

Basic Model Theory of XPath on Data Trees^{*}

Diego Figueira
University of Edinburgh
UK

Santiago Figueira
Universidad de Buenos Aires
and CONICET
Argentina

Carlos Areces
Universidad Nacional de Córdoba
and CONICET
Argentina

ABSTRACT

We investigate model theoretic properties of XPath with data (in)equality tests over the class of data trees, i.e., the class of trees where each node contains a label from a finite alphabet and a data value from an infinite domain.

We provide notions of (bi)simulations for XPath logics containing the **child**, **descendant**, **parent** and **ancestor** axes to navigate the tree. We show that these notions precisely characterize the equivalence relation associated with each logic. We study formula complexity measures consisting of the number of nested axes and nested subformulas in a formula; these notions are akin to the notion of quantifier rank in first-order logic. We show characterization results for fine grained notions of equivalence and (bi)simulation that take into account these complexity measures. We also prove that positive fragments of these logics correspond to the formulas preserved under (non-symmetric) simulations. We show that the logic including the **child** axis is equivalent to the fragment of first-order logic invariant under the corresponding notion of bisimulation. If upward navigation is allowed the characterization fails but a weaker result can still be established. These results hold over the class of possibly infinite data trees and over the class of finite data trees.

Besides their intrinsic theoretical value, we argue that bisimulations are useful tools to prove (non)expressivity results for the logics studied here, and we substantiate this claim with examples.

Categories and Subject Descriptors

F.4.1 [Mathematical Logic]: Model theory; H.2.3 [Languages]: Query Languages; I.7.2 [Document Preparation]: Markup Languages

^{*}This work was partially supported by grant ANPCyT-PICT-2010-688, ANPCyT-PICT-2011-0365, UBACyT 20020110100025 and the FP7-PEOPLE-2011-IRSES Project “Mobility between Europe and Argentina applying Logics to Systems” (MEALS) and the Laboratoire International Associé “INFINIS”.

(c) 2014, Copyright is with the authors. Published in Proc. 17th International Conference on Database Theory (ICDT), March 24-28, 2014, Athens, Greece: ISBN 978-3-89318066-1, on OpenProceedings.org. Distribution of this paper is permitted under the terms of the Creative Commons license CC-by-nc-nd 4.0.

General Terms

Theory, Languages

1. INTRODUCTION

We study the expressive power and model theory of XPath—arguably the most widely used XML query language. Indeed, XPath is implemented in XSLT and XQuery and it is used as a constituent part of many 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 [6].

Core-XPath (term coined in [13]) is the fragment of XPath 1.0 containing the navigational behavior of XPath. It can express properties of the underlying tree structure of the XML document, such as the label (tag name) of a node, but it cannot express conditions on the actual data contained in the attributes. In other words, it only allows to reason about trees over a finite alphabet. Core-XPath has been well studied and its satisfiability problem is known to be decidable even in the presence of DTDs [17, 1]. Moreover, it is known that it is equivalent to FO² (first-order logic with two variables over an appropriate signature on trees) in terms of expressive power [18], and that it is strictly less expressive than PDL with converse over trees [2]. From a database perspective, however, Core-XPath fails to include the single most important construct in a query language: the join. Without the ability to relate nodes based on the actual *data values* of the attributes, the logic’s expressive power is inappropriate for many applications.

The extension of Core-XPath with (in)equality tests between attributes of elements in an XML document is named Core-Data-XPath in [4]. Here, we will call this logic XPath₌. Models of XPath₌ are *data trees* which can be seen as XML documents. A data tree is a tree whose nodes contains a *label* from a finite alphabet and a *data value* from an infinite domain (see Figure 1 for an example). We will relax the condition on finiteness and consider also infinite data trees, although all our results hold also on finite structures.

The main characteristic of XPath₌ is to allow formulas of the form $\langle \alpha = \beta \rangle$, where α, β are *path expressions*, that navigate the tree using *axes*: **descendant**, **child**, **ancestor**, **next-sibling**, etc. and can make tests in intermediate nodes. The formula is true at a node x of a data tree if there are nodes y, z that can be reached by the relations denoted by α, β , respectively, and such that the data value of y is equal to the data value of z .

Recent articles investigate several algorithmic problems

of logics evaluated over data trees. For example, satisfiability and evaluation are discussed in [8, 5]. In particular, all the logics studied in this article have a decidable satisfiability problem [10, 9]; but tools to investigate their *expressive power* are still lacking. There are good reasons for this: in the presence of joins and data values, classical notions such as Ehrenfeucht-Fraïssé games or structural bisimulations are difficult to handle. In this article we take the first steps towards understanding the expressive power and model theory of $\text{XPath}_=$ on data trees.

Contribution: $\text{XPath}_=$ can navigate the data tree by means of its axes: **child** (that we will note \downarrow), **descendant** (\downarrow_*), **parent** (\uparrow), **ancestor** (\uparrow_*), etc. $\text{XPath}_=$ can also navigate the data tree horizontally, by going to a next or previous **sibling** of the current node. However, we focus on the vertical axes that allow downward and upward exploration. In particular, we will discuss the following languages: $\text{XPath}_=^\downarrow$ ($\text{XPath}_=$ with \downarrow); $\text{XPath}_=^{\downarrow\uparrow}$ ($\text{XPath}_=$ with \downarrow and \uparrow); $\text{XPath}_=^{\downarrow\downarrow_*}$ ($\text{XPath}_=$ with \downarrow and \downarrow_*); $\text{XPath}_=^{\downarrow\uparrow\downarrow_*}$ ($\text{XPath}_=$ with \downarrow , \uparrow , \downarrow_* and \uparrow_*); and its positive fragments. Our main contributions can be summarized as follows:

- In §3 and §5 we introduce bisimulation notions for $\text{XPath}_=^\downarrow$, $\text{XPath}_=^{\downarrow\uparrow}$, $\text{XPath}_=^{\downarrow\downarrow_*}$, and $\text{XPath}_=^{\downarrow\uparrow\downarrow_*}$ and show that they precisely characterize the logical equivalence relation of the respective logic. We also consider fine grained versions of these bisimulations that take into account the number of nested axes and subformulas. The notion of bisimulation for $\text{XPath}_=^\downarrow$ relies on a strong normal form which we also introduce.
- In §4 we show that the simulations associated to the defined bisimulations characterize the positive fragments of the logics: a formula is equivalent to a positive formula if and only if it is invariant under simulations.
- In §6 we characterize $\text{XPath}_=^\downarrow$ as the fragment of first-order logic over data trees (over a signature that includes the **child** relation and an equivalence relation) that is invariant under bisimulations. If we consider $\text{XPath}_=^{\downarrow\downarrow_*}$ instead the characterization fails, but a weaker result can still be established.
- Using bisimulations we show (non)expressivity results about $\text{XPath}_=$ in §7. We characterize, for example, in which cases increasing the nesting depth increases the expressive power of $\text{XPath}_=^\downarrow$.
- All results are proved both over the class of arbitrary (possibly infinite) data trees, and over the class of finite data trees.

Related work: The notion of bisimulation was introduced independently by Van Benthem [26] in the context of modal correspondence theory, Milner [19] and Park [23] in concurrency theory, and Forti and Honsell [11] in non-wellfounded set theory (see [25] for a historical outlook). This classical work defines a *standard notion of bisimulation* but this notion has to be suitably adapted for a particular, given logic. The notion of bisimulation for a given logic \mathcal{L} defines when two models are indistinguishable for \mathcal{L} , that is, when there is no formula of \mathcal{L} that is true in one model but false in the other. Bisimulations can also be used to obtain model theoretic characterizations that identifies the expressive power of a logic \mathcal{L}_1 in terms of the bisimulation invariant fragment of a logic \mathcal{L}_2 which, hopefully, is better understood. The challenge, here, is to pinpoint both the appropriate notion

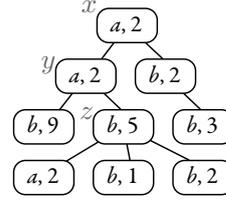


Figure 1: A data tree of $\text{Trees}(\mathbb{A} \times \mathbb{D})$ with $\mathbb{A} = \{a, b\}$ and $\mathbb{D} = \mathbb{N}$.

of bisimulation required and the adequate ‘framework’ logic \mathcal{L}_2 . The classical example of a result of this kind is Van Benthem’s characterization for the basic modal logic as the bisimulation (with the standard notion of bisimulation) invariant fragment of first-order logic [26]. Van Benthem’s original result over arbitrary structures was proved to hold for finite structures by Rosen [24]. The proof was then simplified and unified by Otto [20, 22], and later expanded by Dawar and Otto [7] to other classes of structures.

Logics for semi-structured databases can often be seen as modal logics. In fact, structural characterizations for XPath without equality test were studied in [14], and XPath is known to be captured by PDL [15], whose bisimulation is well-understood [3]. It is then natural to look for an intuitive bisimulation definition for $\text{XPath}_=$.

