } \llbracket \alpha \rrbracket^{\mathcal{T}}$$

$$\llbracket \varepsilon \rrbracket^{\mathcal{T}} = \{(x, x) \mid x \in T\}$$

$$\llbracket \alpha\beta \rrbracket^{\mathcal{T}} = \{(x, z) \mid \exists y. (x, y) \in \llbracket \alpha \rrbracket^{\mathcal{T}}, (y, z) \in \llbracket \beta \rrbracket^{\mathcal{T}}\}$$

$$\llbracket \alpha \cup \beta \rrbracket^{\mathcal{T}} = \llbracket \alpha \rrbracket^{\mathcal{T}} \cup \llbracket \beta \rrbracket^{\mathcal{T}}$$

$$\llbracket \alpha[\varphi] \rrbracket^{\mathcal{T}} = \{(x, y) \in \llbracket \alpha \rrbracket^{\mathcal{T}} \mid y \in \llbracket \varphi \rrbracket^{\mathcal{T}}\}$$

$$\llbracket [\varphi]\alpha \rrbracket^{\mathcal{T}} = \{(x, y) \in \llbracket \alpha \rrbracket^{\mathcal{T}} \mid x \in \llbracket \varphi \rrbracket^{\mathcal{T}}\}$$

$$\llbracket a \rrbracket^{\mathcal{T}} = \{x \in T \mid \sigma(x) = a\}$$

$$\llbracket \langle \alpha \rangle \rrbracket^{\mathcal{T}} = \{x \in T \mid \exists y. (x, y) \in \llbracket \alpha \rrbracket^{\mathcal{T}}\}$$

$$\llbracket \neg\varphi \rrbracket^{\mathcal{T}} = T \setminus \llbracket \varphi \rrbracket^{\mathcal{T}}$$

$$\llbracket \varphi \wedge \psi \rrbracket^{\mathcal{T}} = \llbracket \varphi \rrbracket^{\mathcal{T}} \cap \llbracket \psi \rrbracket^{\mathcal{T}}$$

$$\llbracket \alpha = \beta \rrbracket^{\mathcal{T}} = \{x \in T \mid \exists y, z. (x, y) \in \llbracket \alpha \rrbracket^{\mathcal{T}}, (x, z) \in \llbracket \beta \rrbracket^{\mathcal{T}}, \delta(y) = \delta(z)\}$$

$$\llbracket \alpha \neq \beta \rrbracket^{\mathcal{T}} = \{x \in T \mid \exists y, z. (x, y) \in \llbracket \alpha \rrbracket^{\mathcal{T}}, (x, z) \in \llbracket \beta \rrbracket^{\mathcal{T}}, \delta(y) \neq \delta(z)\}$$

For instance, in the model of Example 1,

$$\llbracket \langle \downarrow^* [b \wedge \downarrow [b] \neq \downarrow [b]] \rangle \rrbracket^{\mathcal{T}} = \{\varepsilon, 1, 12\}.$$

We now state the problem we will address.

Definition 1. The *satisfiability problem* $\text{SAT-}\mathcal{L}$ consists in, given an \mathcal{L} -formula η , to decide whether there exists a data tree \mathcal{T} such that $\llbracket \eta \rrbracket^{\mathcal{T}} \neq \emptyset$.

It turns out that –as we are working with downward-looking fragments of XPath– this is equivalent to asking if there is a model where the formula η is satisfied at its root. Moreover, for this problem we can restrict ourselves to the case where η is a node expression. We remind the reader that although we state the problem in terms of data trees, all our results hold on the class of all XML documents with multiple attributes.

2.3 Main contribution

Our main results are the following.

Theorem 1. *SAT-regXPath* $(\downarrow, =)$ is decidable, with complexity in EXPTIME .

Theorem 2. *SAT-XPath* $(\downarrow^*, =)$ is hard for EXPTIME .

Consequently, for a logic $\mathcal{L} \in \{\text{XPath}(\downarrow^*, =), \text{XPath}(\downarrow^*, \downarrow, =), \text{regXPath}(\downarrow, =)\}$, $\text{SAT-}\mathcal{L}$ is EXPTIME -complete.

The strategy of the proof can be outlined as follows.

1. We introduce a model of automata that captures all the expressivity of $\text{regXPath}(\downarrow, =)$. Automata of this new class are called *Bottom-up Interleaved Path automata* (BIP). The automaton relies on an interaction between the runs of two kinds of automata, which corresponds to the two sorts of formulæ of XPath.
2. We show that the translation from $\text{regXPath}(\downarrow, =)$ to BIP automata can be done in PTIME .
3. The main result is that the emptiness problem for the class of BIP automata is in EXPTIME . We show this by a reduction to the non-emptiness problem of classical bottom-up tree automata over trees with bounded branching width. Given a BIP automaton M , we construct a bottom-up tree automaton \mathcal{A} whose states describe the behavior of M for certain data values. We show it is sufficient to consider both the branching width and the number of data values to be polynomially bounded by the BIP automaton M . We show that M is non-empty iff \mathcal{A} is non-empty.
4. Finally, we prove that $\text{XPath}(\downarrow^*, =)$ is EXPTIME -hard by a reduction from the two-player corridor tiling game. Here, the challenge is to be able to move from one square of the game board to the next one, without counting with the ‘ \downarrow ’ operator in the language.

3. A new class of automata

We introduce a new automaton that we call *Bottom-up Interleaved Path automata* (BIP for short) that in its transition function uses another automaton, the *Pathfinder automaton* that runs over the already executed BIP run. This interleaving mechanism of running one automaton as a condition of the transition function of the other one corresponds exactly to the two sorts of formulæ of XPath. Thanks to this duality, we obtain an automaton that captures the whole expressivity of $\text{regXPath}(\downarrow, =)$. It is worth noting that although the BIP automaton does not contain ‘registers’ *per se*, it can do an

unbounded number of equality and inequality tests between any pair of nodes of the tree. However strong this automaton may appear to be, we prove that the emptiness problem is only in EXPTIME.

3.1 Definitions

A *Pathfinder automaton* is a bottom-up non deterministic automaton that basically can only recognize a *path* from some node to the root and retrieve *one* data value from it. The automaton retrieves a data value d with state k if there is a run that starts in a node with data value d and ends at the root with state k . Its definition is very weak as it can only *retrieve* a data value, but it cannot test it against any other. It runs over a data tree over the alphabet $\Sigma = 2^Q$ (i.e., where the labeling function $\sigma : T \rightarrow 2^Q$ tags each node with a subset of Q), where Q is a finite set of symbols that—as we will see shortly—consists of the states of another automaton (the BIP). It is defined as the tuple $\mathcal{P} = \langle K, k_I, Q, \nu \rangle$ where K is a finite set of states, $k_I \in K$ is a distinguished initial state, and ν is the transition function. At each transition, the automaton can either (1) check the presence of some element of Q in the label of the node (we call this a ‘non-moving’ transition), or (2) move up in the tree (we call this a ‘moving’ transition).

$$\nu : (Q \cup \{\text{up}\}) \times K \rightarrow 2^K$$

For example, the non-moving transition $\nu(q_1, k_1) = \{k_2, k_3\}$ indicates that if we are in state k_1 at some node labeled with a set $S \subseteq Q$ such that $q_1 \in S$, then we can label it with one of the states among k_2, k_3 . On the other hand, if $\nu(\text{up}, k_1) = \{k_4\}$ and if we are in state k_1 , then the *father* of this node can be labeled with state k_4 .

Observe that, although this automaton runs on models labeled with *subsets* of Q , the transition function takes only *one* state of Q at a time, its intended meaning being that it applies to any set that contains the specified state. We do so in order to obtain a polynomial time translation from XPath($\downarrow^*, \downarrow, =$) to a pathfinder automaton and to prove the precise upper-bound of EXPTIME. Otherwise, we would have a translation on models with an exponential number of states. This will become clear in the following.

A *run* ρ of a pathfinder $\mathcal{P} = \langle K, k_I, Q, \nu \rangle$ on a data tree $\mathcal{T} = \langle T, \sigma, \delta \rangle$ is a non-empty list of states with positions $\rho \in (K \times \text{Pos}(\mathcal{T}))^+$. We denote by $\rho(i)$ the i th element of the run, starting from 0. The run must be such that $\rho(0) = (k, p)$ with $k = k_I$, and $\rho(N) = (k', p)$ with $p = \varepsilon$ for $N = |\rho| - 1$. For any two positions $\rho(i) = (k', x')$, $\rho(i+1) = (k, x)$ either (1) $x = x'$ and a ‘non-moving’ transition applies between k and k' for x , or (2) $xn = x'$ for some $n \in \mathbb{N}$ and a ‘moving’ transition applies between k and k' for xn . We define that a ‘non moving’ transition *applies* between k and k' for x iff $k \in \nu(k', q)$ for some $q \in \sigma(x)$, and that a ‘moving’ transition applies iff $k \in \nu(k', \text{up})$.

Runs of pathfinder automata are noted by the symbol ρ . The *output* of a run ρ , denoted by $o(\rho)$, is defined as the pair (k, d) , where $\rho(N) = (k, \varepsilon)$, and $d = \delta(p)$ with $\rho(0) = (k_I, p)$, $N = |\rho| - 1$. We also define the *non-moving closure*

of a pathfinder \mathcal{P} w.r.t. a state k and label $S \subseteq Q$ (noted $cl(k, S)$) as the set of states k' that can be reached by ‘non-moving’ transitions on a node labeled S starting with the state k . Formally, $cl(k, S) := \bigcup_{n \geq 0} (f_S^n(\{k\}))$, with $f_S(K) := \{k_1 \mid k_1 \in \nu(q, k_2), q \in S, k_2 \in K\}$. Observe that this set can be built in time polynomial in the set S and the automaton \mathcal{P} . We do not give accepting conditions because this automaton is used by the BIP automaton as we shall see.

Example 2. Consider $\mathcal{P} = \langle \{k_I, k_1, k_{\downarrow 1}, k_2, k_{\downarrow 2}, k_3\}, k_I, \{q_1, q_2, q_f\}, \nu \rangle$ that recognizes $(q_1 q_2)^+$, that is, where $\nu(k_I, q_2) = \{k_2\}$, $\nu(k_2, \text{up}) = \{k_{\downarrow 2}\}$, $\nu(k_{\downarrow 2}, q_1) = \{k_1\}$, $\nu(k_1, \text{up}) = \{k_{\downarrow 1}\}$, $\nu(k_{\downarrow 1}, q_2) = \{k_2\}$, $\nu(k_I, q_1) = \{k_3\}$, $\nu(k_3, \text{up}) = \{k_3\}$. Any run of \mathcal{P} that ends in k_1 retrieves a data value that can be accessed by a ‘path’ like $Q_1^1 Q_2^2 \dots Q_1^{t-1} Q_2^t$ where for any i , $q_1 \in Q_1^i$ and $q_2 \in Q_2^i$. Any run that ends in k_3 retrieves a data value from a node that is labeled by q_1 .

As we show next, the runs of the pathfinder are the basic means for the BIP automaton to test for data (in)equalities.

A *bottom-up Interleaved Path automaton* (BIP) M is a tuple $\langle \Sigma, Q, \mu, F, \mathcal{P} \rangle$, where Σ is a finite set of symbols, Q is a finite set of states, $F \subseteq Q$ is the set of final states, $\mathcal{P} = \langle K, k_I, 2^Q, \nu \rangle$ is a Pathfinder automaton, and $\mu : Q \rightarrow \text{Form}_M$ is the transition function, where Form_M is defined:

$$\varphi ::= a \mid \varphi \wedge \psi \mid \varphi \vee \psi \mid \neg \varphi \mid \exists(k_1, k_2)^{\otimes}$$

where $\varphi, \psi \in \text{Form}_M$, $a \in \Sigma$, $\otimes \in \{=, \neq\}$, $k_1, k_2 \in K$.

Intuitively, $\exists(k_1, k_2)^{\otimes}$ tests for the values retrieved by pathfinder runs (quantified existentially).

A *run* of the BIP automaton over a data tree \mathcal{T} is a labeling function $\lambda : \text{Pos}(\mathcal{T}) \rightarrow 2^Q$. In general if $\mathcal{T} = \langle T, \sigma, \delta \rangle$ we define as $\lambda(\mathcal{T})$ the data tree $\langle T, \lambda, \delta \rangle$. Remember that $\mathcal{T}|_n$ is the subtree that has n as a root, and let $\lambda|_n(i) = \lambda(ni)$. A run λ must fulfill, for every position n , that $q \in \lambda(n)$ iff $\mathcal{T}|_n, \lambda|_n \models \mu(q)$, where \models is defined as follows

- $\mathcal{T}, \lambda \models a$ iff $\sigma(\varepsilon) = a$, that is, the root’s symbol is a ,
- $\mathcal{T}, \lambda \models \neg \varphi$ iff $\mathcal{T}, \lambda \not\models \varphi$, and all boolean connectors are defined in the standard way, and
- $\mathcal{T}, \lambda \models \exists(k_1, k_2)^{\otimes}$ with $\otimes \in \{=, \neq\}$ iff there exist two runs ρ_1, ρ_2 of \mathcal{P} over the run-labeled tree $\lambda(\mathcal{T})$ such that $o(\rho_1) = (k_1, d)$, $o(\rho_2) = (k_2, d')$ and $d \otimes d'$.