2. PRELIMINARIES

2.1 Notation

Let $\mathbb{N} = \{1, 2, 3, \dots\}$ and let $[n] := \{1, \dots, n\}$ for $n \in \mathbb{N}$. We use the symbol \mathbb{A} to denote a finite alphabet, and \mathbb{D} to denote an infinite domain (e.g., \mathbb{N}) of **data values**. In our examples we will consider $\mathbb{D} = \mathbb{N}$. We write $X \sim Y$ to say that X is the result of replacing every data value $d \in \mathbb{D}$ from Y by $f(d)$ where $f : \mathbb{D} \rightarrow \mathbb{D}$ is some arbitrary bijection, for any objects X, Y . We write λ for the empty string.

2.2 Data trees

Let $\text{Trees}(A)$ be the set of ordered and unranked trees over an arbitrary alphabet A . We say that \mathcal{T} is a **data tree** if it is a tree from $\text{Trees}(\mathbb{A} \times \mathbb{D})$ where \mathbb{A} is a finite set of **labels** and \mathbb{D} is an infinite set of **data values**. Figure 1 shows an example of a (finite) data tree. A data tree is **finitely branching** if every node has finitely many children. For any given data tree \mathcal{T} , we denote by T its set of nodes. We use letters x, y, z, v, w as variables for nodes. Given a node $x \in T$ of \mathcal{T} , we write $\text{label}(x) \in \mathbb{A}$ to denote the node’s label, and $\text{data}(x) \in \mathbb{D}$ to denote the node’s data value.

Given two nodes $x, y \in T$ we write $x \rightarrow y$ if y is a child of x , and $x \xrightarrow{n} y$ if y is a descendant of x at distance n . In particular, $\xrightarrow{1}$ is the same as \rightarrow , and $\xrightarrow{0}$ is the identity relation. $(x \xrightarrow{n})$ denotes the set of all descendants of x at distance n , and $(\xrightarrow{n} y)$ denotes the sole ancestor of y at distance n (assuming it has one).

For any binary relation R over elements of data trees, we say that a property P is **R -invariant** whenever the following condition holds: for every data tree \mathcal{T} and $u \in T$, if (\mathcal{T}, u) satisfies P and (\mathcal{T}, u) is R -related to (\mathcal{T}', u') then (\mathcal{T}', u') satisfies P .

2.3 XPath

$\llbracket \downarrow \rrbracket^{\mathcal{T}} = \{(x, y) \mid x \rightarrow y\}$	$\llbracket \downarrow^* \rrbracket^{\mathcal{T}} = \text{reflexive transitive closure of } \llbracket \downarrow \rrbracket^{\mathcal{T}}$
$\llbracket \uparrow \rrbracket^{\mathcal{T}} = \{(x, y) \mid y \rightarrow x\}$	$\llbracket \uparrow^* \rrbracket^{\mathcal{T}} = \text{reflexive transitive closure of } \llbracket \uparrow \rrbracket^{\mathcal{T}}$
$\llbracket \varepsilon \rrbracket^{\mathcal{T}} = \{(x, x) \mid x \in T\}$	$\llbracket a \rrbracket^{\mathcal{T}} = \{x \in T \mid \text{label}(x) = a\}$
$\llbracket [\varphi] \rrbracket^{\mathcal{T}} = \{(x, x) \mid x \in \llbracket [\varphi] \rrbracket^{\mathcal{T}}\}$	$\llbracket [\alpha\beta] \rrbracket^{\mathcal{T}} = \{(x, z) \mid (\exists y \in T) (x, y) \in \llbracket [\alpha] \rrbracket^{\mathcal{T}}, (y, z) \in \llbracket [\beta] \rrbracket^{\mathcal{T}}\}$
$\llbracket \neg\varphi \rrbracket^{\mathcal{T}} = T \setminus \llbracket [\varphi] \rrbracket^{\mathcal{T}}$	$\llbracket \langle \alpha \rangle \rrbracket^{\mathcal{T}} = \{x \in T \mid (\exists y \in T) (x, y) \in \llbracket [\alpha] \rrbracket^{\mathcal{T}}\}$
$\llbracket [\alpha \cup \beta] \rrbracket^{\mathcal{T}} = \llbracket [\alpha] \rrbracket^{\mathcal{T}} \cup \llbracket [\beta] \rrbracket^{\mathcal{T}}$	$\llbracket \langle \alpha = \beta \rangle \rrbracket^{\mathcal{T}} = \{x \in T \mid (\exists y, z \in T) (x, y) \in \llbracket [\alpha] \rrbracket^{\mathcal{T}}, (x, z) \in \llbracket [\beta] \rrbracket^{\mathcal{T}}, \text{data}(y) = \text{data}(z)\}$
$\llbracket [\varphi \wedge \psi] \rrbracket^{\mathcal{T}} = \llbracket [\varphi] \rrbracket^{\mathcal{T}} \cap \llbracket [\psi] \rrbracket^{\mathcal{T}}$	$\llbracket \langle \alpha \neq \beta \rangle \rrbracket^{\mathcal{T}} = \{x \in T \mid (\exists y, z \in T) (x, y) \in \llbracket [\alpha] \rrbracket^{\mathcal{T}}, (x, z) \in \llbracket [\beta] \rrbracket^{\mathcal{T}}, \text{data}(y) \neq \text{data}(z)\}$

Table 1: Semantics of XPath₌ for a data tree \mathcal{T} .

We introduce the query language XPath adapted to data trees as abstractions of XML documents. We work with a simplification of XPath, stripped of its syntactic sugar. We consider fragments of XPath that correspond to the navigational part of XPath 1.0 with data equality and inequality. XPath₌ is a two-sorted language, with **path expressions** (that we write α, β, γ) and **node expressions** (that we write φ, ψ, η). The fragment XPath₌(\mathcal{O}), with $\mathcal{O} \subseteq \{\downarrow, \downarrow^*, \uparrow, \uparrow^*\}$, is defined by mutual recursion as follows:

$$\begin{aligned} \alpha, \beta &::= o \mid [\varphi] \mid \alpha\beta \mid \alpha \cup \beta & o &\in \mathcal{O} \cup \{\varepsilon\} \\ \varphi, \psi &::= a \mid \neg\varphi \mid \varphi \wedge \psi \mid \varphi \vee \psi \mid \langle \alpha \rangle \mid \\ &\langle \alpha = \beta \rangle \mid \langle \alpha \neq \beta \rangle & a &\in \mathbb{A} \end{aligned}$$

A **formula** of XPath₌(\mathcal{O}) is either a node expression or a path expression. To save space, we use XPath₌[↓] for XPath₌(\downarrow); XPath₌[↑] for XPath₌(\uparrow); XPath₌^{↓*} for XPath₌(\downarrow, \downarrow^*); and XPath₌^{↑*} for XPath₌(\uparrow, \uparrow^*).

We formally define the semantics of XPath₌ in Table 1. As an example, if \mathcal{T} is the data tree shown in Figure 1, then $\llbracket \langle \downarrow^* [b \wedge \downarrow [b] \neq \downarrow [b]] \rangle \rrbracket^{\mathcal{T}} = \{x, y, z\}$, where the formula reads: “there is a descendant node labeled b , with two children labeled b with different data values.” For a data tree \mathcal{T} and $u \in T$, we write $\mathcal{T}, u \models \varphi$ to denote $u \in \llbracket [\varphi] \rrbracket^{\mathcal{T}}$, and we say that \mathcal{T}, u **satisfies** φ . We say that the formulas φ, ψ of XPath₌ are **equivalent** (notation: $\varphi \equiv \psi$) iff $\llbracket [\varphi] \rrbracket^{\mathcal{T}} = \llbracket [\psi] \rrbracket^{\mathcal{T}}$ for all data trees \mathcal{T} . Similarly, path expressions α, β of XPath₌ are **equivalent** (notation: $\alpha \equiv \beta$) iff $\llbracket [\alpha] \rrbracket^{\mathcal{T}} = \llbracket [\beta] \rrbracket^{\mathcal{T}}$ for all data trees \mathcal{T} .

We call **downward XPath** to XPath₌[↓] and **vertical XPath** to XPath₌[↑].

In terms of expressive power, it is easy to see that \cup is unessential: every XPath₌ node expression φ has an equivalent φ' with no \cup in its path expressions. φ' can be computed in exponential time without incrementing the number of nested axes or the number of nested subformulas. It is enough to use the following equivalences to eliminate occurrences of \cup

$$\begin{aligned} \langle \alpha \circ \beta \rangle &\equiv \langle \beta \circ \alpha \rangle \\ \langle \beta(\alpha \cup \alpha')\beta' \rangle &\equiv \langle \beta\alpha\beta' \rangle \vee \langle \beta\alpha'\beta' \rangle \\ \langle \gamma \circ \beta(\alpha \cup \alpha')\beta' \rangle &\equiv \langle \gamma \circ \beta\alpha\beta' \rangle \vee \langle \gamma \circ \beta\alpha'\beta' \rangle \end{aligned}$$

where $\circ \in \{=, \neq\}$. We will henceforth assume that formulas do not contain union of path expressions.