Note that the run λ of M on \mathcal{T} is unique by definition. The run is *accepting* if $\lambda(\varepsilon) \cap F \neq \emptyset$.

Example 3. Consider the BIP automaton $\langle \Sigma, \{q_1, q_2, q_f\}, \mu, \{q_f\}, \mathcal{P} \rangle$, where $\Sigma = \{a, b\}$, and \mathcal{P} is the one defined in Example 2. We can define μ to accept the trees that contain two elements accessible by a $(ab)^+$ path from the root, with different data values: $\mu(q_f) = \exists(k_{\downarrow 1}, k_{\downarrow 1})^{\neq} \wedge \neg \exists(k_I, k_3)^{\neq}$, $\mu(q_1) = a$, $\mu(q_2) = b$. Note that it corresponds to the XPath formula $(\downarrow [a] \downarrow [b])^+ \neq (\downarrow [a] \downarrow [b])^+ \wedge \neg \varepsilon \neq \downarrow^* [a]$, and that it accepts the tree of Example 1.

3.2 From regXPath to BIP automata

Theorem 3. *Given a node expression $\eta \in \text{regXPath}(\downarrow, =)$, there is a BIP automaton M such that for any model \mathcal{T} , $\varepsilon \in \llbracket \eta \rrbracket^{\mathcal{T}}$ iff M accepts \mathcal{T} . Moreover, this automaton can be obtained in PTIME.*

PROOF. Let η be a formula of $\text{regXPath}(\downarrow, =)$. We build the BIP automaton $M = \langle \Sigma, Q, \mu, F, \mathcal{P} \rangle$ where $\Sigma = \{a \mid a \text{ a label in } \eta\} \cup \{a_{\perp}\}$, $Q = \{q_{\psi} \mid \psi \text{ is a node expression in } \text{sub}(\eta)\} \cup \{q_{\top}\}$ where $\text{sub}(\eta)$ is the set of subformulae of η , and $F = \{q_{\eta}\}$. Intuitively, for every state q_{ψ} where ψ is a node expression, μ associates it to a formula that is –exactly as ψ – a boolean combination of symbol and data equality tests. For formulae whose principal operator is a boolean connector, the transition is straightforward. E.g., $\mu(q_{\psi_1 \wedge \psi_2}) = \mu(q_{\psi_1}) \wedge \mu(q_{\psi_2})$. If ψ is of the form $\alpha \otimes \alpha'$, then $\mu(q_{\psi}) = \exists(k_{\alpha}, k_{\alpha'})^{\otimes}$, and if $\psi = \langle \alpha \rangle$, $\mu(q_{\psi}) = \exists(k_{\alpha}, k_{\alpha})^{\exists}$.

On the other hand, any path expression $\alpha \in \text{sub}(\eta)$ can be seen as a regular expression over the alphabet $\Sigma_{\eta} = \{e \mid e \text{ is a node expression of } \eta\} \cup \{\downarrow\}$. Then, for any path α we can make the standard PTIME translation of α^r into a NFA over Σ , where α^r stands for the *reverse* of α (it is enough to exactly reverse the symbols, as path expressions are closed under reversal). This is necessary because although XPath path expressions name the path from the root to the leaves, the pathfinder automaton reads the branch from the leaves to the root. We thus obtain a NFA $A_{\alpha} = \langle K_{\alpha}, \Sigma_{\eta}, k_{\alpha}^0, F_{\alpha}, \delta_{\alpha} \rangle$ for each path expression $\alpha \in \text{sub}(\eta)$ where we name the states $K_{\alpha} = \{k_{\alpha}^0, k_{\alpha}^1, \dots\}$. We then define $\mathcal{P} = \langle K, k_I, Q, \nu \rangle$ the pathfinder where $K = \{k_I\} \cup \bigcup_{\alpha} (\{k_{\alpha}\} \cup K_{\alpha})$, and

$$\begin{aligned} \nu(k_I, q_{\varphi}) &= \{k_{\alpha}^0 \mid k_{\alpha}^0 \in K\}, \\ \nu(k_{\alpha}^i, q_{\varphi}) &= \delta_{\alpha}(\varphi, k_{\alpha}^i) \cup \{k_{\alpha} \mid \delta_{\alpha}(\varphi, k_{\alpha}^i) \cap F_{\alpha} \neq \emptyset\}, \\ \nu(k_{\alpha}^i, \text{up}) &= \delta_{\alpha}(\downarrow, k_{\alpha}^i) \cup \{k_{\alpha} \mid \delta_{\alpha}(\varphi, k_{\alpha}^i) \cap F_{\alpha} \neq \emptyset\}. \end{aligned}$$

We can check that all the runs that end in k_{α} –for α some path subformula of η – retrieve a data value of a path accessed via α , and conversely that it can retrieve all of them.

It is not surprising that the translation is so direct, as the BIP automaton mimicks closely XPath semantics. \square

BIP automata are more expressive than $\text{regXPath}(\downarrow, =)$, but if we restrict the definition of BIP to have a bounded number of mutual recursions between the BIP and the pathfinder, we precisely characterize $\text{regXPath}(\downarrow, =)$: for each BIP there is an equivalent $\text{regXPath}(\downarrow, =)$ formula, and vice-versa. In a sense, the restriction disallows the existence of two states q, k (the former of BIP, the latter of pathfinder) in mutual recursion, where k is named in $\mu(q)$ and q in $\nu(k, \cdot)$.

4. Main result

We devote this section to prove that emptiness of BIP automata is in EXPTIME and that the satisfiability problem for $\text{XPath}(\downarrow^*, =)$ is EXPTIME-hard. In this way we obtain that the satisfiability problem for $\text{XPath}(\downarrow^*, =)$, $\text{XPath}(\downarrow^*, \downarrow, =)$ and $\text{regXPath}(\downarrow, =)$ are all EXPTIME-complete.

4.1 Upper bound

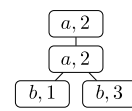
Theorem 4. *Emptiness of BIP automata is in EXPTIME.*

PROOF. The proof consists in a reduction from the BIP emptiness problem into the emptiness of a classical bottom-up non deterministic tree automaton over trees with bounded branching width. Moreover, the tree automaton has at most an exponential number of states, and therefore its emptiness problem can be solved in EXPTIME. A state of this automaton (that we call *extended state* to avoid confusion with states of BIP automata) contains the *description* of the behavior of the BIP automaton w.r.t. some data values. More precisely, each data value d is described by the set of states of the pathfinder by which d can be reached (cf. Ex. 4 below).

We guarantee that for each extended state reachable by the tree automaton, there is a witnessing data tree that consists in any tree that reaches this state at the root, together with an assignment for data values. For the extended states that correspond to the leaves, the witnessing data tree is the leaf with an arbitrary data value. For an extended state of an inner node, the witnessing tree is constructed bottom-up, by identifying the data values that are described in the states of the witnessing subtrees inductively obtained, and ‘merging’ them according to an equivalence relation associated to the transition.

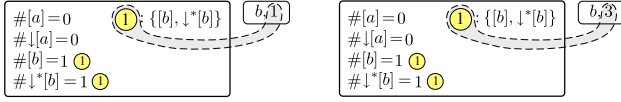
In Proposition 1, we show that given an accepting run of the tree automaton, we can easily build a run of the BIP automaton on a witnessing data tree, and in Proposition 2 we show that if there is an accepting run on the BIP automaton, there must be a model with an accepting run in the tree automaton. As a byproduct of this reduction, we obtain a small model property. We establish that if a BIP automaton accepts at least one data tree, then it accepts in particular a model with polynomial branching width and exponential height, such that for any pair of disjoint subtrees there are only a polynomial number of data values in common.

Example 4. We give the main idea for checking emptiness of BIP automata by abstracting runs. Let M be the BIP that accepts the models satisfying $\varphi = \langle \downarrow [a] \rangle \wedge \neg \langle [a] \rangle \neq \downarrow [a] \wedge \downarrow^* [b] \neq \downarrow^* [b]$, and consider the model below that verifies φ .



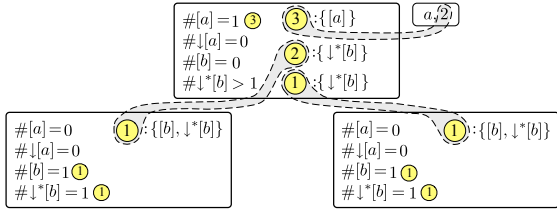
Below, we give an intuition about the nature of the extended states of the tree automaton \mathcal{A}_M built from M , and what would be an accepting run for this particular data tree.

To each subtree, we associate an extended state that contains the description of some data values. Each data value d (here represented by a balloon like Ⓢ) is described by a set of states of the pathfinder automaton. Here, for simplicity’s sake, such states are represented by (sub)path expressions that can reach d . Additionally, the extended state specifies, for each path expression, whether it can retrieve (a) only one data value, (b) more than one, or (c) no data values. In the case (a), we specify which is the only data value retrieved.

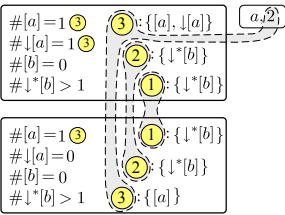


Both leaves have the same extended state. We read the state as follows: there is only one data value reachable by the path $\downarrow^* [b]$, and this data value is denoted by ①. There is also only one data value reached by $[b]$, which is also denoted by ① (hence, they are the same). There are no data values reachable by a path $\downarrow [a]$ (as it is a leaf) or $[a]$ (as it is labeled ‘b’). Finally, we describe the data value denoted by ① with the set of path expressions that may reach it: $\{[b], \downarrow^* [b]\}$. Intuitively, ① makes reference to the current node’s data value, that is why ① and the datum are surrounded by the same \ominus area, to denote that they are equal.

In the next step, we apply a transition, which specifies how the data values between the balloons are merged. In this case we demand that the balloon of the left leaf and that of the right must be *different*, otherwise they would be surrounded in the same \ominus grey area. Consequently, we will have that there are more than one data value reached by $\downarrow^* [b]$, which is reflected in the extended state.



In this transition, we define balloons that describe the same data values as the ones of the leaves (the case of ① and ②), and we define ③ as denoting the data value of the current element $\langle a, 2 \rangle$. Observe that the extended state of the root depends only on (1) the root’s symbol, (2) the descriptions of the data values of its children, and (3) the way of ‘merging’ these data values.



Here, we show how the last transition leads to a final state. Once this transition is performed we can see that it implies that the formula φ is satisfied. In the following, we show how to build this tree automaton systematically.

Let M be a BIP automaton $M = \langle \Sigma, Q, \mu, F, \mathcal{P} \rangle$ with $\mathcal{P} = \langle K, k_I, Q, \nu \rangle$. In the development below we make use of two parameters t_0 (related to the number of data values described in each extended state) and u_0 (the maximum branching width of the witnessing tree) which we assume to be bounded by two polynomials on $|K|$. As usual, we write ‘ \otimes ’ to denote any element of $\{=, \neq\}$. We define a tree automaton $\mathcal{A}_M = \langle \Sigma, Q_{\mathcal{A}}, \tau, F_{\mathcal{A}} \rangle$ where $Q_{\mathcal{A}}$ is the set of extended states, and $F_{\mathcal{A}}$ is the set of final states. Finally, $\tau \subseteq 2_{\leq u_0}^{Q_{\mathcal{A}}} \times \Sigma \times Q_{\mathcal{A}}$ (where $2_{\leq u_0}^{Q_{\mathcal{A}}}$ stands for the set of subsets of $Q_{\mathcal{A}}$ with at most u_0 elements) is the transition function that, given the root’s

symbol and the set of states of its children, labels the root with a state.

Abstracting runs.

An *extended state* is the building block to abstract the runs of M . It is a pair $c = \langle v, D \rangle$ where

- v is a *valuation*, i.e., a function $v : \text{atForm}_M \rightarrow \{0, 1\}$ that specifies which of the formulæ of M hold at the abstracted node. Here, atForm_M is the (finite) subset of *atomic formulæ* (i.e., with no boolean connectors) of Form_M .
- $D = \langle D^\ominus, D^\circ \rangle$ consists in two descriptions of a polynomial number of data values. $D^\circ \in (2^K)^{|K|}$ describes at most $|K|$ data values, and $D^\ominus \in (2^K)^{t_0}$ at most t_0 . We will later detail exactly what this represents.

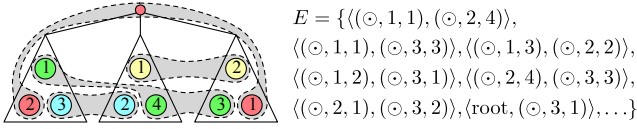
It is easy to check that $|Q_{\mathcal{A}}|$ is exponential in $|M|$. In Example 4, D° is represented by the path expressions α s.t. $\# \alpha = 1$ and the data descriptions associated to them, and D^\ominus by the remaining data descriptions. Each extended state represents a class of models. Before presenting the abstraction, let us fix some notation.