3. BISIMULATION

3.1 Downward XPath

We write $\text{dd}(\varphi)$ to denote the **downward depth** of φ , defined in Table 2. Let $\ell\text{-XPath}_{\perp}^{\downarrow}$ be the fragment of XPath₌[↓] consisting of all formulas φ with $\text{dd}(\varphi) \leq \ell$.

Let \mathcal{T} and \mathcal{T}' be data trees, and let $u \in T, u' \in T'$. We say that \mathcal{T}, u and \mathcal{T}', u' are **equivalent for XPath₌[↓]** (notation: $\mathcal{T}, u \equiv_{\perp}^{\downarrow} \mathcal{T}', u'$) iff for all formulas $\varphi \in \text{XPath}_{\perp}^{\downarrow}$, we have $\mathcal{T}, u \models \varphi$ iff $\mathcal{T}', u' \models \varphi$. We say that \mathcal{T}, u and \mathcal{T}', u' are **ℓ -equivalent for XPath₌[↓]** (notation: $\mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'$) iff for all $\varphi \in \ell\text{-XPath}_{\perp}^{\downarrow}$, we have $\mathcal{T}, u \models \varphi$ iff $\mathcal{T}', u' \models \varphi$.

For every ℓ , there are finitely many different formulas φ of $\text{dd}(\varphi) \leq \ell$ up to logical equivalence.

PROPOSITION 3.1. $\equiv_{\ell}^{\downarrow}$ has finite index.

COROLLARY 3.2. $\{\mathcal{T}', u' \mid \mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'\}$ is definable by an $\ell\text{-XPath}_{\perp}^{\downarrow}$ -formula $\chi_{\ell, \mathcal{T}, u}$.

3.1.1 Bisimulation and ℓ -bisimulation

Let \mathcal{T} and \mathcal{T}' be two data-trees. We say that $u \in T$ and $u' \in T'$ are **bisimilar for XPath₌[↓]** (notation: $\mathcal{T}, u \leftrightarrow^{\downarrow} \mathcal{T}', u'$) iff there is a relation $Z \subseteq T \times T'$ such that uZu' and for all $x \in T$ and $x' \in T'$ we have

- **Harmony:** If xZx' then $\text{label}(x) = \text{label}(x')$.
- **Zig (Figure 2):** If $xZx', x \xrightarrow{n} v$ and $x \xrightarrow{m} w$ then there are $v', w' \in T'$ such that $x' \xrightarrow{n} v', x' \xrightarrow{m} w'$ and
 1. $\text{data}(v) = \text{data}(w) \Leftrightarrow \text{data}(v') = \text{data}(w')$,
 2. $(\overset{i}{\rightarrow} v)Z(\overset{i}{\rightarrow} v')$ for all $0 \leq i < n$, and
 3. $(\overset{i}{\rightarrow} w)Z(\overset{i}{\rightarrow} w')$ for all $0 \leq i < m$.
- **Zag:** If $xZx', x' \xrightarrow{n} v'$ and $x' \xrightarrow{m} w'$ then there are $v, w \in T$ such that $x \xrightarrow{n} v, x \xrightarrow{m} w$ and items 1, 2 and 3 above are verified.

For a data tree \mathcal{T} and $u \in T$, let $\mathcal{T}|u$ denote the subtree of \mathcal{T} induced by $\{v \in T \mid (\exists n) u \xrightarrow{n} v\}$. Observe that the root of $\mathcal{T}|u$ is u . The following results are straightforward consequences of the definition of bisimulation:

PROPOSITION 3.3. $\mathcal{T}, u \leftrightarrow^{\downarrow} (\mathcal{T}|u), u$.

PROPOSITION 3.4. If \mathcal{T} is a subtree of \mathcal{T}' and $u \in T$ then $\mathcal{T}, u \leftrightarrow^{\downarrow} \mathcal{T}', u$.

We say that $u \in T$ and $u' \in T'$ are **ℓ -bisimilar for XPath₌[↓]** (notation: $\mathcal{T}, u \leftrightarrow_{\ell}^{\downarrow} \mathcal{T}', u'$) if there is a family of relations $(Z_j)_{j \leq \ell}$ in $T \times T'$ such that $uZ_{\ell}u'$ and for all $j \leq \ell$, $x \in T$ and $x' \in T'$ we have

$dd(a) = 0$	$vd(a) = (0, 0)$	$nd(a) = 0$
$dd(\varphi \wedge \psi) = \max\{dd(\varphi), dd(\psi)\}$	$vd(\varphi \wedge \psi) = \max\{vd(\varphi), vd(\psi)\}$	$nd(\varphi \wedge \psi) = \max\{nd(\varphi), nd(\psi)\}$
$dd(\neg\varphi) = dd(\varphi)$	$vd(\neg\varphi) = vd(\varphi)$	$nd(\neg\varphi) = nd(\varphi)$
$dd(\langle\alpha\rangle) = dd(\alpha)$	$vd(\langle\alpha\rangle) = vd(\alpha)$	$nd(\langle\alpha\rangle) = nd(\alpha)$
$dd(\langle\alpha \odot \beta\rangle) = \max\{dd(\alpha), dd(\beta)\}$	$vd(\langle\alpha \odot \beta\rangle) = \max\{vd(\alpha), vd(\beta)\}$	$nd(\langle\alpha \odot \beta\rangle) = \max\{nd(\alpha), nd(\beta)\}$
$dd(\lambda) = 0$	$vd(\lambda) = (0, 0)$	$nd(\alpha\beta) = \max\{nd(\alpha), nd(\beta)\}$
$dd(\varepsilon\alpha) = dd(\alpha)$	$vd(\varepsilon\alpha) = vd(\alpha)$	$nd(\varepsilon) = 0$
$dd([\varphi]\alpha) = \max\{dd(\varphi), dd(\alpha)\}$	$vd([\varphi]\alpha) = \max\{vd(\varphi), vd(\alpha)\}$	$nd([\varphi]) = 1 + nd(\varphi)$
$dd(\downarrow\alpha) = 1 + dd(\alpha)$	$vd(\downarrow\alpha) = \max\{(0, 0), vd(\alpha) + (1, -1)\}$	$nd(\downarrow) = 0$
	$vd(\uparrow\alpha) = \max\{(0, 0), vd(\alpha) + (-1, 1)\}$	$nd(\uparrow) = 0$
Downward depth	Vertical depth	Nesting depth

Table 2: Definitions of downward depth, vertical depth and nesting depth. ($a \in \mathbb{A}$, $\odot \in \{=, \neq\}$, ‘+’ and ‘max’ are performed component-wise, α is any path expression or the empty string λ .)

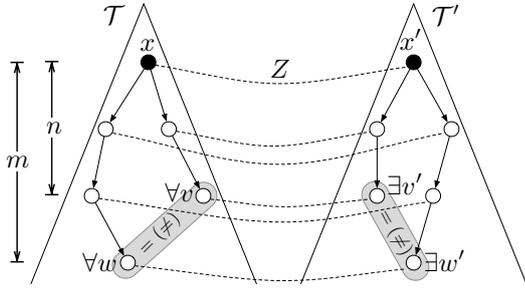


Figure 2: Zig clause of bisimulation for $\text{XPath}_{\downarrow}^{\downarrow}$.

- **Harmony:** If xZ_jx' then $label(x) = label(x')$.
- **Zig:** If xZ_jx' , $x \xrightarrow{n} v$ and $x \xrightarrow{m} w$ with $n, m \leq j$ then there are $v', w' \in T'$ such that $x' \xrightarrow{n} v'$, $x' \xrightarrow{m} w'$ and
 1. $data(v) = data(w) \Leftrightarrow data(v') = data(w')$,
 2. $(\xrightarrow{i}v)Z_{j-n+i}(\xrightarrow{i}v')$ for all $0 \leq i < n$, and
 3. $(\xrightarrow{i}w)Z_{j-m+i}(\xrightarrow{i}w')$ for all $0 \leq i < m$.
- **Zag:** If xZ_jx' , $x' \xrightarrow{n} v'$ and $x' \xrightarrow{m} w'$ with $n, m \leq j$ then there are $v, w \in T$ such that $x \xrightarrow{n} v$, $x \xrightarrow{m} w$ and items 1, 2 and 3 above are verified.

Clearly if $\mathcal{T}, u \Leftrightarrow_{\ell}^{\downarrow} \mathcal{T}', u'$ then $\mathcal{T}, u \Leftrightarrow_{\ell}^{\downarrow} \mathcal{T}', u'$ for all ℓ . For a data tree \mathcal{T} and $u \in T$, let $\mathcal{T}|_{\ell}u$ denote the subtree of \mathcal{T} induced by $\{v \in T \mid (\exists n \leq \ell) u \xrightarrow{n} v\}$.

3.1.2 Equivalence and bisimulation

We now show that $\Leftrightarrow_{\ell}^{\downarrow}$ coincides with $\equiv_{\ell}^{\downarrow}$ on finitely branching data trees, and that $\Leftrightarrow_{\ell}^{\downarrow}$ coincides with $\equiv_{\ell}^{\downarrow}$.

THEOREM 3.5.