Let $\mathcal{T} = \langle T, \sigma, \delta \rangle$ be a data tree and λ be the unique run of M on \mathcal{T} . We denote by $\lambda(\mathcal{T})$ the data tree labeled with the run $\langle T, \lambda, \delta \rangle$. Given a valuation $v : \text{atForm}_M \rightarrow \{0, 1\}$, we define $\mathcal{C}(v) = \{q \mid q \in Q \text{ with } \mu(q) \text{ true under the valuation } v\}$. If $d \in \Delta$, let $\text{Reach}(d) = \{k \mid \text{there is } \rho \text{ run of } \mathcal{P} \text{ on } \lambda(\mathcal{T}), \text{ such that } o(\rho) = (k, d)\}$. We fix χ to be any bijection $\chi : K \rightarrow [1..|K|]$, and $\chi(k) = i$ stands just as a correlation between $k \in K$ and the element described in position i of the tuple D° (we note this as $D^\circ(i)$). We then say that an extended state $c = \langle v, D \rangle$ *abstracts* the run λ of M on \mathcal{T} (notation: $\mathcal{T} \triangleright c$) iff $\lambda(\varepsilon) = \mathcal{C}(v)$ and the following two conditions are met.

- For every $k \in K$, either $D^\circ(\chi(k)) = \text{Reach}(d_0)$ for $d_0 \in \Delta$ such that for every run ρ of \mathcal{P} on $\lambda(\mathcal{T})$, if $o(\rho) = (k, d)$ then $d = d_0$; or $D^\circ(\chi(k)) = \emptyset$ if there is no such d_0 .
- We only state that D^\ominus describes some data values, with the possibility of even have ‘empty’ descriptions as well if $D^\ominus(i) = \emptyset$. More formally, there are at most t_0 data values $d_1, \dots, d_{t_0} \in \Delta$ where the i th component $D^\ominus(i)$ is the empty set, or $D^\ominus(i) = \text{Reach}(d_i)$, provided that d_i is not already described in D° . I.e., we must ensure that $(\bigcup_i D^\ominus(i)) \cap \{k \mid D^\circ(\chi(k)) \neq \emptyset\} = \emptyset$.

The set $Q_{\mathcal{A}}$ consists of *all* the exponentially many possible extended states, and $F_{\mathcal{A}} = \{\langle v, D \rangle \in Q_{\mathcal{A}} \mid \mathcal{C}(v) \cap F \neq \emptyset\}$. We next describe the transition function τ .

We start with the transitions that correspond to the leaves. Let $t = \langle \emptyset, a, c \rangle$ with c an extended state and a a symbol. We define that $t \in \tau$ iff $\langle a, 1 \rangle \triangleright c$, where $\langle a, 1 \rangle$ is the singleton tree with symbol a and datum 1 (it could be any).



$$E = \{ \langle (\odot, 1, 1), (\odot, 2, 4) \rangle, \langle (\odot, 1, 1), (\odot, 3, 3) \rangle, \langle (\odot, 1, 3), (\odot, 2, 2) \rangle, \langle (\odot, 1, 2), (\odot, 3, 1) \rangle, \langle (\odot, 2, 4), (\odot, 3, 3) \rangle, \langle (\odot, 2, 1), (\odot, 3, 2) \rangle, \langle \text{root}, (\odot, 3, 1) \rangle, \dots \}$$

Figure 1: An example of merging, where D° may be assumed to be empty.

Now we show how to obtain the recursive transitions. We explain how, from $u \leq u_0$ extended states, we can construct the state that is supposed to represent the root of the tree.

For convenience of notation, by c_i we refer to the extended state $\langle v_i, D_i \rangle$. Suppose we have a tuple $t = \langle \{c_1 \dots c_u\}, a, c_0 \rangle$ with $u > 0$. We now show how to check if $t \in \tau$. Remember that \mathcal{A} is *non-deterministic* and the different possibilities for c_0 depend on the way the data values in $c_1 \dots c_u$ are merged. We must then describe how these descriptions are merged together, and ensure that transitions are consistent with the merging. For example, if c_0 describes a datum that can be reached in two steps with label a , or in at least one step with label b (e.g. because we want to check $\Downarrow[a] = \Downarrow^*[b]$) then some c_i must describe a datum accessible through $\Downarrow[a]$, some c_j must describe a datum accessible through $\Downarrow^*[b]$, and these two descriptions must be of the *same* data value.

We define that $t \in \tau$ iff there exists a *merging* \equiv_E such that c_0 is *coherent* w.r.t. \equiv_E and $\{c_1 \dots c_u\}$. We next define what is a ‘merging’ and which are the ‘coherence’ conditions.

Merging data values.

We describe the ways of merging the $u(t_0 + |K|)$ data values described by the (non-empty) elements of $D_1 \dots D_u$. For this purpose, we consider an equivalence relation \equiv_E on $\{(\odot, i, j) \mid D_i^\odot(j) \neq \emptyset\} \cup \{(\circ, i, j) \mid D_i^\circ(j) \neq \emptyset\} \cup \{\text{root}\}$ that describes exactly *how* these data values are going to be collapsed between them and w.r.t. the root’s data value. However we check one condition in \equiv_E . If $k_1 \in D_i^\circ(\chi(k_2))$ then either $D_i^\circ(\chi(k_1)) = \emptyset$ or $(\circ, i, \chi(k_1)) \equiv_E (\circ, i, \chi(k_2))$, as they make reference to the exact same data value (by definition of D°). An example of merging is shown in Fig. 1.

We consider that the *only* data values that can be merged among the trees represented by the u extended states are those described in \equiv_E . In other words, we can consider that the witnessing data tree for c_0 is composed of the witnessing trees for each c_i with the following property. For every c_i we associate a data value to each description, and two data values from two different subtrees i, j are equal if and only if (1) they are described in c_i, c_j respectively, and (2) both descriptions are in same equivalence class of \equiv_E . Observe that this is a strong restriction, the emptiness algorithm of BIP relies on at most $u_0(t_0 + |K|)$ data values at every point of a branch. However, we remark that the automaton M can make at any step, any number of comparisons between any number of data values that can be found in the subtree.

Checking coherence of c_0 with respect to \equiv_E .

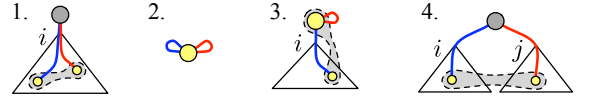


Figure 2: The 4 cases for making $\exists(k_1, k_2)^\circ$ true.

We require that $\langle \{c_1 \dots c_u\}, a, c_0 \rangle \in \tau$ exactly when there exists a merging \equiv_E such that the conditions below hold true. We first define the following relation: $\text{step-up}(k', k)$ iff $k'' \in \nu(k', \text{up})$ and $k \in \text{cl}(k'', \mathcal{C}(v_0))$.

To start with, we must verify that v_0 is correct. For $a' \in \Sigma$, $v_0(a') = 1$ iff $a' = a$. And $v_0(\exists(k_1, k_2)^\circ) = 1$ iff any of the 4 conditions below holds (they are depicted in Fig. 2).

1. There are some k'_1, k'_2 states that retrieve equal data in some subtree, and moving one step up with \mathcal{P} we obtain k_1 and k_2 . I.e., for some $i, v_i(\exists(k'_1, k'_2)^\circ) = 1$, $\text{step-up}(k'_1, k_1)$, and $\text{step-up}(k'_2, k_2)$.
2. Both k_1, k_2 can be obtained as runs that start and end at the root (and hence both carry root’s data value) $k_1, k_2 \in \text{cl}(k_I, \mathcal{C}(v_0))$.
3. k_2 is like in the preceding case, and k_1 retrieves a data value declared in \equiv_E to be equal to the root. For some $i, k_2 \in \text{cl}(k_I, \mathcal{C}(v_0))$, $\exists \ell, \alpha (\alpha, i, \ell) \equiv_E \text{root}, k'_1 \in D_i^\alpha(\ell)$, $\text{step-up}(k'_1, k_1)$ (or the converse, swapping k_1 and k_2).
4. k_1 and k_2 retrieve data values from different subtrees that are equal according to the merging \equiv_E . For some i, j , there exist m, ℓ, α, β s.t. $k'_1 \in D_i^\alpha(m)$, $k'_2 \in D_j^\beta(\ell)$, $(\alpha, i, m) \equiv_E (\beta, j, \ell)$, $\text{step-up}(k'_1, k_1)$, $\text{step-up}(k'_2, k_2)$.

And $v_0(\exists(k_1, k_2)^\circ) = 1$ iff any of the similar conditions previously described in 1, 3, 4 holds (changing $=$ by \neq), or

4. One of the states is *not* described in D° , and hence it retrieves at least 2 different data values, which means that the \neq constraint can be met: $D_i^\circ(\chi(k_1)) = \emptyset$ and $v_0(\exists(k_2, k_2)^\circ) = 1$ (or the converse).

We must check now the coherence of D_0° and D_0° . We concentrate on the former as it is the most involved. Let B_k be the set of all the descriptions of data values that can be reached at the root with state k , $B_k = \{(\alpha, i, j) \mid k' \in D_i^\alpha(j), \text{step-up}(k', k)\}$. We require that $D_0^\circ(\chi(k)) \neq \emptyset$ iff:

- There is no $\langle \odot, i, j \rangle \in B_k$, as this would mean that k retrieves at least *two* data values. Note that any state k' in $D_i^\odot(j)$ retrieves more than one data value (otherwise the datum would be described in $D_i^\circ(\chi(k'))$).
- Any pair of data values that may be reached with state k at the root must be equal: for all $a, b \in B_k, a \equiv_E b$.
- $D_0^\circ(\chi(k))$ consists of every state from the descriptions of B_k , or from any other description that is declared to be equal in \equiv_E . $D_0^\circ(\chi(k)) = \{k' \mid a \in B_k, a \equiv_E \langle \alpha, i, j \rangle, k'' \in D_i^\alpha(j), \text{step-up}(k'', k')\}$.

With respect to D_0° we simply need to state that it cannot contain elements already in D_0° . This completes the definition of τ . It is easy to see that checking all the preceding

conditions consumes at most an exponential amount of time and hence that \mathcal{A}_M can be built in EXPTIME .

We have completely described the tree automaton \mathcal{A}_M . As it contains an exponential amount of extended states, the emptiness problem for this automaton can be computed in EXPTIME . We have that M is empty iff \mathcal{A}_M is empty. \square

Proposition 1. (*Soundness*) For every accepting run of M on \mathcal{T} there is an accepting run of M on some \mathcal{T}' .

PROOF. It is easy to show by induction that for any run of \mathcal{A}_M on \mathcal{T} we can define data values for all the nodes of \mathcal{T} such that the data tree defined corresponds to an *abstraction* of the extended states in the run. The way of merging the data values in the inductive step is completely described by the relation \equiv_E . \square

Proposition 2. (*Completeness*) For every accepting run of M on \mathcal{T} there is an accepting run of \mathcal{A}_M on some \mathcal{T}' .

PROOF. The proof can be sketched as follows. We first show that BIP automata are closed under subtree duplication. More concretely, the data tree $\langle a, d \rangle(t_1, t_2)$ (where $\langle a, d \rangle$ is the root and t_1, t_2 its two immediate subtrees) is accepted by a BIP M iff $\langle a, d \rangle(t_1, t_2, t_2)$ is accepted by M .

We then show a bounded-branching model property:

1. For each atomic formula $\exists(k_1, k_2)^\neq$ that holds at a node z , we mark (at most) two of their immediate subtrees that ‘witness’ this fact. It could be that we only need to mark one or none, if one of the two components k_1 or k_2 is directly witnessed at z . Each marking consists in a label that indicates a state and data value necessary to witness the formula, for example ‘ (k', d') ’. We proceed similarly for $\exists(k_1, k_2)^\equiv$.
2. Moreover, this can be done in such a way that we don’t mark twice the same subtree. We can ensure this by duplicating sibling subtrees if necessary.
3. In this way, we have that each node marks at most some bounded number of subtrees (say N), and at the same time it may be marked by its father. It is easy to see that N is polynomially bounded by the number of states $|K|$. Given a marking of a node z (coming from its father) it could be that (1) the marking is ‘witnessed’ at z , or (2) that it actually needs a subtree. In the case of (2) we add the marking to the corresponding subtree, always making sure that no subtree is marked twice.
4. We then have that each node has marked (at most) $N+1$ immediate subtrees. The rest of the subtrees that are not marked can be safely removed from the tree. This can be done with a top-down algorithm, where the root is the only node to mark at most N nodes (as it has no father).

Finally, we associate an extended state of \mathcal{A}_M to each node. Consider the marking just explained. For each node z we build the extended state, where D^\equiv is completely determined by the tree $\mathcal{T}|_z$, and we use D^\neq to ensure that for each of

the markings (k', d') generated by z , d' is described. There are at most $N+1$ such markings, and we especially choose a sufficiently large size of D^\neq (i.e., of t_0) to be able to accommodate all of them. We can then check that this is indeed a correct accepting run of \mathcal{A}_M . \square

Corollary 1. $\text{SAT-regXPath}(\downarrow, =)$ is in EXPTIME .

PROOF. A direct consequence of Theorems 3 and 4. \square

In the presence of document type definitions.

The BIP automaton can be extended to have transitions that may demand conditions on the states of the child nodes stating, for example, that a node can be labeled by state q_1 if it has a child tagged with state q_2 . It is actually easy to see that this can be simulated using the pathfinder automaton. Furthermore, consider the extension of its formulæ by *positive occurrences* of $\#q \geq n$ where $n \in \mathbb{N}$ is a constant, with the intended meaning that it is verified whenever there are *at least* n child nodes labeled with state q . Similarly, $\#q = 0$ states that there are no children with state q , but formulæ $\#q \leq n$ are not allowed.

It can be checked that a similar emptiness algorithm can be applied, the only difference being that the maximum branching width of the algorithm depends on the greatest constant n_0 used in the definition of the automaton. We obtain then an algorithm of time exponential in n_0 .

So, in the case where n_0 is fixed, or where we consider constraints $\#q \geq n$ with n encoded with a unary representation, we still have an EXPTIME algorithm for emptiness.

Consider the document types definable with a tree automaton on unranked trees with this ‘zero/many’ counting ability, where at each transition we can check either that there is no child with a certain state, that there are at least n with a certain state for some n , or conjunctions and disjunctions of these kind of conditions. It is possible thus to check the satisfiability of any $\text{regXPath}(\downarrow, =)$ under these document types in EXPTIME .

To verify this, we should mention that intersection of two BIP automata M_1 and M_2 can be simply obtained by a BIP with the $Q_{M_1} \times Q_{M_2}$ an defining $\mu_{M_1 \cap M_2}(q_1, q_2) = \mu_{M_1}(q_1) \wedge \mu_{M_2}(q_2)$. Observe that all positive occurrences of the $\#q \geq n$ formulæ remain positive.

Inclusion and equivalence problems.

We can finally mention that as $\text{regXPath}(\downarrow, =)$ is closed under negation and boolean operations, we also get a decision procedure for the equivalence and the inclusion problems for *node expressions* ($\varphi \subset \psi$ iff $\varphi \wedge \neg\psi$ is not satisfiable). However, we cannot solve the problems of *path expressions* inclusion or equivalence with this kind of automata.

4.2 Lower bound

In this section we prove EXP_{TIME}-hardness of satisfiability of $XPath(\downarrow^*, =)$. Remarkably, this logic cannot express a *one step* down in the tree as it does not possess the \downarrow axis, and this will be the major obstacle in the coding.

Theorem 5. *SAT- $XPath(\downarrow^*, =)$ is EXP_{TIME}-hard.*

PROOF. The proof is by reduction from the *two-player corridor tiling game*. An instance of this game consists in a size of the corridor n (encoded in unary), a set of tiles $T = \{T_1, \dots, T_s\}$, a special winning tile T_s , the set of initial tiles $\{T_1^0 \dots T_n^0\}$, and the horizontal and vertical tiling relations $H, V \subseteq T \times T$. The game is played in an $n \times \mathbb{N}$ board where the initial configuration of the first row is given by $T_1^0 \dots T_n^0$. At any moment during the game any pair of horizontal consecutive tiles must be in the relation H and every pair of vertical consecutive tiles in the relation V . The game is played by two players: Abelard and Eloise. Each player takes turn in placing a tile of his choice, filling the board from left to right, from bottom to top, always respecting the horizontal and vertical constraints H and V . Eloise is the first to play, and she wins iff during the game the winning tile T_s is placed on the board. If the game ends without this configuration being reached, or if it runs infinitely, the game is won by Abelard. It is known that deciding whether Eloise has a winning strategy is EXP_{TIME}-complete. For more details on this game we refer the reader to [3].

Representation of a winning strategy.

It is easy to see that in this game Eloise has a winning strategy iff she has a strategy to win before the row s^n of the board is reached (where s is the number of tiles). Then each game between Eloise and Abelard can be coded as a succession of at most s^n rows of n tiles each. Wlog we assume that n is an even number, and hence all odd positions are played by Eloise, while even ones by Abelard. We can then represent a winning strategy as a tree, where at each even position there exists one branch for every possible play of Abelard and where all branches of the tree contain the winning tile T_s .

We must now come up with a way to encode all possible games for all possible choices of Abelard in $XPath(\downarrow^*, =)$, and verify that all of them are won by Eloise and hence that they consist in a winning strategy for Eloise.

Our alphabet consists in the symbols $I_1 \dots I_n$ that indicate the current column of the corridor, the symbols $b_0 \dots b_m$ where $m = \lceil (n+1) \cdot \log(s) \rceil$ that act as *bits* to count from 0 to s^n (it is enough that they count *at least* up to s^n), and the symbols $T_1 \dots T_s$ to code the tile placed at each move. The coding makes use of a symbol $\#$ to separate rows, and an extra symbol $\$$ whose role will be explained later.

Each block of nodes between two consecutive $\#$ codes the evolution of the game for a particular row. Each node labeled I_i has a tile associated, coded as a descendant node T_j with the same data value. In Fig. 3 the first column I_1 of the current row is associated to the tile T_3 , because $\langle T_3, 1 \rangle$ is a descendant of $\langle I_1, 1 \rangle$ with the same data value. Similarly, each occurrence of $\#$ is associated to a number, coded by the

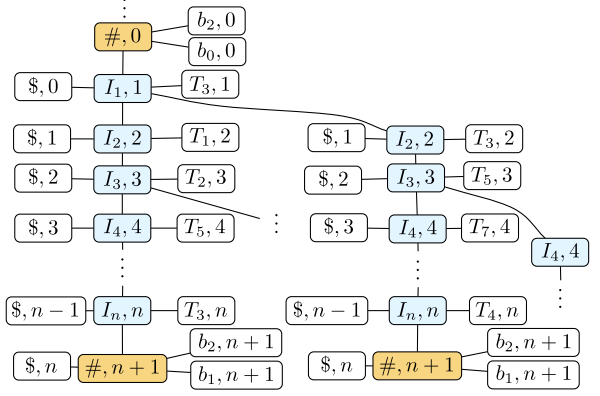


Figure 3: Part of the model coding all the plays of row 5, which is between the $\#$ -element associated to 5 (101 in binary), and the element $\#$ with number 6 (110).

b_i elements. In the example, $\langle \#, 0 \rangle$ is associated to the bits b_0 and b_2 that give the binary number 101.

Finally, the symbol $\$$ is used to delimit the region where the next element of the coding must appear, this will be our way of thinking the *next step* of the coding. Intuitively, between I_i and the $\$$ with the same data value, only a I_{i+1} may appear. This mechanism of coding a very relaxed ‘one step’ is the building block of our coding. As the logic lacks the \downarrow axis, we need to restrict the appearance of the next move of the game to a limited fragment of the model. By means of this element $\$$, we can state, for example, that whenever we are in a I_2 element, then in this restricted portion I_3 must be true by stating $\varepsilon = \downarrow^* [I_3] \downarrow^* [\$]$. In a similar way we can demand that *all* elements verify I_3 (except, perhaps, a prefix of I_2 elements).

However, we cannot avoid having more than one element before the $\$$ as shown in the figure. We may have ‘repeated’ elements or extra branches, but this does not spoil the coding.

We are actually able to force properties for *all* branches and all possible extra elements that the tree may contain. Intuitively, any extra element or branching induces more copies of winning strategies for Eloise.

In Fig. 3 we show an example of a possible extract of the tree between the $\#$ associated to the counting of 5 until the next $\#$ of counting 6. The coding forces a branching as it contains all possible answers of Abelard at even positions.

Building up the coding.

Let us define some useful predicates. $s_\sigma^k(\varphi)$ evaluates φ at a node at k -steps (with our way of coding a step as we have seen before) from the current point of evaluation, given that the current symbol is σ . For this purpose we first define $next(I_i) := I_{i+1}$ (if $i < n$), $next(I_n) := \#$ and $next(\#) := I_1$. Hence, for $a \in \{\#, I_1, \dots, I_n\}$,

$$s_a^0(\varphi) := a \wedge \varphi \quad s_a^{k+1}(\varphi) := a \wedge \varepsilon = \downarrow^* [s_{next(a)}^k(\varphi)] \downarrow^* [\$]$$

Similarly, t_j checks that the tile of the current node I corre-

sponds to T_j , bit_i checks that the i -bit of the counter's binary encoding of a $\#$ -node is one (1), and G forces a property to hold at all nodes of the tree.

$$t_i := \varepsilon = \downarrow^* [T_i] \quad G(\varphi) := \neg \langle \downarrow^* [\neg \varphi] \rangle \quad bit_i := \varepsilon = \downarrow^* [b_i]$$

We now describe all the conditions to force the aforementioned encoding. We also exhibit the XPath formula counterparts of the non-trivial conditions.

1. Every I_i , T_i , and $\#$ along the tree has different data value. As we have only the transitive closure axis, we actually express that whenever there are two elements with the same symbol a such that there is a third element in the middle with another symbol different from a (and then they can be effectively distinguished), they must have different data value. Actually, the fact that there could be a sequence of elements with equal label does not cause any problem. Let us see the case for I_i : $\neg \downarrow^* [I_i \wedge \varepsilon = \downarrow^* [\neg I_i] \downarrow^* [I_i]]$.

2. Every I_i has a next element, unless it contains the winning tile, $G(I_i \wedge \neg t_s \rightarrow s_{I_i}^1(\top)) \wedge G(\# \rightarrow s_{\#}^1(\top))$.

3. $\$$ are leaves, in the sense that no other symbol may appear as descendant: $\neg \langle \downarrow^* [\$ \wedge \downarrow^* [\neg \$]] \rangle$.

4. Every I_i has its corresponding $\$$: $G(I_i \rightarrow \varepsilon = \downarrow^* [\$])$.

5. Each I_i has a unique tile: $G(\neg(t_\ell \wedge t_j))$ for $\ell \neq j$.

6. All I_{i+1} inside a step along a branch must have the same tile. (And a similar condition for I_1 .) That is, for every $i < n$ and $j \neq k$, $G(I_i \rightarrow \neg \varepsilon = \downarrow^* [I_{i+1} \wedge t_j] \downarrow^* [I_{i+1} \wedge t_k] \downarrow^* [\$])$.

7. Between I_i , $i < n$ and its corresponding $\$$ only I_{i+1} may appear. (And also for I_n and $\#$, and for $\#$ and I_1 .) That is, for any $i < n$ and $j \notin \{i, i+1\}$, $G(I_i \rightarrow \neg \varepsilon = \downarrow^* [I_j] \downarrow^* [\$])$, and $G(I_i \rightarrow \neg \varepsilon = \downarrow^* [I_{i+1}] \downarrow^* [I_i] \downarrow^* [\$])$.

8. The tiles match horizontally: for every k and T_i, T_j such that $\neg H(T_i, T_j)$, $\neg \langle \downarrow^* [I_k \wedge t_i \wedge s_{I_k}^1(t_j)] \rangle$. Moreover, the tiles match vertically, for every k and T_i, T_j such that $\neg V(T_i, T_j)$, $\neg \langle \downarrow^* [I_k \wedge t_i \wedge s_{I_k}^{n+1}(t_j)] \rangle$.

9. All the elements corresponding to the *first* row match with $T_1^0 \dots T_n^0$. That is, for all $i \in [1..n]$ and tile T_j such that $\neg V(T_i^0, T_j)$, then $\neg s_{\#}^i(t_j)$ must hold at the root.

10. All possible moves of Abelard are taken into account. For every triple of tiles T_i, T_j, T_k such that $H(T_j, T_k)$, $V(T_i, T_k)$, each time Abelard can play T_k , he *must* play it.

$$\neg \langle \downarrow^* [I_{2\ell} \wedge t_i \wedge s_{I_{2\ell}}^n(I_{2\ell-1} \wedge t_j \wedge \neg s_{I_{2\ell-1}}^1(t_k)) \rangle$$

11. There is no $\#$ element that has all the b_i bits in 1. Because that would mean that Eloise was not able to put a T_s tile in less than s^n rounds.

12. The data value of a $\#$ element is associated to a counter. It is easy to code that the first $\#$ is all-zero. The increment of the counter between two $\#$ is coded as $G(\# \wedge flip(i) \rightarrow zero_{<i} \wedge turn_i \wedge copy_{>i})$, where

$$flip(i) = \bigwedge_{j < i} \neg bit_j \wedge bit_i$$

$$zero_{<i} = \bigwedge_{j < i} \neg s_{\#}^{n+1}(bit_j)$$

$$turn_i = \neg s_{\#}^{n+1}(\neg bit_i)$$

$$copy_{>i} = \bigwedge_{j > i} (bit_j \wedge \neg s_{\#}^{n+1}(\neg bit_j)) \vee (\neg bit_j \wedge \neg s_{\#}^{n+1}(bit_j))$$

This completes the coding. It is easy to see that each one of the formulæ described has a polynomial length on s and n . It can be shown then that Eloise has a winning strategy in the two-player corridor tiling game iff the conjunction of the formulæ just described is satisfiable. \square

5. PSpace fragments

We now turn to some other downward fragments of XPath. We complete the picture of the complexity for all possible combinations of downward axis in the presence and in the absence of data values. We first need to introduce a basic definition that we use throughout the section.