1. $\mathcal{T}, u \Leftrightarrow_{\ell}^{\downarrow} \mathcal{T}', u'$ implies $\mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'$. The converse also holds when \mathcal{T} and \mathcal{T}' are finitely branching.
2. $\mathcal{T}, u \Leftrightarrow_{\ell}^{\downarrow} \mathcal{T}', u'$ iff $\mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'$.

The Theorem above is a consequence of the next two propositions:

PROPOSITION 3.6. $\mathcal{T}, u \Leftrightarrow_{\ell}^{\downarrow} \mathcal{T}', u'$ implies $\mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'$.

PROOF. We actually show that if $\mathcal{T}, u \Leftrightarrow_{\ell}^{\downarrow} \mathcal{T}', u'$ via $(Z_i)_{i \leq \ell}$ then for all $0 \leq n \leq j \leq \ell$, for all φ with $dd(\varphi) \leq j$, and for all α with $dd(\alpha) \leq j$:

1. If xZ_jx' then $\mathcal{T}, x \models \varphi$ iff $\mathcal{T}', x' \models \varphi$;
2. If $x \xrightarrow{n} v$, $x' \xrightarrow{n} v'$ and $(\xrightarrow{i}v)Z_{j-n+i}(\xrightarrow{i}v')$ for all $0 \leq i \leq n$, then $(x, v) \in \llbracket \alpha \rrbracket^{\mathcal{T}}$ iff $(x', v') \in \llbracket \alpha \rrbracket^{\mathcal{T}'}$.

We show 1 and 2 by induction on $|\varphi| + |\alpha|$.

Let us see item 1. The base case is $\varphi = a$ for some $a \in \mathbb{A}$. By Harmony, $label(x) = label(x')$ and then $\mathcal{T}, x \models \varphi$ iff $\mathcal{T}', x' \models \varphi$. The Boolean cases for φ are straightforward.

Suppose $\varphi = \langle\alpha = \beta\rangle$. We show $\mathcal{T}, x \models \varphi \Rightarrow \mathcal{T}', x' \models \varphi$, so assume $\mathcal{T}, x \models \varphi$. Suppose there are $v, w \in T$ and $n, m \leq j$ such that $x \xrightarrow{n} v$, $x \xrightarrow{m} w$, $(x, v) \in \llbracket \alpha \rrbracket^{\mathcal{T}}$, $(x, w) \in \llbracket \beta \rrbracket^{\mathcal{T}}$ and $data(v) = data(w)$. By Zig, there are $v', w' \in T'$ such that $x' \xrightarrow{n} v'$, $x' \xrightarrow{m} w'$, $(\xrightarrow{i}v)Z_{j-n+i}(\xrightarrow{i}v')$ for all $0 \leq i \leq n$, $(\xrightarrow{i}w)Z_{j-m+i}(\xrightarrow{i}w')$ for all $0 \leq i \leq m$, and $data(v') = data(w')$. By inductive hypothesis 2 (twice), $(x', v') \in \llbracket \alpha \rrbracket^{\mathcal{T}'}$ and $(x', w') \in \llbracket \beta \rrbracket^{\mathcal{T}'}$. Hence $\mathcal{T}', x' \models \varphi$. The implication $\mathcal{T}', x' \models \varphi \Rightarrow \mathcal{T}, x \models \varphi$ is analogous. The case $\varphi = \langle\alpha \neq \beta\rangle$ is shown similarly. The case $\varphi = \langle\alpha\rangle$ is similar (and simpler) to the previous case.

Let us now analyze item 2. We only show the ‘only if’ direction. The base case is when $\alpha \in \{\varepsilon, \downarrow\}$. If $\alpha = \varepsilon$ then $v = x$ and so $n = 0$. Since $v' = x'$, we conclude $(x', v') \in \llbracket \alpha \rrbracket^{\mathcal{T}'}$. If $\alpha = \downarrow$ then $x \rightarrow v$ in \mathcal{T} , and so $n = 1$. Since $x' \rightarrow v'$, we have $(x', v') \in \llbracket \alpha \rrbracket^{\mathcal{T}'}$.

For the inductive step, let

$$x_0, \dots, x_n \in T \quad \text{and} \quad x'_0, \dots, x'_n \in T'$$

be such that

$$\begin{aligned} x &= x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_n = v && \text{in } \mathcal{T}, \\ x' &= x'_0 \rightarrow x'_1 \rightarrow x'_2 \rightarrow \dots \rightarrow x'_n = v' && \text{in } \mathcal{T}', \end{aligned}$$

and $x_i Z_{j-i} x'_i$ for all $0 \leq i \leq n$. Assume, for contradiction, that $(x', v') \notin \llbracket \alpha \rrbracket^{\mathcal{T}'}$. Then, there is a subformula φ of α and $k \in \{0, \dots, n\}$ such that $\mathcal{T}, x_k \models \varphi$ and $\mathcal{T}', x'_k \not\models \varphi$. This contradicts the inductive hypothesis 1. \square

PROPOSITION 3.7. $\mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'$ implies $\mathcal{T}, u \leftrightarrow_{\ell}^{\downarrow} \mathcal{T}', u'$.

PROOF. Fix $u \in T$ and $u' \in T'$ such that $\mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'$. Define $(Z_i)_{i \leq \ell}$ by

$$xZ_i x' \text{ iff } \mathcal{T}, x \equiv_i^{\downarrow} \mathcal{T}', x'.$$

We show that Z is an ℓ -bisimulation between \mathcal{T}, u and \mathcal{T}', u' . By hypothesis, $uZ_{\ell} u'$. Fix $h \leq \ell$, by construction, Z_h satisfies Harmony. Let us see that Z_h satisfies Zig (the case for Zag is analogous). Suppose $xZ_h x'$,

$$\begin{aligned} x &= v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_n = v & \text{in } \mathcal{T}, \\ x &= w_0 \rightarrow w_1 \rightarrow \dots \rightarrow w_m = w & \text{in } \mathcal{T}', \end{aligned}$$

and $data(v) = data(w)$ (the case $data(v) \neq data(w)$ is shown in a similar way), where $m, n \leq h$. Let $P \subseteq T'^2$ be defined by

$$P = \{(v', w') \mid x' \xrightarrow{n} v' \wedge x' \xrightarrow{m} w' \wedge data(v') = data(w')\}.$$

Since $\mathcal{T}, x \equiv_h^{\downarrow} \mathcal{T}', x'$, $dd(\langle \downarrow^n = \downarrow^m \rangle) \leq h$ and $\mathcal{T}, x \models \langle \downarrow^n = \downarrow^m \rangle$, we conclude that $P \neq \emptyset$. We next show that there exists $(v', w') \in P$ such that

- i. $x' = v'_0 \rightarrow v'_1 \rightarrow \dots \rightarrow v'_n = v'$ in \mathcal{T}' ,
- ii. $x' = w'_0 \rightarrow w'_1 \rightarrow \dots \rightarrow w'_m = w'$ in \mathcal{T}' ,
- iii. $(\forall i \in \{0, \dots, n\}) \mathcal{T}, v_i \equiv_{h-i}^{\downarrow} \mathcal{T}', v'_i$, and
- iv. $(\forall j \in \{0, \dots, m\}) \mathcal{T}, w_j \equiv_{h-j}^{\downarrow} \mathcal{T}', w'_j$,

and hence Zig is satisfied by Z_h . By way of contradiction, assume that for all $(v', w') \in P$ satisfying i and ii we have either

- (a) $(\exists i \in \{0, \dots, n\}) \mathcal{T}, v_i \not\equiv_{h-i}^{\downarrow} \mathcal{T}', v'_i$, or
- (b) $(\exists j \in \{0, \dots, m\}) \mathcal{T}, w_j \not\equiv_{h-j}^{\downarrow} \mathcal{T}', w'_j$.

Fix \top as any tautology such that $dd(\top) = 0$. For each $(v', w') \in P$ we define two families of formulas,

$$\varphi_{v', w'}^0, \dots, \varphi_{v', w'}^n \text{ and } \psi_{v', w'}^0, \dots, \psi_{v', w'}^m,$$

satisfying that $dd(\varphi_{v', w'}^i) \leq h - i$ for all $i \in \{0, \dots, n\}$ and $dd(\psi_{v', w'}^j) \leq h - j$ for all $j \in \{0, \dots, m\}$ as follows:

- Suppose that (a) holds and that i is the smallest number such that $\mathcal{T}, v_i \not\equiv_{h-i}^{\downarrow} \mathcal{T}', v'_i$. Let $\varphi_{v', w'}^i$ be such that $dd(\varphi_{v', w'}^i) \leq h - i$ and $\mathcal{T}, v_i \models \varphi_{v', w'}^i$ but $\mathcal{T}', v'_i \not\models \varphi_{v', w'}^i$. For $k \in \{0, \dots, n\} \setminus \{i\}$, let $\varphi_{v', w'}^k = \top$, and for $k \in \{0, \dots, m\}$, let $\psi_{v', w'}^k = \top$.
- Suppose that (a) does not hold. Then (b) holds. Let j be the smallest number such that $\mathcal{T}, w_j \not\equiv_{h-j}^{\downarrow} \mathcal{T}', w'_j$. Let $\psi_{v', w'}^j$ be such that $dd(\psi_{v', w'}^j) \leq h - j$ and $\mathcal{T}, w_j \models \psi_{v', w'}^j$ but $\mathcal{T}', w'_j \not\models \psi_{v', w'}^j$. For $k \in \{0, \dots, m\} \setminus \{j\}$, let $\psi_{v', w'}^k = \top$, and for $k \in \{0, \dots, n\}$, let $\varphi_{v', w'}^k = \top$.

For each $i \in \{0, \dots, n\}$ and $j \in \{0, \dots, m\}$, let

$$\Phi^i = \bigwedge_{(v', w') \in P} \varphi_{v', w'}^i \text{ and } \Psi^j = \bigwedge_{(v', w') \in P} \psi_{v', w'}^j. \quad (1)$$

Since $dd(\varphi_{v', w'}^i) \leq h - i$, by Proposition 3.1, there are finitely many non-equivalent formulas $\varphi_{v', w'}^i$; the same applies to $\psi_{v', w'}^j$. Hence, both infinite conjunctions in (1) are

equivalent to finite ones, and we may assume that Φ^i and Ψ^j are well-formed formulas. Finally, let

$$\alpha = [\Phi^0] \downarrow [\Phi^1] \downarrow \dots \downarrow [\Phi^n] \text{ and } \beta = [\Psi^0] \downarrow [\Psi^1] \downarrow \dots \downarrow [\Psi^m].$$

By construction, $dd(\alpha), dd(\beta) \leq h$ and so $dd(\langle \alpha = \beta \rangle) \leq h$. Furthermore, $\mathcal{T}, x \models \langle \alpha = \beta \rangle$ and $\mathcal{T}', x' \not\models \langle \alpha = \beta \rangle$. This contradicts $\mathcal{T}, x \equiv_h^{\downarrow} \mathcal{T}', x'$. \square

3.2 Vertical XPath

We now study bisimulation for $XPath_{\downarrow}^{\downarrow}$. Interestingly, the notion we give is simpler than the one for $XPath_{\downarrow}$ due to a normal form enjoyed by the logic.

In the downward fragment of $XPath_{\downarrow}$ we used $dd(\varphi)$ to measure the maximum depth from the current point of evaluation that the formula can access. For the vertical fragment of $XPath_{\downarrow}$, we need to define both the maximum distance r going downward and the maximum distance s going upward that the formula can reach. We call the pair (r, s) the vertical depth of a formula. Formally, the **vertical depth** of a formula φ (notation: $vd(\varphi)$) is the pair $vd(\varphi) \in \mathbb{Z}_{\geq 0}^2$ defined in Table 2.

The **nesting depth** of a formula φ (notation: $nd(\varphi)$) is the maximum number of nested $[\]$ appearing in φ . See Table 2 for the formal definition.

Let (r, s, k) - $XPath_{\downarrow}^{\downarrow}$ be the set of all formulas φ in $XPath_{\downarrow}^{\downarrow}$ with $vd(\varphi) \leq (r, s)$ and $nd(\varphi) \leq k$.

Let \mathcal{T} and \mathcal{T}' be data trees, let $u \in T$ and $u' \in T'$. We say that \mathcal{T}, u and \mathcal{T}', u' are **equivalent for $XPath_{\downarrow}^{\downarrow}$** (notation: $\mathcal{T}, u \equiv^{\downarrow} \mathcal{T}', u'$) iff for all $\varphi \in XPath_{\downarrow}^{\downarrow}$, we have $\mathcal{T}, u \models \varphi$ iff $\mathcal{T}', u' \models \varphi$. \mathcal{T}, x and \mathcal{T}', x' are **(r, s) -equivalent** [resp. **(r, s, k) -equivalent**] for $XPath_{\downarrow}^{\downarrow}$, and we note it $\mathcal{T}, x \equiv_{r, s}^{\downarrow} \mathcal{T}', x'$ [resp. $\mathcal{T}, x \equiv_{r, s, k}^{\downarrow} \mathcal{T}', x'$] if they satisfy the same $XPath_{\downarrow}^{\downarrow}$ formulas φ so that $vd(\varphi) \leq (r, s)$ [resp. $vd(\varphi) \leq (r, s)$ and $nd(\varphi) \leq k$].

3.2.1 Normal form

We define a useful normal form for $XPath_{\downarrow}^{\downarrow}$ that will be implicitly used in the definition of bisimulation in the section. For $n \geq 0$, let \downarrow^n denote the concatenation of n symbols \downarrow . I.e., \downarrow^0 is the empty string λ , $\downarrow^1 = \downarrow$, and $\downarrow^{n+1} = \downarrow \downarrow^n$ (similarly for \uparrow^n).

A path expression α of $XPath_{\downarrow}^{\downarrow}$ is **downward** [resp. **upward**] if it is of the form $\downarrow^n[\varphi]$ [resp. $[\varphi]\uparrow^n$] for some $n \geq 0$ with $\varphi \in XPath_{\downarrow}^{\downarrow}$. For example, $\downarrow[\langle \uparrow \rangle]$ is a downward expression whereas $\downarrow[\langle \downarrow \rangle]$ is not. An **up-down** expression is any expression of the form $\varepsilon, \alpha^{\uparrow}, \alpha^{\downarrow}$ or $\alpha^{\uparrow} \alpha^{\downarrow}$ where α^{\uparrow} is upward and α^{\downarrow} is downward. Henceforth we will use $\alpha^{\uparrow}, \beta^{\uparrow}, \gamma^{\uparrow}$ to denote upward expressions and $\alpha^{\downarrow}, \beta^{\downarrow}, \gamma^{\downarrow}$ to denote downward expressions and $\alpha^{\uparrow \downarrow}, \beta^{\uparrow \downarrow}, \gamma^{\uparrow \downarrow}$ to denote up-down expressions. Note that in particular any downward or upward expression is an up-down expression. An $XPath_{\downarrow}^{\downarrow}$ formula or expression is in **up-down normal form** if every path expression contained in it is up-down and every data test is of the form $\langle \varepsilon \odot \alpha^{\uparrow \downarrow} \rangle$ with $\odot \in \{=, \neq\}$.

PROPOSITION 3.8. Let $\varphi \in (r, s, k)$ - $XPath_{\downarrow}^{\downarrow}$. There is $\varphi^{\uparrow \downarrow} \in XPath_{\downarrow}^{\downarrow}$ in up-down normal form such that

1. $\varphi^{\uparrow \downarrow} \equiv \varphi$;
2. $vd(\varphi^{\uparrow \downarrow}) = (r, s)$; and
3. $nd(\varphi^{\uparrow \downarrow}) \leq k \cdot (r + s + 2)$.

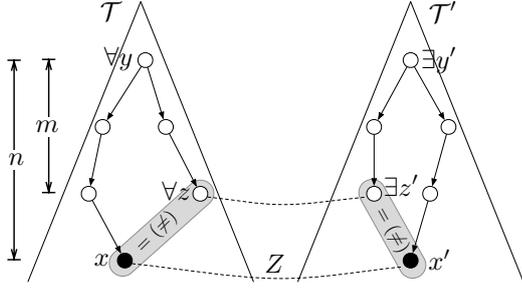


Figure 3: Zig clause of bisimulation for $\text{XPath}_{\perp}^{\dagger}$

3.2.2 Finite index

Contrary to the case of $\text{XPath}_{\perp}^{\downarrow}$ (cf., Proposition 3.1), the logical equivalence relation restricted to $\text{XPath}_{\perp}^{\dagger}$ -formulas of bounded vertical depth has infinitely many equivalence classes.

PROPOSITION 3.9. *If $r + s \geq 2$ then $\equiv_{r,s}^{\dagger}$ has infinite index.*

In the proof of the above proposition we need to use formulas with unbounded nesting depth. In fact, when restricted to bounded nesting depth there are only finitely many formulas up to logical equivalence, as stated next.

PROPOSITION 3.10. $\equiv_{r,s,k}^{\dagger}$ has finite index.

COROLLARY 3.11. $\{\mathcal{T}', u' \mid \mathcal{T}, u \equiv_{r,s,k}^{\dagger} \mathcal{T}', u'\}$ is definable by an (r, s, k) - $\text{XPath}_{\perp}^{\dagger}$ -formula.