Definition 2. We say that the logic \mathcal{L} has the *poly-depth model property* if there exists a polynomial f such that for every formula $\varphi \in \mathcal{L}$, φ is satisfiable iff φ is satisfiable in a data tree model of depth at most $f(|\varphi|)$.

We can now prove the following statement that we will use to show PSPACE-completeness for XPath($\downarrow, =$).

Theorem 6. *Every fragment \mathcal{L} of regXPath($\downarrow, =$) with the poly-depth model property is in PSPACE.*

PROOF. Suppose that if a formula $\eta \in \mathcal{L}$ is satisfiable in a model, then it is satisfiable in a model of height h with $h \leq f(|\eta|)$ where f is a polynomial.

We can then translate η into a BIP automaton M . We show next how to modify the emptiness algorithm to make it work in non-deterministic polynomial space by means of f . We define an algorithm by recursion on the height of the tree h . The algorithm receives *three* parameters: (1) a BIP automaton M , (2) an extended state c_0 , and (3) h , the maximum height of the tree to reach the extended state. The algorithm must verify that c_0 can be reached in a tree of height at most h in non-deterministic polynomial space in h . Now the emptiness algorithm for \mathcal{A}_M must be done *on the fly*, that is, we do not build the set of all possible extended states. The base case is when $h = 1$. In this case we can easily check the existence of a singleton tree that satisfies c_0 in polynomial space.

Suppose now the height is $h = n + 1$. The algorithm guesses $c_1 \dots c_u$ extended states corresponding to the immediate subtrees and verifies that c_0 and $c_1 \dots c_u$ are in a transition of \mathcal{A}_M . To do this, we guess a relation \equiv_E and test that all the conditions described in the construction of

\mathcal{A}_M seen before are satisfied. Finally, by inductive hypothesis we can check that each one of $c_1 \dots c_u$ extended states are satisfied in a model of depth at most n , and this test can be done in space polynomial in n . In terms of space complexity this algorithm uses (a) the space to store $c_1 \dots c_u$ where u is a polynomial in M and the space required to store each c_i is polynomial in M , (b) the space to store the relation \equiv_E to check their correctness, that can be bounded by $2.u.(2|K|^2 + 3|K| + 2)$ and also remains polynomial in M , and (c) some polynomially bounded space on n (call it $S(n)$) to do the u recursive calls. Then the space required is $S(n+1) = (a) + (b) + S(n)$. It is then immediate that the algorithm is in NPSpace .

The main algorithm can be sketched as follows. Given η , we compile η into M in PTime , we guess an extended state c_0 that contains a final state of M and we check the guessing is correct by calling the algorithm just described. Thus, as $\text{NPSpace} = \text{PSpace}$ the theorem follows. \square

Proposition 3. *SAT-XPath($\downarrow, =$) is PSPACE-complete.*

PROOF. XPath(\downarrow) is shown to be PSPACE-hard in [1]. For the upper bound, we show the poly-depth model property. It is easy to show that if η is satisfiable in \mathcal{T} , then it is satisfiable in $\mathcal{T} \upharpoonright n$ where n is the maximum quantity of nested \downarrow of the formula, and $\mathcal{T} \upharpoonright n$ is the submodel of \mathcal{T} consisting of all the nodes that are at distance at most n from the root. Hence, by Theorem 6, XPath(\downarrow) is in PSPACE. \square

So far we have that, in the presence of data values, the ability to have the descendant axis (\downarrow^*) produces an increase in the complexity from PSPACE to EXPTIME¹. However, we argue that it is not the ability to test for data equality of distant elements what produces this raise in complexity. It is, as a matter of fact, in the ability to test data values against that of the root in formulæ like $\varepsilon = \downarrow^* [a]$. We show that if we actually eliminate this kind of data tests, we can prove the resulting logic to be only in PSPACE.

Definition 3. We denote by XPath($\downarrow^*, =$) $\setminus \varepsilon$ the fragment of XPath($\downarrow^*, =$) where the ε path formulæ are forbidden, and in general where there are no ε -testing in a path (like in $[\varphi] \downarrow^*$), $\alpha ::= \downarrow^* | \alpha[\varphi] | \alpha\beta | \alpha \cup \beta$.

Proposition 4. *SAT-XPath($\downarrow^*, =$) $\setminus \varepsilon$ is PSPACE-complete.*

PROOF (SKETCH). The proof is done by proving the poly-depth model property. The key observation is that any XPath($\downarrow^*, =$) $\setminus \varepsilon$ path expression that is satisfied at a node n of a tree, is also satisfied in any ancestor of n , this is basically because all path expressions start with a \downarrow^* axis. This means that for any pair of nodes n, n' such that n is an ancestor of n' , the set of formulæ of the type $\langle p \rangle, p = p'$ or $p \neq p'$ (with p, p' path expressions) that are satisfied in n' is a *subset* of those that are satisfied in n . Thus, if φ is a formula satisfied in \mathcal{T} , for a given branch there is only a polynomial number of configurations of the (sub)paths in φ verified in each of its nodes. Long branches with a repeated description can actually be pruned into a shorter branch, preserving satisfiability of φ in \mathcal{T} . \square

¹In the case $\text{PSpace} \neq \text{EXPTIME}$.

\downarrow	\downarrow^*	=	Complexity	Details
+			PSpace-complete	Prop 3
	+		PSpace-complete	Prop 5
+	+		EXPTIME-complete	[1]
+		+	PSpace-complete	Prop 3 and [1]
	+	+	EXPTIME-complete	Cor 1, Th 5
+	+	+	EXPTIME-complete	Cor 1, Th 5
regXPath($\downarrow, =$)			EXPTIME-complete	Cor 1, Th 5
XPath($\downarrow^*, =$) $\setminus \varepsilon$			PSpace-complete	Prop 4

All the results hold also in the absence of path unions.

Figure 4: Summary of results.

Proposition 5. *SAT-XPath(\downarrow^*) is PSPACE-complete.*

PROOF (SKETCH). The lower bound is shown by coding the QBF problem. The upper bound is shown via the poly-depth model property. It is slightly involved and requires to show a normal form of the model with the following property. If for some node both $\langle \alpha \rangle$ and $\langle \beta \rangle$ hold, then α and β are witnessed by two *disjoint* branches of polynomial depth. \square

6. Concluding remarks

We have shown the complexity of various downward fragments of XPath. The highest complexity class we obtained is EXPTIME. In the presence of data equality tests, this is a well behaved fragment considering that in the presence of all the axes XPath is undecidable. One important reason for this, is the absence of sibling axis. Actually, in the presence of arbitrary DTDs we can show that the satisfiability problem of the downward fragment is either undecidable, or decidable with a non-primitive recursive algorithm. We have shown however that we can evaluate some restricted fragment of DTDs that cannot express sibling order nor limit the quantity of occurrences of nodes of a certain type, but that can demand, for example, that any a has at least five b children and no c child. By solving the satisfiability problem we are also able to solve the inclusion and equivalence problems of node expressions for free. We leave open the question of whether the inclusion of path expressions (as binary relations) is also decidable in EXPTIME.

We introduced the new class of BIP automata that capture all the expressivity of regXPath($\downarrow, =$). By the proof of decidability, we conclude that there is a very strong normal form of the model for this logic. If a formula $\eta \in \text{regXPath}(\downarrow, =)$ is satisfiable, then it is satisfiable in a model of exponential height and polynomial branching width, whose data values are such that only a polynomial number of data values can be shared between any two disjoint subtrees. This small model property is reflected by the fact that the emptiness of the automaton only depends on a polynomial number of data values at every point of a branch. However, there is no syntactic restriction in the automaton, it can retrieve and compare any number of data values between them and the root's data value at each step of the execution of M .

Acknowledgments. I am grateful to Luc Segoufin and Stéphane Demri for helpful discussions and for critically reading this manuscript.

7. References

- [1] M. Benedikt, W. Fan, and F. Geerts. XPath satisfiability in the presence of DTDs. *J.ACM*, 55(2), 2008.
- [2] M. Bojańczyk, C. David, A. Muscholl, T. Schwentick, and L. Segoufin. Two-variable logic on data trees and XML reasoning. In *PODS*, pages 10–19. ACM, 2006.
- [3] B. S. Chlebus. Domino-tiling games. *J. Comput. Syst. Sci.*, 32(3):374–392, 1986.
- [4] J. Clark and S. DeRose. XML path language (XPath). Website, November 1999. W3C Recommendation. <http://www.w3.org/TR/xpath>.
- [5] F. Geerts and W. Fan. Satisfiability of XPath queries with sibling axes. In *DBPL*, volume 3774, pages 122–137. Springer, 2005.
- [6] G. Gottlob, C. Koch, and R. Pichler. Efficient algorithms for processing XPath queries. *ACM Trans. Database Syst.*, 30(2):444–491, 2005.
- [7] M. Jurdziński and R. Lazić. Alternating automata on data trees and XPath satisfiability. *CoRR*, abs/0805.0330, 2008.
- [8] M. Marx. First order paths in ordered trees. In *ICDT*, volume 3363, pages 114–128. Springer, 2005.
- [9] B. ten Cate. The expressivity of XPath with transitive closure. In *PODS*, pages 328–337. ACM Press, 2006.
- [10] B. ten Cate and L. Segoufin. XPath, transitive closure logic, and nested tree walking automata. In *PODS*, pages 251–260. ACM Press, 2008.

APPENDIX

A. From XML to data tree

As outlined before, XML documents may have multiple attributes with data values on each element, while data trees can only have one. Let us consider now that our finite set of symbols $\Sigma = \Sigma_{attr} \cup \Sigma_{elem}$ is divided between the names for attributes and the symbols of the XML elements. We can now have the models as $\langle T, \sigma, \delta \rangle$ where $\delta \subseteq T \times \Sigma_{attr} \times \Delta$ is a relation that may have a data value for some attribute symbols on any node. Let us consider then the extension of the languages where different attributes may be compared, where node expressions are defined

$$\varphi ::= a \mid \neg\varphi \mid \varphi \wedge \psi \mid \langle \alpha \rangle \mid \alpha @ attr1 \otimes \beta @ attr2$$

where $\otimes \in \{=, \neq\}$, $a \in \Sigma$ and $attr1, attr2 \in \Sigma_{attr}$. Let us call this logic attrXPath. It is easy to see that this language with the expected semantics can well encode any XPath request on an XML document.

However, as already mentioned, each XML document can be coded in a data tree by adding one child for each attribute with its corresponding value. We can force this kind of model using XPath($\downarrow^*, \downarrow, =$) by stating that all the nodes with a symbol from Σ_{attr} are leaves.

$$\varphi_{struct} = \neg \langle \downarrow^* [\bigvee_{s \in \Sigma_{attr}} \wedge \downarrow] \rangle$$

In addition we can transform any XPath formula with attributes like ' $\downarrow^* [a] \downarrow @attr1 = \downarrow [b] @attr2$ ' into a formula ' $\downarrow^* [a] \downarrow \downarrow [attr1] = \downarrow [b] \downarrow \downarrow [attr2]$ ', let us call this translation tr .

We can then decide the satisfiability of a formula ψ of attrXPath on trees with multiple attributes by transforming it to an equivalent one on data trees by ' $tr(\psi) \wedge \varphi_{struct}$ '. We have then an EXPTIME decidability procedure for the full downward fragment of attrXPath even with the Kleene star operator (as the translation is clearly in PTIME).

On the other hand any XPath formula on data trees can be thought of a attrXPath formula that uses at most one attribute. We can then deduce the EXPTIME-hardness result of attrXPath($\downarrow^*, =$) from that of XPath($\downarrow^*, =$).

For the case of attrXPath($\downarrow, =$) we can do the same translation the only difference being that for a formula $\psi \in attrXPath(\downarrow, =)$ the structure can be forced by

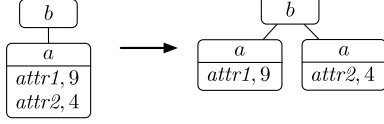
$$\varphi_{struct} = \bigwedge_{0 \leq n \leq d+1} \neg \langle \downarrow^n [\bigvee_{s \in \Sigma_{attr}} \wedge \downarrow] \rangle$$

where d is the maximum quantity of nested \downarrow of ψ . It is easy to see that this forces the requested property for all the portion of the data tree that we are interested in. That is, for the whole region where $tr(\psi)$ can access. This is associated with the poly-depth model property of the logic. We then have that attrXPath($\downarrow, =$) is PSPACE-complete.

Finally, for the case of attrXPath($\downarrow^*, =$) $\setminus \varepsilon$ we are not able

to force the wanted structure of the data tree. However it can be shown that if a formula $\varphi \in \text{attrXPath}(\downarrow^*, =) \setminus \varepsilon$ is satisfiable, then it is satisfiable in a model where each node has at most one attribute.