3.2.3 Bisimulation and (r, s, k) -bisimulation

The advantage of the normal form presented in Section 3.2.1, is that it makes it possible to use a very simple notion of bisimulation. The disadvantage is that, since it does not preserve nesting depth, $\leftrightarrow_{r,s,k}^{\dagger}$ does not correspond precisely to $\equiv_{r,s,k}^{\dagger}$, although $\leftrightarrow^{\dagger}$ corresponds precisely to \equiv^{\dagger} . Nonetheless, we obtain, for all r, s, k ,

$$\leftrightarrow_{r,s,k}^{\dagger} \subseteq \equiv_{r,s,k}^{\dagger} \subseteq \leftrightarrow_{r,s,k \cdot (r+s+2)}^{\dagger}.$$

Let \mathcal{T} and \mathcal{T}' be two data-trees. We say that $u \in T$ and $u' \in T'$ are **bisimilar for $\text{XPath}_{\perp}^{\dagger}$** (notation: $\mathcal{T}, u \leftrightarrow^{\dagger} \mathcal{T}', u'$) iff there is a relation $Z \subseteq T \times T'$ such that uZu' and for all $x \in T$ and $x' \in T'$ we have

- **Harmony:** If xZx' then $\text{label}(x) = \text{label}(x')$,
- **Zig (Figure 3):** If xZx' , $y \xrightarrow{n} x$ and $y' \xrightarrow{m} x'$ then there are $y', z' \in T'$ such that $y' \xrightarrow{n} x'$, $y' \xrightarrow{m} z'$, $\text{data}(z) = \text{data}(x) \Leftrightarrow \text{data}(z') = \text{data}(x')$, and zZz' .
- **Zag:** If xZx' , $y' \xrightarrow{n} x'$ and $y' \xrightarrow{m} z'$ then there are $y, z \in T$ such that $y \xrightarrow{n} x$, $y \xrightarrow{m} z$, $\text{data}(z) = \text{data}(x) \Leftrightarrow \text{data}(z') = \text{data}(x')$, and zZz' .

Observe that contrary to the definition of $\leftrightarrow^{\downarrow}$, the conditions above do not require intermediate nodes to be related by Z . This is a direct consequence of the up-down normal form (Proposition 3.8).

We say that $u \in T$ and $u' \in T'$ are (r, s, k) -**bisimilar for $\text{XPath}_{\perp}^{\dagger}$** (notation: $\mathcal{T}, u \leftrightarrow_{r,s,k}^{\dagger} \mathcal{T}', u'$) if there is a family of relations $(Z_{\hat{r}, \hat{s}}^{\hat{k}})_{\hat{r} + \hat{s} \leq r+s, \hat{k} \leq k}$ in $T \times T'$ such that $uZ_{r,s}^k u'$ and for all $\hat{r} + \hat{s} \leq r + s$, $\hat{k} \leq k$, $x \in T$ and $x' \in T'$ we have that the following conditions hold.

- **Harmony:** If $xZ_{\hat{r}, \hat{s}}^{\hat{k}} x'$ then $\text{label}(x) = \text{label}(x')$.
- **Zig:** If $xZ_{\hat{r}, \hat{s}}^{\hat{k}} x'$, $y \xrightarrow{n} x$ and $y' \xrightarrow{m} x'$ with $n \leq \hat{s}$ and $m \leq \hat{r} + n$ then there are $y', z' \in T'$ such that $y' \xrightarrow{n} x'$, $y' \xrightarrow{m} z'$, and the following hold
 - (1) $\text{data}(z) = \text{data}(x) \Leftrightarrow \text{data}(z') = \text{data}(x')$,
 - (2) if $\hat{k} > 0$, $zZ_{\hat{r}', \hat{s}'}^{\hat{k}-1} z'$ for $\hat{r}' = \hat{r} + n - m$, $\hat{s}' = \hat{s} - n + m$.
- **Zag:** If $xZ_{\hat{r}, \hat{s}}^{\hat{k}} x'$, $y' \xrightarrow{n} x'$ and $y' \xrightarrow{m} z'$ with $n \leq \hat{s}$ and $m \leq \hat{r} + n$ then there are $y, z \in T$ such that $y \xrightarrow{n} x$, $y \xrightarrow{m} z$, and items (1) and (2) above are verified.

OBSERVATION 3.12. *If $xZ_{\hat{r}, \hat{s}}^{\hat{k}} x'$, $y \xrightarrow{n} x$ and $y' \xrightarrow{n} x'$ then it follows that $yZ_{\hat{r}', \hat{s}'}^{\hat{k}-1} y'$, for $\hat{r}' = \hat{r} + n$, $\hat{s}' = \hat{s} - n$. The same occurs with Z instead of $Z_{\hat{r}, \hat{s}}^{\hat{k}}$ for the case of bisimilarity.*

For a data tree \mathcal{T} and $u \in T$, let $\mathcal{T}|_r^s u$ denote the subtree of \mathcal{T} induced by

$$\{v \in T \mid (\exists m \leq s) (\exists n \leq r + m) (\exists w \in T) w \xrightarrow{m} u \wedge w \xrightarrow{n} v\}.$$

3.2.4 Equivalence and bisimulation

The next result says that $\leftrightarrow^{\dagger}$ coincides with \equiv^{\dagger} on finitely branching data trees, and states precisely in what way $\leftrightarrow_{r,s,k}^{\dagger}$ is related to $\equiv_{r,s,k}^{\dagger}$.

THEOREM 3.13.

1. $\mathcal{T}, u \leftrightarrow^{\dagger} \mathcal{T}', u'$ implies $\mathcal{T}, u \equiv^{\dagger} \mathcal{T}', u'$. The converse also holds when \mathcal{T} and \mathcal{T}' are finitely branching.
2. $\mathcal{T}, u \leftrightarrow_{r,s,k \cdot (r+s+2)}^{\dagger} \mathcal{T}', u'$ implies $\mathcal{T}, u \equiv_{r,s,k}^{\dagger} \mathcal{T}', u'$.
3. $\mathcal{T}, u \equiv_{r,s,k}^{\dagger} \mathcal{T}', u'$ implies $\mathcal{T}, u \leftrightarrow_{r,s,k}^{\dagger} \mathcal{T}', u'$.

COROLLARY 3.14. $\leftrightarrow_{r,s,k}^{\dagger}$ has finite index.

4. SIMULATION

In this section we define notions of directed (non-symmetric) simulations for $\text{XPath}_{\perp}^{\downarrow}$ and $\text{XPath}_{\perp}^{\dagger}$, as it is done, e.g., in [16] for some modal logics. We obtain results similar to Theorems 3.5 and 3.13 but relating each simulation notion with the corresponding logical implication.

We say that an XPath_{\perp} formula is **positive** if it contains no negation \neg and no inequality data tests $\langle \alpha \neq \beta \rangle$. For \mathcal{L} one of $\text{XPath}_{\perp}^{\downarrow}$, $\text{XPath}_{\perp}^{\dagger}$, $\text{XPath}_{\perp}^{\downarrow+}$, or $\text{XPath}_{\perp}^{\dagger+}$, we write \mathcal{L}^+ for the positive fragment of \mathcal{L} .

A **simulation for $\text{XPath}_{\perp}^{\downarrow}$** [resp. for $\text{XPath}_{\perp}^{\dagger}$] is simply a bisimulation from which the *Zag* clause and half of the first condition in the *Zig* clause have been omitted. Observe that simulations need not be symmetric.

Formally, we say that $u \in T$ is **similar to $u' \in T'$ for $\text{XPath}_{\perp}^{\downarrow}$** (notation: $\mathcal{T}, u \rightarrow^{\downarrow} \mathcal{T}', u'$) iff there is a relation $Z \subseteq T \times T'$ such that uZu' and for all $x \in T$ and $x' \in T'$ we have

- **Harmony:** If xZx' then $label(x) = label(x')$.
- **Zig:** If xZx' , $x \xrightarrow{n} v$ and $x \xrightarrow{m} w$ then there are $v', w' \in T'$ such that $x' \xrightarrow{n} v'$, $x' \xrightarrow{m} w'$ and
 1. $data(v) = data(w) \Rightarrow data(v') = data(w')$,
 2. $(\xrightarrow{i} v) Z (\xrightarrow{i} v')$ for all $0 \leq i < n$, and
 3. $(\xrightarrow{i} w) Z (\xrightarrow{i} w')$ for all $0 \leq i < m$.

$u \in T$ is **similar to** $u' \in T'$ for $\mathbf{XPath}_{\leq}^{\downarrow}$ (notation: $\mathcal{T}, u \xrightarrow{\dagger} \mathcal{T}', u'$) iff there is a relation $Z \subseteq T \times T'$ such that uZu' and for all $x \in T$ and $x' \in T'$ we have

- **Harmony:** If xZx' then $label(x) = label(x')$.
- **Zig:** If xZx' , $y \xrightarrow{n} x$ and $y \xrightarrow{m} z$ then there are $y', z' \in T'$ such that $y' \xrightarrow{n} x'$, $y' \xrightarrow{m} z'$, zZz' , and if $data(z) = data(x)$ then $data(z') = data(x')$.

Relations $\xrightarrow{\downarrow}$ and $\xrightarrow{\downarrow}_{r,s,k}$ are defined accordingly. We define one-way (non-symmetric) logical implication between models as follows. We write $\mathcal{T}, u \Rightarrow^{\downarrow} \mathcal{T}', u'$ for

$$(\forall \varphi \in \mathbf{XPath}_{\leq}^{\downarrow}) [\mathcal{T}, u \models \varphi \Rightarrow \mathcal{T}', u' \models \varphi].$$

Define \Rightarrow^{\downarrow} , \Rightarrow^{\downarrow} , and $\Rightarrow^{\downarrow}_{r,s,k}$ in an analogous way for ℓ - $\mathbf{XPath}_{\leq}^{\downarrow}$, $\mathbf{XPath}_{\leq}^{\downarrow}$, (r, s, k) - $\mathbf{XPath}_{\leq}^{\downarrow}$, respectively. As for bisimulation, we have that \Rightarrow coincides with \Rightarrow .