The basic idea is that any node with multiple attribute can be ‘split’ into several nodes each one of them with only one attribute.



It can be checked that this transformation preserves the satisfaction of all $\text{attrXPath}(\downarrow^*, =) \setminus \varepsilon$ formulæ. We then have a trivial transformation into data trees and the same result of PSPACE-completeness of $\text{attrXPath}(\downarrow^*, =) \setminus \varepsilon$.

B. Characterizing $\text{regXPath}(\downarrow, =)$

We now show that we can easily identify the class of BIP automata that correspond to the logic of $\text{regXPath}(\downarrow, =)$. Intuitively, we need to make certain that there are no mutual unbounded recursion between the main BIP automaton and its pathfinder.

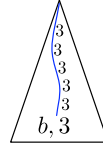
Definition 4. We say that a BIP automaton $M = \langle \Sigma, Q, \mu, F, \mathcal{P} \rangle$ with $\mathcal{P} = \langle K, k_I, Q, \mu_{\mathcal{P}} \rangle$ has the *bounded interleaving property* if there is a partition $Q = Q_1 \cup \dots \cup Q_n$ where $\forall i \neq j : Q_i \cap Q_j = \emptyset$ and $K = K_1 \cup \dots \cup K_m$ where $\forall i \neq j : K_i \cap K_j = \emptyset$, such that for every $i \in [1..n]$, if $q \in Q_i$ then $\mu(q) \in \text{Form}(\bigcup_{j \leq i} K_j)$, and on the other hand for every $i \in [1..m]$ if $k \in K_i$ then $\nu(k, q) = \emptyset \forall q \in \bigcup_{j \geq i} Q_j$.

That is, if we do not have mutual recursion between a state q from the BIP automaton and a state k of the pathfinder automaton it contains.

Proposition 6. *For every BIP automaton M with the bounded interleaving property there exists a $\text{regXPath}(\downarrow, =)$ -formula η such that for every model \mathcal{T} , M accepts \mathcal{T} iff $\llbracket \eta \rrbracket^{\mathcal{T}} \neq \emptyset$.*

PROOF. It is easy to see that the set of runs leading to a state $k \in K_i$ from the pathfinder automaton can be coded into a regular expression over the alphabet $Q_1 \dots Q_{i-1}$. Conversely, every state $q \in Q_i$ of the BIP automaton can be straightforwardly coded into a node expression where, instead of data-expressions like $\downarrow [a] \downarrow [b] = (\downarrow [c])^*$, we have states from the pathfinder of the form $\downarrow k_5 = \downarrow k_7$.

It is then just a matter of replacing states of BIP automaton by node expressions in the path expressions, and vice-versa replacing states of the pathfinder automaton by path expressions in the node expressions. The bounded interleaving property forces that we only need $n + m$ replacement steps. The formula η then consists in the disjunction of all node expressions associated to final BIP states. \square



Remark. Note that the expressivity of the BIP automaton is indeed very close to that of $\text{regXPath}(\downarrow, =)$. The essence of what BIP can express that $\text{regXPath}(\downarrow, =)$ cannot do is the recursiveness between node and path expressions, that would correspond to a ‘star’ operator on the *nesting* of the formula. More concretely, BIP automata can express, for instance, the disjunction of all the formulæ defined by the grammar $A ::= \varepsilon \mid \downarrow [A] \mid b$. That is, the property of having a ‘chain’ of equal data until a b , which is not something expressible in $\text{regXPath}(\downarrow, =)$.

C. Completeness of emptiness of BIP

PROOF OF PROPOSITION 2. We show that whenever there is a model \mathcal{T} accepted by M , there is another one \mathcal{T}' such that it is both accepted by M and by \mathcal{A}_M . Note that \mathcal{A}_M can run on data trees by simply ignoring the values of the nodes. We will show that not only there is an accepting run of \mathcal{A}_M on \mathcal{T}' , but there is an accepting run with some extra information, where we can associate a concrete data value of the tree \mathcal{T}' to each of the descriptions of the extended state. We will call this a ‘data run’ of \mathcal{A}_M .

We first need to introduce the following definition of equivalence between trees.

Definition 5. [Strong S-equivalence] A context C is a data tree with a hole, we denote by $C[\mathcal{T}]$ the substitution of the hole by \mathcal{T} in C . We say that two models \mathcal{T} and \mathcal{T}' are *strongly S-equivalent* for some $S \subseteq \Delta$ if and only if for every context C such that $\delta(C) \cap S = \emptyset$ and for every BIP automaton M , M accepts $C[\mathcal{T}]$ iff M accepts $C[\mathcal{T}']$.

We now consider a special kind of run of \mathcal{A}_M on a tree for that run, that we call ‘data run’. Together with the run we maintain a witness tree that is a representation of the extended states indicated in the run. The proof then consists in giving an accepting data run α and a witness tree \mathcal{T} for every time the automaton M is non-empty (i.e., every time there is a tree \mathcal{T}' that is accepted by M).

We will make use of the two parameters u_0 and t_0 which are bounds, polynomial in the BIP: $t_0 = 2|K|^2 + 2$ and $u_0 = (2|K|^2 + |K| + 2)|K|$. u_0 is a bound for the maximum branching width of \mathcal{A}_M , and t_0 is the maximum number of data descriptions that D° may have at any extended state.

Definition 6. A *data run* on \mathcal{T} is a labeling function α such that for every $p \in \text{Pos}(\mathcal{T})$, $\alpha(p) = (c_p, cd_p^{\bar{=}}, cd_p^\circ)$ where c_p is an extended state, $cd_p^\circ : [1..t_0] \rightarrow \Delta$ is a partial mapping from the D° elements of c_p to the represented data values in the model \mathcal{T} and $cd_p^{\bar{=}} : [1..|K|] \rightarrow \Delta$ a partial mapping for the $D^{\bar{=}}$. And:

- $\mathcal{T}|_p \triangleright c_p$,
- for all $i, j \in \mathbb{N}$ s.t. $i \neq j$, $pi, pj \in \text{Pos}(\mathcal{T})$, if $d \in \delta(\mathcal{T}|_{pi}) \cap \delta(\mathcal{T}|_{pj})$, then $d \in (Im(cd_{pi}^{\bar{=}}) \cup Im(cd_{pi}^\circ)) \cap (Im(cd_{pj}^{\bar{=}}) \cup Im(cd_{pj}^\circ))$, where Im stands for the *image* of the function.
- for every $i \in \mathbb{N}$, $pi \in \text{Pos}(\mathcal{T})$, if $\delta(p) \in \delta(\mathcal{T}|_{pi})$, then $\delta(p) \in Im(cd_{pi}^{\bar{=}}) \cup Im(cd_{pi}^\circ)$,
- $(Im(cd_p^{\bar{=}}) \cup Im(cd_p^\circ)) \subseteq \bigcup_{pi \in \text{Pos}(\mathcal{T})} (Im(cd_{pi}^{\bar{=}}) \cup Im(cd_{pi}^\circ)) \cup \{\delta(p)\}$
- for each atomic formula $\exists(k_1, k_2)^\circ \in \text{atForm}_M$ such that $v_p(\exists(k_1, k_2)^\circ) = 1$, there exist at most two data values in $Im(cd_p^{\bar{=}}) \cup Im(cd_p^\circ)$ that witness the formula. That is, for example, if $\exists(k_1, k_2)^\circ$, then there are two data values $d, d' \in Im(cd_p^{\bar{=}}) \cup Im(cd_p^\circ)$ and

two runs ρ, ρ' of the pathfinder automaton such that $o(\rho) = (k_1, d)$, $o(\rho') = (k_2, d')$ and $d \neq d'$.

If we also have that \mathcal{T} has its rank bounded by the polynomial used by \mathcal{A}_M ($2|K|^2 + |K| + 2$), it is then easy to see that for every internal node $p \in \text{Pos}(\mathcal{T})$ the extended state corresponds to a transition of \mathcal{A}_M . We say that an extended state c is *incomplete* if $D_\varepsilon^\circ(i) = \emptyset$ for some $i \in [1..t_0]$ (or *complete* otherwise). We say that $(c, cd^{\bar{=}}, cd^\circ)$ is a *d-completion* of $(c', cd'^{\bar{=}}, cd'^\circ)$ for $d \in \Delta$ if either (i) all the components are equal ($c = c', cd^{\bar{=}} = cd'^{\bar{=}}$ and $d \in Im(cd^{\bar{=}}) \cup Im(cd^\circ)$), or (ii) $cd^{\bar{=}} = cd'^{\bar{=}}$, $cd^\circ = cd'^\circ[r \mapsto d]$ for some r such that $cd'^\circ(r)$ is undefined, and c and c' differ only in that $D^\circ(r) \neq \emptyset$ in c' .

We now show that given an accepting run of M for \mathcal{T} , we can find a model \mathcal{T}' that is also accepted by M together with a *data run* on it. Moreover we show that we can assume that M' is ranked. In this way we explicitly present the steps that \mathcal{A}_M needs to perform to reach the desired final extended state. We can see that having a data run on a ranked tree implies having a run in \mathcal{A}_M and then in the emptiness algorithm, as the \equiv_E -relation needed at each step can be deduced from the cd functions of the involved nodes.

Let λ be the run of the automaton M on a model \mathcal{T} . We show that there exists a model \mathcal{T}' and data run α such that for every $p \in \text{Pos}(\mathcal{T}')$, $\alpha(p) = (c_p, cd_p^{\bar{=}}, cd_p^\circ)$ and

1. \mathcal{T}' is S -strongly equivalent to \mathcal{T} for some S ,
2. $\mathcal{C}(v_\varepsilon) = \lambda(\varepsilon)$ for v_ε the valuation of c_ε ,
3. for every p , $\delta(p) \in Im(cd_p^{\bar{=}}) \cup Im(cd_p^\circ)$,
4. $\alpha(\varepsilon)$ is incomplete,
5. \mathcal{T} and \mathcal{T}' have the same height,
6. $S \cap \delta(\mathcal{T}) = \emptyset$.

We proceed by induction on the height of the tree \mathcal{T} . The general strategy is first to apply inductive hypothesis on all the immediate subtrees of the root. Then, we select the data values necessary to validate all the formulae at the root. For example, if in the root of \mathcal{T} the formula $\exists(k_1, k_2)^\circ$ is true, then we select two data values that ‘witness’ this formula in \mathcal{T}' . Equivalently, if for a certain state k all the runs ending in k retrieve the same data value d , then the component $D^{\bar{=}}(\chi(k))$ of the extended state at the root must describe the datum d . In this case, we also select the data value (d) to witness this. For each subtree that contains one of these data values, we make sure that it is described by the extended state by a ‘Completion Lemma’ we will later prove.

The base case is immediate by definition. So let us suppose we are in a model \mathcal{T} such that the root has u immediate descendants $1 \dots u$ for any $u \in \mathbb{N}$. Suppose by inductive hypothesis that we have data runs $\alpha_1 \dots \alpha_u$ for some models $\mathcal{T}_1 \dots \mathcal{T}_u$ such that they are all incomplete in the root and \mathcal{T}_i is X_i -strongly equivalent to $\mathcal{T}|_i$ for some $X_i \subset \Delta$. Wlog we can assume that $\bigcup_i X_i \cap \delta(\mathcal{T}) = \emptyset$ and that $\forall i \neq j : X_i \cap X_j = \emptyset$ (it is an easy exercise using condition 6 to see that otherwise we can do a data transformation via a bijection).

Let $v_\varepsilon : \text{Form} \rightarrow \{0, 1\}$ be the valuation witnessed in the root of \mathcal{T} . In the following construction we use a set G which initially is $\{d\}$, with d the data value of the root. The idea is that depending on the valuation v_ε we collect in G all the data values of \mathcal{T} necessary to witness all existential formulæ and place them all in the D -component of c_0 . Here, using the strong S -equivalence property of $\mathcal{T}_1 \dots \mathcal{T}_u$ wrt $\mathcal{T}|_1 \dots \mathcal{T}|_u$, we know that if $d \in \delta(\mathcal{T}|_i)$ is a witness in $\mathcal{T}|_i$, then $d \in \delta(\mathcal{T}_i)$ and it is also a witness.

If $v_\varepsilon(\exists(k_1, k_2)^\otimes) = 1$, then there are two (or one) data values d_1, d_2 involved in the satisfaction of this formula. We pick any two, and add them to G . For every immediate subtree \mathcal{T}_i and data run that contains d_1 we can apply the Completion Lemma we will see next, extending this model into a $S_i^{d_1}$ -strongly equivalent one $\mathcal{T}_i^{d_1}$ (for some $S_i^{d_1}$) such that there exists a data run α'_i where if $\alpha'_i(\varepsilon) = (c_i, cd_i^{\bar{=}}, cd_i^{\circ})$ then $d_1 \in \text{Im}(cd_i^{\bar{=}}) \cup \text{Im}(cd_i^{\circ})$. We repeat the process for d_2 and for all formulæ $\exists(k_1, k_2)^\otimes$ that hold in the root. We do this for all the existential formulæ whose valuation is 1.