THEOREM 4.1.

1. Let $\dagger \in \{\downarrow, \uparrow\}$. $\mathcal{T}, u \xrightarrow{\dagger} \mathcal{T}', u'$ implies $\mathcal{T}, u \Rightarrow^{\dagger} \mathcal{T}', u'$. The converse holds when \mathcal{T}' is finitely branching.
2. $\mathcal{T}, u \xrightarrow{\downarrow} \mathcal{T}', u'$ iff $\mathcal{T}, u \Rightarrow^{\downarrow} \mathcal{T}', u'$.
3. $\mathcal{T}, u \xrightarrow{\downarrow}_{r,s,k \cdot (r+s+2)} \mathcal{T}', u'$ implies $\mathcal{T}, u \Rightarrow^{\downarrow}_{r,s,k} \mathcal{T}', u'$.
4. $\mathcal{T}, u \Rightarrow^{\downarrow}_{r,s,k} \mathcal{T}', u'$ implies $\mathcal{T}, u \xrightarrow{\downarrow}_{r,s,k} \mathcal{T}', u'$.

We say that \mathcal{T}' is a **substructure** of \mathcal{T} if \mathcal{T}' is a data tree which results from removing some nodes of \mathcal{T} , i.e., $T' \subseteq T$ and for all $u, v \in T'$ we have: 1) $u \rightarrow v$ on \mathcal{T} iff $u \rightarrow v$ on \mathcal{T}' ; 2) $label(u)$ on \mathcal{T}' equals $label(u)$ on \mathcal{T} ; and 3) $data(u)$ on \mathcal{T}' equals $data(u)$ on \mathcal{T} . Equivalently, seen as σ -structures, \mathcal{T}' is the σ -substructure of \mathcal{T} induced by $T' \subseteq T$. One can verify that the identity on T' is a simulation for $\mathbf{XPath}_{\leq}^{\downarrow}$ from \mathcal{T}' to \mathcal{T} .

LEMMA 4.2. If \mathcal{T}' is a substructure of \mathcal{T} and $u' \in T'$ then $\mathcal{T}', u' \xrightarrow{\dagger} \mathcal{T}, u'$.

We obtain that the formulas of $\mathbf{XPath}_{=}$ invariant under simulations are, precisely, the positive ones.

THEOREM 4.3.

1. $\varphi \in \mathbf{XPath}_{\leq}^{\downarrow}$ is $\xrightarrow{\downarrow}$ -invariant [resp. $\xrightarrow{\downarrow}_{\ell}$] iff it is equivalent to a formula of $\mathbf{XPath}_{\leq}^{\downarrow}$ [resp. ℓ - $\mathbf{XPath}_{\leq}^{\downarrow}$].
2. $\varphi \in \mathbf{XPath}_{\leq}^{\downarrow}$ is $\xrightarrow{\downarrow}$ -invariant iff it is equivalent to a formula of $\mathbf{XPath}_{\leq}^{\downarrow}$.
3. If $\varphi \in \mathbf{XPath}_{\leq}^{\downarrow}$ is $\xrightarrow{\downarrow}_{r,s,k}$ -invariant then it is equivalent to a formula of (r, s, k) - $\mathbf{XPath}_{\leq}^{\downarrow}$.
4. If $\varphi \in \mathbf{XPath}_{\leq}^{\downarrow}$ is equivalent to a formula of (r, s, k) - $\mathbf{XPath}_{\leq}^{\downarrow}$ then φ is $\xrightarrow{\downarrow}_{r,s,k'}$ -invariant, for $k' = k \cdot (r+s+2)$.

5. ADDING TRANSITIVITY

As it happens, for example, with the basic modal logic and propositional dynamic logic, the same notion of bisimulation [resp. simulation] of each logic captures the logical equivalence [resp. logical implication] for the corresponding fragments including the reflexive-transitive closure of the axes which are present. Intuitively, this occurs because \downarrow_* is an infinite union of compositions of \downarrow , and similarly for \uparrow .

Let $\equiv^{\downarrow*}$ and $\equiv^{\uparrow*}$ be the logical equivalence relation for $\mathbf{XPath}_{\leq}^{\downarrow*}$ and $\mathbf{XPath}_{\leq}^{\uparrow*}$ respectively, and let $\Rightarrow^{\downarrow*}$ and $\Rightarrow^{\uparrow*}$ be the logical implication for $\mathbf{XPath}_{\leq}^{\downarrow*}$ and $\mathbf{XPath}_{\leq}^{\uparrow*}$ respectively.

THEOREM 5.1. Let $\dagger \in \{\downarrow, \uparrow\}$.

1. $\mathcal{T}, u \xrightarrow{\dagger} \mathcal{T}', u'$ implies $\mathcal{T}, u \equiv^{\dagger} \mathcal{T}', u'$. The converse also holds when \mathcal{T}' is finitely branching.
2. $\mathcal{T}, u \xrightarrow{\dagger} \mathcal{T}', u'$ implies $\mathcal{T}, u \Rightarrow^{\dagger} \mathcal{T}', u'$. The converse also holds when \mathcal{T}' is finitely branching.

6. CHARACTERIZATION

In §6.1 we show that there is a truth-preserving translation from $\mathbf{XPath}_{\leq}^{\downarrow}$ to first-order logic over an appropriate signature. In §6.2 we characterize $\mathbf{XPath}_{\leq}^{\downarrow}$ as the fragment of first-order logic $\xrightarrow{\downarrow}$ -invariant over data trees. In §6.3 we show that this result fails for $\mathbf{XPath}_{\leq}^{\downarrow}$ in general, but a weaker result can still be proved.

6.1 Translating to first-order logic

We say that an $\mathbf{XPath}_{\leq}^{\downarrow}$ -path expression α is in **simple normal form** if it is of the form

$$[\varphi_0]o_1[\varphi_1]o_2 \cdots o_n[\varphi_n],$$

for $n \geq 0$, $\varphi_i \in \mathbf{XPath}_{\leq}^{\downarrow}$, and $o_i \in \{\downarrow, \uparrow\}$. For any $\mathbf{XPath}_{\leq}^{\downarrow}$ -[resp. $\mathbf{XPath}_{\leq}^{\downarrow}$ -] path expression α there is an equivalent $\mathbf{XPath}_{\leq}^{\downarrow}$ -[resp. $\mathbf{XPath}_{\leq}^{\downarrow}$ -] path expression α' in simple normal form. Further, α' can be computed in polynomial time from α .¹ We say that an $\mathbf{XPath}_{\leq}^{\downarrow}$ -formula φ is in **simple normal form** if each path expression α occurring in φ is in simple normal form.

Fix the signature σ with binary relations \rightsquigarrow and \approx , and a unary predicate P_a for each $a \in \mathbb{A}$. Any data tree \mathcal{T} can be seen as a first-order σ -structure such that

$$\begin{aligned} \rightsquigarrow^{\mathcal{T}} &= \{(x, y) \in T^2 \mid y \text{ is a child of } x\}; \\ \approx^{\mathcal{T}} &= \{(x, y) \in T^2 \mid data(x) = data(y)\}; \\ P_a^{\mathcal{T}} &= \{x \in T \mid label(x) = a\}. \end{aligned}$$

We define the following translation Tr mapping $\mathbf{XPath}_{\leq}^{\downarrow}$ formulas in simple normal form to first-order σ -formulas:

$$\begin{aligned} \text{Tr}_x(a) &= P_a(x) & (a \in \mathbb{A}) \\ \text{Tr}_x(\varphi \dagger \psi) &= \text{Tr}_x(\varphi) \dagger \text{Tr}_x(\psi) & (\dagger \in \{\wedge, \vee\}) \\ \text{Tr}_x(\neg \varphi) &= \neg \text{Tr}_x(\varphi) \\ \text{Tr}_x(\langle \alpha \rangle) &= (\exists \bar{y})(x = y_0 \wedge \text{Tr}_{\bar{y}}(\alpha)) \\ \text{Tr}_x(\langle \alpha = \beta \rangle) &= (\exists \bar{y})(\exists \bar{z})(x = y_0 \wedge x = z_0 \wedge y_n \approx z_m \wedge \\ &\quad \text{Tr}_{\bar{y}}(\alpha) \wedge \text{Tr}_{\bar{z}}(\beta)) \\ \text{Tr}_x(\langle \alpha \neq \beta \rangle) &= (\exists \bar{y})(\exists \bar{z})(x = y_0 \wedge x = z_0 \wedge y_n \not\approx z_m \wedge \end{aligned}$$

¹Note that this proposition holds only for paths expressions without union.

$$\text{Tr}_{\bar{y}}(\alpha) \wedge \text{Tr}_{\bar{z}}(\beta)$$

$$\text{Tr}_{\bar{y}}(\alpha) = \bigwedge_{i=0}^{n-1} o_{i+1}(y_i, y_{i+1}) \wedge \bigwedge_{i=0}^n \text{Tr}_{y_i}(\varphi_i),$$

where $\bar{y} = y_0, \dots, y_n$ and $\bar{z} = z_0, \dots, z_m$, and are fresh when quantified in the fourth and fifth definition;

$$\alpha = [\varphi_0]o_1[\varphi_1]o_2[\varphi_2]o_3 \cdots o_n[\varphi_n];$$

$$\beta = [\psi_0]o'_1[\psi_1]o'_2[\psi_2]o'_3 \cdots o'_m[\psi_m];$$

$o_i, o'_i \in \{\downarrow, \uparrow\}$; $o_j(u, v)$ represents $u \rightsquigarrow v$ if $o_j = \downarrow$, and $v \rightsquigarrow u$ otherwise. For $\varphi \in \text{XPath}_{\pm}^{\downarrow}$ we have $\mathcal{T}, u \models \varphi$ iff $\mathcal{T} \models \text{Tr}_x(\varphi)(u)$.