It is easy to see that we can assume that all the S_i^d are disjoint, and that they are also disjoint from $\delta(\mathcal{T})$. This is due to the fact that $S_i^d \cap (\text{Im}(cd_i^{\bar{=}}) \cup \text{Im}(cd_i^{\circ})) = \emptyset$ by the Completion Lemma.

We then build the model \mathcal{T}' as $\langle s, d \rangle(\mathcal{T}_1 \dots \mathcal{T}_z)$ where $\mathcal{T}_1 \dots \mathcal{T}_z$ are all the models \mathcal{T}_i^d obtained before, and $\langle s, d \rangle$ is the root of \mathcal{T} . We can verify then that the height between \mathcal{T} and \mathcal{T}' remains unchanged (condition 5). Wlog we assume that for any of the data values ' d ' that were considered before (i.e., the ones in G), $d \in \delta(\mathcal{T}_j^{d'})$ iff $d = d'$. We can always apply a data bijection to the subtrees to be sure that this is the case.

We define α as follows: $\alpha(\varepsilon) = (c_\varepsilon, cd_\varepsilon^{\bar{=}}, cd_\varepsilon^{\circ})$ where $c_\varepsilon = (v, \langle D^{\bar{=}}, D^{\circ} \rangle)$ and $v, D^{\bar{=}}$ and $cd_\varepsilon^{\bar{=}}$ are the ones inferred from the model \mathcal{T}' (that is, the *only ones* such that $\mathcal{T}' \triangleright c_\varepsilon$). In order to obtain D_ε° and cd_ε° , we use the set G that contains the data values that witness the existential formulæ. If $G = \{d_1, \dots, d_r\}$, define $D_\varepsilon^{\circ}(i) = \text{Reach}(d_i)$ if $i \leq r$ and $d_i \notin \text{Im}(cd_\varepsilon^{\bar{=}})$ and also define $cd_\varepsilon^{\circ}(i) = d_i$. Otherwise, $D_\varepsilon^{\circ}(i) = \emptyset$ and $cd_\varepsilon^{\circ}(i)$ is undefined. It is easy to see that as $|G| \leq 2|K|^2 + 2|K| + 1 < t_0$, $\alpha(\varepsilon)$ is clearly incomplete (condition 4). Finally, define $S = \bigcup_i X_i \cup \bigcup_i S_i$. By this definition of S condition 6 is valid.

We can see that condition 2 holds and that α is a well-defined data run. All existential formulæ $\exists(k_1, k_2)^\otimes$ with true valuation in v_ε at the root of \mathcal{T} are witnessed by the D -components of the immediate subtrees of the root and hence continue to be true. On the other hand, the formulæ such that $\exists(k_1, k_2)^\neq$ with *false* valuation in \mathcal{T} are taken care of by the $D^{\bar{=}}$ -component. Finally, the formulæ $\exists(k_1, k_2)^{\bar{=}}$ with false valuation obviously continue to be false because we are collapsing only data values of the D -component. Note that the difference is that now we are exactly under the hypothesis of the algorithm, where *only* the data values of the D -component are merged, according to the cd functions. We can then check that α is a well defined data run.

We have then showed the existence of a data run and a

witnessing tree for every model accepted by M . In order to finally show that there is a run of algorithm that reaches the extended state of the root of the data run, we must show that all this can be done with a tree with bounded rank. Moreover, with the rank used in the emptiness algorithm $(2|K|^2 + |K| + 2)|K|$. This is shown with a simple argument by the Bounded Rank Lemma. This concludes the proof of correctness of the emptiness algorithm. \square

Lemma 1. [Incompletion Lemma] For any data run α on \mathcal{T} there is another one α' on the same model \mathcal{T} such that the extended state of the root is incomplete.

PROOF. If $\alpha(\varepsilon)$ is complete, it means that all $2|K|^2 + 2$ registers of D° are defined, where we use 'register' to refer to each of the components of the D° tuple. It is easy to see that in order to witness all the existential formulæ that hold in v_ε we need only –at most– $2|K|^2$. If the data value of the root is also preserved, we need not more than $2|K|^2 + 1$ data values. Then there must be an element $D(i)$ that is either the empty set, or that is not necessary to verify all the existential formulæ of v_ε . We can then replace $D' = D[i \mapsto \emptyset]$, $cd' = cd[i \mapsto \perp]$. \square

Lemma 2. [Completion Lemma] For any data run α on \mathcal{T} with an incomplete extended state in the root, and for any $d \in \delta(\mathcal{T})$ there is another run α' on \mathcal{T}' and a set $S \subseteq \Delta$ such that

1. \mathcal{T} and \mathcal{T}' are strongly S -equivalent,
2. \mathcal{T} and \mathcal{T}' have the same height,
3. $\delta(\mathcal{T}) \cap S = \emptyset$, $d \notin S$,
4. $\alpha'(\varepsilon)$ is a d -completion of $\alpha(\varepsilon)$,

PROOF. Suppose that r is such that $D^\circ(r) = \emptyset$ in the root and that α is the data run on \mathcal{T} such that $\alpha(p) = (c_p, cd_p^{\bar{=}}, cd_p^{\circ})$ for every $p \in \text{Pos}(\mathcal{T})$. If d was already in $\text{Im}(cd_\varepsilon^{\bar{=}}) \cup \text{Im}(cd_\varepsilon^{\circ})$, there is nothing to be done. If d is in some $\text{Im}(cd_i^a)$ $a \in \{\bar{=}, \circ\}$, $i \in \mathbb{N}$, then we only need to define $cd_\varepsilon^{\circ}(r) = d$, $D_\varepsilon^{\circ}(r) = \text{Reach}(d)$ and $S = \emptyset$.

Otherwise, d is included in one (and *only one*) subtree i among the u immediate subtrees of the root. The uniqueness is due to the fact that we are working with a data run, under the hypothesis that data values that can be shared among the subtrees are those declared in the cd functions (which is the case we have already ruled out).

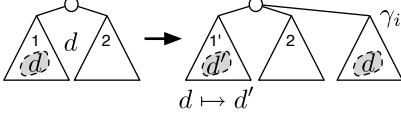
If $\alpha(i)$ is complete, applying the Incompletion Lemma we obtain β_i an incomplete data run on $\mathcal{T}|_i$, and in this way we 'make room' for d . We now apply inductive hypothesis on $\mathcal{T}|_i, \beta_i$ with the datum d . We obtain a data run γ_i , the model \mathcal{T}'_i and a data set S_i . We can assume wlog that S_i is different from $\delta(\mathcal{T})$ and that \mathcal{T}'_i has a data set different from \mathcal{T} except perhaps in the data values of the cd functions of γ_i . So, $\delta(\mathcal{T}'_i) \cap \delta(\mathcal{T}) \subseteq \text{Im}(cd_i^{\bar{=}}) \cup \text{Im}(cd_i^{\circ})$, where cd_i^a are the ones defined in $\gamma_i(\varepsilon)$.

Let d' be a fresh data value not in $\delta(\mathcal{T})$ nor in S_i and define the data transformation $f = id[d \mapsto d']$, and $\mathcal{T}'_i^b = f(\mathcal{T}'_i)$.

Then the final model is,

$$\mathcal{T}' = \langle s_\varepsilon, d_\varepsilon \rangle (\mathcal{T}|_1 \dots \mathcal{T}|_{i-1}, \mathcal{T}_i^b, \mathcal{T}|_{i+1} \dots \mathcal{T}|_u, \mathcal{T}'_i)$$

and $S = S_i \cup \{d'\} \cup (\delta(\mathcal{T}'_i) \setminus G)$, where $G = (Im(cd_\gamma^\ominus) \cup Im(cd_\gamma^\ominus)) \cap \delta(\mathcal{T})$ with the maps cd_γ^\ominus and cd_γ^\ominus defined in $\gamma(\varepsilon)$; and $(s_\varepsilon, d_\varepsilon)$ is the root symbol and datum of \mathcal{T} .



By this definition, S is different from all the data values of \mathcal{T} and from d , and then condition 3 is then satisfied. The data run α' for this model is defined by $\alpha'(\varepsilon) = (c'_\varepsilon, cd_\varepsilon^\ominus, cd_\varepsilon^\ominus[r \mapsto d])$ where c'_ε is equal to c_ε except that D_ε^\ominus has an extra definition $D_\varepsilon^\ominus(r) = Reach(d)$; $\alpha'(ip) = (c_{ip}, f \circ cd_{ip}^\ominus, f \circ cd_{ip}^\ominus)$ where $\alpha(ip) = (c_{ip}, cd_{ip}^\ominus, cd_{ip}^\ominus)$; and the other subtrees $j \in [1..u] \setminus \{i\}$ are preserved, $\alpha'(jp) = \alpha(jp)$. Finally, $\alpha'((u+1)p) = \gamma_i(p)$.

It is immediate that α' is a completion of α (condition 4). It can be verified that by construction α' is a well-defined data run on \mathcal{T}' , and that \mathcal{T} and \mathcal{T}' are equivalent for contexts with no data values in S . This is mainly because when the copy of the tree is made, we are careful enough not to change the data values contained in the cd functions (especially cd^\ominus). It is evident that the height was preserved (condition 2). \square

Lemma 3. (Bounded rank) Given a data run α on \mathcal{T} , there is another one α' on \mathcal{T}' such that

- $\alpha(\varepsilon) = \alpha'(\varepsilon)$,
- the rank of \mathcal{T}' is bounded by $(2|K|^2 + |K| + 2)|K|$,
- \mathcal{T}' is a subtree of \mathcal{T} .

PROOF. In $\alpha(\varepsilon)$ we have the description of at most $2|K|^2 + |K| + 2$ data values, each one with its description in D^\ominus or D^\ominus . Let d be one of the data values of $Im(cd_\varepsilon^\ominus) \cup Im(cd_\varepsilon^\ominus)$, and let $S \subseteq K$ be the description of it. It is easy to see that in order to preserve the description we need to keep at most $|S|$ immediate subtrees, as in the worst case each one of them would contribute with one state of S . Doing this for all data values we obtain that we just need to select at most $(2|K|^2 + |K| + 2)|K|$ immediate subtrees to preserve $D_\varepsilon^\ominus, D_\varepsilon^\ominus$. As all the data values to satisfy v_ε are included in $Im(cd_\varepsilon^\ominus) \cup Im(cd_\varepsilon^\ominus)$ (by definition of a data run), we only need to keep these subtrees. This can be clearly done for any inner node of the tree. As a result \mathcal{T}' has a branching width bounded by $(2|K|^2 + |K| + 2)|K|$. \square

D. PSpace of $XPath(\downarrow^*, =) \setminus \varepsilon$

Proposition 7. $XPath(\downarrow^*, =) \setminus \varepsilon$ has the poly-depth model property.

PROOF. For any formula η and data tree we show that we can describe each element of a given branch using only

a polynomial number of descriptions. As a result, we prove that there is a polynomial bound on the height of the tree. Otherwise, we show show long branches can be 'shortened' preserving satisfiability of η .

For any two positions $p_1, p_2 \in Pos(\mathcal{T})$, let us call $p_1 \prec p_2$ iff p_1 is a prefix of p_2 , which means that p_2 is a descendant of p_1 . The key observation for this proof is that

- If a formula of the type $\downarrow^* p_1 \otimes \downarrow^* p_2$, $\otimes \in \{=, \neq\}$ is true at a certain position $\nu \in Pos(\mathcal{T})$, then it is true at all positions $\mu \prec \nu$. That is, at all nodes that occur before ν (i.e., closer to the root).
- If a formula of the type $\neg(\downarrow^* p_1 \otimes \downarrow^* p_2)$, $\otimes \in \{=, \neq\}$ is true at a certain position $\nu \in Pos(\mathcal{T})$, then it is true at all positions $\mu \in Pos(\mathcal{T})$ such that $\mu \succ \nu$. That is, at all the nodes that are in the subtree generated by ν .

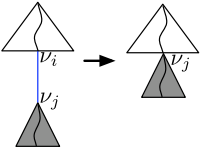
We can see that this fact also extends to unions and then to any path formula of the fragment, as any of the basic bricks of the path start always with \downarrow^* .

Let η be a formula that is satisfied in a model \mathcal{T} . Let $B = \nu_1 \dots \nu_k$ be a maximal branch of this model where ν_1 is the root and ν_k is a leaf. We show that this branch can be bounded by a polynomial on $|\eta|$. We order all the path subformulae of η in a sequence E whose elements we denote by e_i and are of the form p or $p \otimes p'$ (with p, p' path expressions). Each one of these formulae has associated a position in the branch, $pos_B(e) = i$ iff ν_i satisfies e and there is no ν_j with $j < i$ that satisfies it. We assume then that $E = (e_1 \dots e_t)$ is ordered according to this notion of first or last appearance given by pos_B . Here, t is clearly bounded by the number of subformulae of η , bounded by $2|\eta|$. We can then label each node ν of the branch by an index i ($0 \leq i \leq t$) meaning that $\{e_1 \dots e_i\}$ is the set of path subformulae of η that are verified in ν .

On the other hand, for every $e_i \in E$, we must consider the succession of nodes necessary for witnessing this formula in this model. For instance, the formula $\downarrow^* [\varphi] \downarrow^* [\psi]$ may need two nodes to be satisfied. In general, a path expression p that contains n appearances of \downarrow^* needs at most n 'witness' nodes to be satisfied, and as in our case p is a subformula of η we can safely bound n by $|\eta|$. The strategy is to preserve all the nodes necessary to witness e_i in the node $\nu_{pos_B(e_i)}$. The idea is that by preserving all the nodes necessary to satisfy the formula at this node, we are actually making sure that all precedent nodes also satisfy the formula. Some of them will be in the (sub-)branch $\nu_{pos_B(e_i)} \dots \nu_k$, some not. For those that are on the branch we label each element with the tag 'wit : e_i ', for those that are not on the branch, we tag their projection on the branch (i.e., the closest ancestor on B) with 'wit : e_i '. Summing up, for each $1 \leq i \leq t$ we have at most $|\eta|$ elements of B tagged by 'wit : e_i ', and then at most $2|\eta|^2$ number of witnesses of any formula along the branch.

We then have as tag of each element of B a triple of (1) an index i of path subformulae that are satisfied, (2) the symbol of the node $a \in \Sigma$, and (3) perhaps the label $wit : e$ if it is witness of a path formula.

We call that a sub-branch of $B \nu_i \dots \nu_j$ is a ‘safe zone’ iff it is a maximal subbranch such that it does not contain elements with label ‘ $wit : e$ ’ for some e . As a direct consequence of the bound on the number of witness elements, there are at most $2|\eta|^2 + 1$ safe zones in B .



If two elements $\nu_i, \nu_j, i < j$ inside the same safe zone are tagged with the same label, we claim that the branch can be shortened preserving the satisfaction of all the subformulae of η at the nodes $\nu_1 \dots \nu_i$. The operation consists in replacing in \mathcal{T} the subtree of ν_i by that of ν_j .

It is easy to see that the satisfaction of e_i -formulae is preserved as all its witness are maintained. It is easy to see that there will not be new elements that satisfy a path subformula, because we only *erased* nodes from the tree, preserving the descendant relation. This is due to the fact that any path expression that is true in the modified model, was true in the previous one.

We can conclude that we can assume that in a safe zone there are no more than $|\eta|^3$ different tags, and different elements (otherwise we can shorten the branch). Thus, we have that a branch has at most $2|\eta|^2$ witnesses and $2|\eta|^2 + 1$ safe zones of $|\eta|^3$ elements each. Then we can bound the length of any branch by $2|\eta|^2 + (2|\eta|^2 + 1) \cdot |\eta|^3$ and thereby the height of the tree is polynomially bounded. \square

Corollary 2. $XPath(\downarrow^*, =) \setminus \varepsilon$ is PSPACE-complete

PROOF. The membership in PSPACE is an immediate consequence of Theorem 6 and Proposition 7. We can prove that $XPath(\downarrow^*)$ is hard for PSPACE by a reduction from the QBF-validity problem that can be found in Appendix E. \square

E. Complexity of $XPath(\downarrow^*)$

Proposition 8. $XPath(\downarrow^*)$ is hard for PSPACE.

PROOF. The proof consists in coding an instance of the QBF validity problem in satisfiability of $XPath(\downarrow^*)$.

Let $\varphi = Q_1 p_1 \dots Q_n p_n \cdot \psi$ where p_i are propositional variables (pairwise distinct), $Q_i \in \{\forall, \exists\}$ and ψ is a formula of the propositional calculus in CNF.

The idea is to force a model in which every branch contains a full valuation for the variables p_1, \dots, p_n where we force that all valuations present in the tree are contained in those specified by the quantifiers. The alphabet of this tree is $\Sigma = \{p_1, \dots, p_n, \bar{p}_1, \dots, \bar{p}_n, X\}$, and every branch lists a valuation in order, that is, first there is a node with a label in $\{p_1, \bar{p}_1\}$, then another in $\{p_2, \bar{p}_2\}$, etc. The label X simply marks the ending of a valuation in a branch. Although there *could* be more valuations on $p_1 \dots p_n$ after an X , these are redundant, as there has been defined before X as well. After this marking we build the tree that satisfies ψ and finally we check that there are no inconsistencies with respect to its valuation.

Let v_i be the formula that specifies that the node is a valuation for the propositional variable p_i : $v_i = p_i \vee \bar{p}_i$.

- If $Q_1 = \forall$, then $f_1 = \langle \downarrow^* [p_1] \rangle \wedge \langle \downarrow^* [\bar{p}_1] \rangle$.
If $Q_1 = \exists$, then $f_1 = \langle \downarrow^* [p_1] \rangle \vee \langle \downarrow^* [\bar{p}_1] \rangle$.
- ($i > 1$) If $Q_i = \forall$, then $f_i = \neg \langle \downarrow^* [v_{i-1} \wedge \neg(\langle \downarrow^* [p_i] \rangle \wedge \langle \downarrow^* [\bar{p}_i] \rangle)] \rangle$.
If $Q_i = \exists$, then $f_i = \neg \langle \downarrow^* [v_{i-1} \wedge \neg(\langle \downarrow^* [p_i] \rangle \vee \langle \downarrow^* [\bar{p}_i] \rangle)] \rangle$.
- φ_X forces that the label X always appears once the valuation for all propositions has been defined.
 $\varphi_X = \neg \langle \downarrow^* [v_1] \downarrow^* [v_2] \dots \downarrow^* [v_{n-1}] \downarrow^* [v_n \wedge \neg \langle \downarrow^* [X] \rangle] \rangle$
- For all X we build the formula for $\psi = C_1 \wedge \dots \wedge C_t$ where $C_i = t_1 \vee \dots \vee t_{j_i}$ and each t is a valuation for some p_i . That is:

$$\tau = \bigwedge_{C_i} \bigvee_{t \in C_i} \langle \downarrow^* [t] \rangle$$

where t is p_j or \bar{p}_j for some j . And this must hold for all X -valued node:

$$\varphi_\psi = \neg \langle \downarrow^* [X \wedge \neg \tau] \rangle$$

- Finally, we must check that no inconsistencies are to be found between the p_i .

$$\varphi_{inc} = \bigwedge_{i=1}^n \neg \langle \downarrow^* [p_i] \downarrow^* [\bar{p}_i] \rangle \wedge \neg \langle \downarrow^* [\bar{p}_i] \downarrow^* [p_i] \rangle$$

The final formula is then

$$\varphi_F = \bigwedge_{i=1}^n f_i \wedge \varphi_X \wedge \varphi_\psi \wedge \varphi_{inc}$$

Lemma 4. φ is QBF-valid iff φ_F is $XPath(\downarrow^*)$ -satisfiable.

PROOF. The proof is straightforward and is left to the reader. \square

Proposition 9. $XPath(\downarrow^*)$ is in PSPACE

PROOF. For any path formula $p \in XPath(\downarrow^*)$ satisfied in \mathcal{T} , by k witnesses of a branch $\nu_1 \prec \nu_2 \prec \dots \prec \nu_k$. By induction on k using Lemma 5 we can easily show that it is then satisfiable in a model \mathcal{T}' where the witness ν_i is at depth i for all $1 \leq i \leq k$. This can be done for all path formulae that hold at the root.

On the other hand, by Lemma 6 below we can also assume that every internal node is the witness of –at most– *one* path subformula of an ancestor node. The normal form for a model \mathcal{T} consists in:

1. For all nodes $\nu, \nu', \nu'' \in Pos(\mathcal{T})$ and path formulae $p, p' \in XPath(\downarrow^*)$, if $\nu \in wit_{\mathcal{T}}(\nu', p), \nu \in wit_{\mathcal{T}}(\nu'', p')$, then $\nu' = \nu''$.
2. For every node $\nu \in Pos(\mathcal{T})$ and path p , if $wit_{\mathcal{T}}(\nu, p) = \dots \nu_i, \nu_{i+1} \dots$, then $\nu_{i+1} = \nu_i n$ with $n \in \mathbb{N}$.

Let us define the set of subformulae sub as follows

$$\begin{aligned} \text{sub}(\downarrow^* p) &= \{\downarrow^* p\} \cup \text{sub}(p) \\ \text{sub}([\varphi]p) &= \{[\varphi]p\} \cup \text{sub}(p) \cup \text{sub}(\varphi) \\ \text{sub}(\varepsilon p) &= \text{sub}(p) \\ \text{sub}((p_1 \cup p_2)p_3) &= \{(p_1 \cup p_2)p_3\} \cup \text{sub}(p_1) \cdot p_3 \\ &\quad \cup \text{sub}(p_2) \cdot p_3 \cup \text{sub}(p_3) \end{aligned}$$

Given a formula $\varphi \in \text{sub}(\eta)$ satisfied in a node $\nu \in \text{Pos}(\mathcal{T})$, it can either be

- completely satisfied in ν in the case $\llbracket \varphi \rrbracket^{\mathcal{T}} \ni (\nu, \nu)$
- partially satisfied in ν and witnessed in another node $\nu' \succ \nu$ such that ν' satisfies a formula of the *next step* of φ , $ns(\varphi)$.

$$\begin{aligned} ns((p_1 \cup p_2)p_3) &= ns(p_1) \cdot p_3 \cup ns(p_2 \cdot p_3) \cup ns(p_3) \\ ns(p_1 p_2) &= ns(p_1) \cdot p_2 \cup ns(p_2) \\ ns(\varepsilon p) &= ns(p) \\ ns([\psi]p) &= ns(p) \\ ns(\downarrow^* p) &= \{p\} \end{aligned}$$

It is easy to see that $ns(p) \subseteq \text{sub}(p)$. We can then assume that there exists a function $wit_{\mathcal{T}}(\nu) \subseteq \text{Pos}(\mathcal{T}) \times \text{sub}(\eta)$ that describes for each node ν , a finite set of *witnesses*. For example, if $(\downarrow^* [a] \cup \varepsilon) \downarrow^* [b]$ holds in ν , then it could be that $(\nu', [a] \downarrow^* [b]) \in wit_{\mathcal{T}}(\nu)$, or $(\nu', \downarrow^* [b]) \in wit_{\mathcal{T}}(\nu)$ for some node $\nu' \succ \nu$.

By the normal form we have that there will not be a repeated node in $wit_{\mathcal{T}}(\nu)$, and that all nodes will be of the form ν_i , with $i \in \mathbb{N}$.

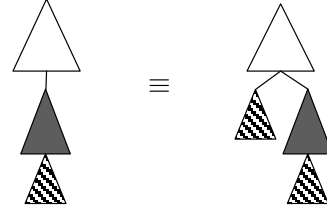
Suppose now that we have a branch $B = \nu_1 \dots \nu_z$ such that there exists a path formula p where $(\nu_i, p) \in wit_{\mathcal{T}}(\nu_{i-1})$, $(\nu_j, p) \in wit_{\mathcal{T}}(\nu_{j-1})$, that is, the path p is witnessed at two nodes $\nu_i \prec \nu_j$ along the branch B . Then we can safely replace $\mathcal{T}|_i$ by $\mathcal{T}|_j$. \square

Lemma 5. (Subtree copy.) $\text{XPath}(\downarrow^*)$ is closed under subtree copy. That is, for every pair of contexts C_1, C_2 , tree \mathcal{T} and $\eta \in \text{XPath}(\downarrow^*)$,

$$C_1[C_2[\mathcal{T}]] \models \eta \text{ iff } C_1[\mathcal{T}, C_2[\mathcal{T}]] \models \eta$$

where $C[\mathcal{T}_1, \mathcal{T}_2]$ is the operation of replacing the hole by a forest, in this case of two trees.

PROOF. It is easy to see that all path formulae that hold at the root of $C_1[C_2[\mathcal{T}]]$, hold also in $C_1[\mathcal{T}, C_2[\mathcal{T}]] \models \eta$ as it is an *extension* of the tree. On the other hand, any path formula that is satisfied by a succession of nodes in a branch in $C_1[\mathcal{T}, C_2[\mathcal{T}]] \models \eta$, can also be found in $C_1[C_2[\mathcal{T}]]$.



In other words, the logic $\text{XPath}(\downarrow^*)$ is closed under *subtree copy*. \square

Corollary 3. $\text{XPath}(\downarrow^*)$ is closed under replication of subtrees.

Lemma 6. (Demand splitting) If η is satisfiable, then it is satisfiable in a model such that all the path subformulae $p \in \text{sub}(\eta)$ are satisfied in incomparable subtrees. More formally, if $p, p' \in \text{sub}(\eta)$ such that $p, p' \models \mathcal{T}$ with $wit_{\mathcal{T}}(p) = \nu_1 \dots$, $wit_{\mathcal{T}}(p') = \nu'_1 \dots$, then $\nu_1 \not\prec \nu'_1$ and $\nu'_1 \not\prec \nu_1$.

PROOF. It can easily be seen by using Corollary 3. \square