6.2 Downward XPath

Let $\text{FO}(\sigma)$ be the set of first-order formulas over a given signature σ , and let \mathcal{C} be a class of σ -models. An $\text{FO}(\sigma)$ -formula $\varphi(x)$ is ℓ -local if for all data trees \mathcal{T} and $u \in T$, we have $\mathcal{T} \models \varphi(u) \Leftrightarrow \mathcal{T}|_{\ell u} \models \varphi(u)$. Finally, for $\varphi \in \text{FO}(\sigma)$ let $\text{qr}(\varphi)$ be its quantifier rank, i.e., the depth of nesting of its quantifiers.

Observe that the following result has two readings: one classical, and one restricted to finite models.

THEOREM 6.1 (CHARACTERIZATION). *Let $\varphi(x) \in \text{FO}(\sigma)$. The following are equivalent:*

- (i) φ is $\Leftrightarrow^{\downarrow}$ -invariant over [finite] data-trees;
- (ii) φ is logically equivalent over [finite] data-trees to an ℓ -XPath $_{\pm}^{\downarrow}$ -formula, where $\ell = 2^{\text{qr}(\varphi)} - 1$.

PROOF. The implication (ii) \Rightarrow (i) follows straightforwardly from Theorem 3.5. The proof of (i) \Rightarrow (ii) goes as follows: First, we show that any $\Leftrightarrow^{\downarrow}$ -invariant $\varphi(x) \in \text{FO}(\sigma)$ is ℓ -local for $\ell = 2^{\text{qr}(\varphi)} - 1$ (Proposition 6.2). Then, we prove that any $\Leftrightarrow^{\downarrow}$ -invariant $\varphi(x) \in \text{FO}(\sigma)$ that is ℓ -local is $\Leftrightarrow_{\ell}^{\downarrow}$ -invariant. Finally, we show that any $\text{FO}(\sigma)$ -definable property which is $\Leftrightarrow_{\ell}^{\downarrow}$ -invariant is definable in ℓ -XPath $_{\pm}^{\downarrow}$. \square

PROPOSITION 6.2. *Any $\Leftrightarrow^{\downarrow}$ -invariant $\varphi(x) \in \text{FO}(\sigma)$ over [finite] data-trees is ℓ -local for $\ell = 2^{\text{qr}(\varphi)} - 1$.*

PROOF. We follow Otto's proof [20]. Assume that $\varphi(x) \in \text{FO}(\sigma)$ is $\Leftrightarrow^{\downarrow}$ -invariant, let $q = \text{qr}(\varphi)$, and put $\ell = 2^q - 1$. Given a data tree \mathcal{T} and $u \in T$ it suffices to show the existence of data trees \mathcal{T}' and \mathcal{T}'' , with corresponding elements $u' \in T'$ and $u'' \in T''$ such that

- (a) $\mathcal{T}', u' \Leftrightarrow^{\downarrow} \mathcal{T}, u$,
- (b) $\mathcal{T}'', u'' \Leftrightarrow^{\downarrow} (\mathcal{T}|_{\ell u}), u$, and
- (c) $\mathcal{T}', u' \equiv_q \mathcal{T}'', u''$.

Indeed, from the above conditions it follows that

$$\begin{aligned} \mathcal{T} \models \varphi(u) \text{ iff } \mathcal{T}' \models \varphi(u') & \quad \text{((a) and } \Leftrightarrow^{\downarrow}\text{-inv. of } \varphi) \\ \text{iff } \mathcal{T}'' \models \varphi(u'') & \quad \text{(c)} \\ \text{iff } (\mathcal{T}|_{\ell u}) \models \varphi(u), & \quad \text{((b) and } \Leftrightarrow^{\downarrow}\text{-inv. of } \varphi) \end{aligned}$$

and hence φ is ℓ -local. By Proposition 3.3 one may assume that $u \in T$ is the root of \mathcal{T} .

We define \mathcal{T}' and \mathcal{T}'' , as structures that are disjoint copies of sufficiently many isomorphic copies of \mathcal{T} and $\mathcal{T}|_{\ell u}$, respectively, all tied together by some common root. Both

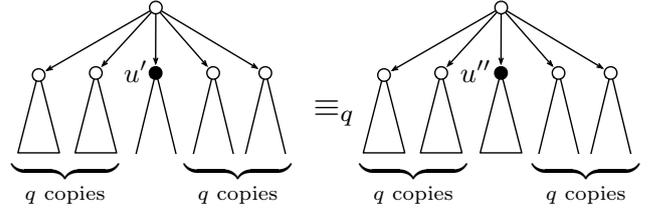


Figure 4: Definition of \mathcal{T}', u' and \mathcal{T}'', u'' .

structures have q isomorphic copies of both \mathcal{T} and $\mathcal{T}|_{\ell u}$, and only distinguish themselves by the nature of the one extra subtree, in which u' and u'' live, respectively: u' is the root of one of the copies of \mathcal{T} and u'' is the root of one of the copies of $\mathcal{T}|_{\ell u}$. We indicate the two structures in the diagram of Figure 4, with distinguished elements u' and u'' marked by \bullet ; the open cones stand for copies of \mathcal{T} , the closed cones for copies of $\mathcal{T}|_{\ell u}$. The new isomorphic copies have the same data values as the original one. The new root has an arbitrary, fixed, data value and label.

By Proposition 3.4, it is straightforward that conditions (a) and (b) are satisfied. Condition (c) is true because one can exhibit a strategy for player **II** in the q -round Ehrenfeucht-Fraïssé game on structures \mathcal{T}' and \mathcal{T}'' . The strategy is exactly the same used in [20]. \square

6.3 Vertical XPath

The analog of Theorem 6.1 fails for XPath $_{\pm}^{\downarrow}$:

LEMMA 6.3. *The $\text{FO}(\sigma)$ -formula*

$$(\exists x) P_a(x)$$

is $\Leftrightarrow^{\downarrow}$ -invariant though not logically equivalent over [finite] data-trees to any XPath $_{\pm}^{\downarrow}$ -formula.

Hence XPath $_{\pm}^{\downarrow}$ is not the fragment of $\text{FO}(\sigma)$ which is $\Leftrightarrow^{\downarrow}$ -invariant over [finite] data-trees. However, the following proposition still holds for the case of XPath $_{\pm}^{\downarrow}$:

PROPOSITION 6.4. *Let $k' = k \cdot (r+s+2)$. If $\varphi(x) \in \text{FO}(\sigma)$ is $\Leftrightarrow_{r,s,k'}^{\downarrow}$ -invariant over [finite] data-trees, then there is $\psi \in (r, s, k)$ -XPath $_{\pm}^{\downarrow}$ such that $\text{Tr}_x(\psi)$ is logically equivalent to φ over [finite] data-trees.*

Notice that the counterexample in Lemma 6.3 is an unrestricted, existential formula. One may wonder if it might be possible to extend the expressive power of XPath $_{\pm}^{\downarrow}$ to account for unrestricted quantification. The natural candidate would be the modal operator E (usually known as the existential modality) which, intuitively, let us express that there is some node in the model where a formula holds. But even with the additional expressive power provided by E the analog of Theorem 6.1 fails. Formally, consider the logic XPath $_{\pm}^{\downarrow E}$, which results from adding the operator E to XPath $_{\pm}^{\downarrow}$ with the following semantics: $\llbracket E\varphi \rrbracket^{\mathcal{T}} = T$ if $\llbracket \varphi \rrbracket^{\mathcal{T}} \neq \emptyset$, and $\llbracket E\varphi \rrbracket^{\mathcal{T}} = \emptyset$ otherwise.

The following lemma shows a counterexample to the analog of Theorem 6.1, showing that XPath $_{\pm}^{\downarrow E}$ is not the fragment of $\text{FO}(\sigma)$ $\Leftrightarrow^{\downarrow}$ -invariant over [finite] data-trees.

LEMMA 6.5. *The FO(σ)-formula*

$$(\exists y, z) [y \approx z \wedge P_a(y) \wedge P_b(z)]$$

is $\leftrightarrow^\downarrow$ -invariant though not logically equivalent over [finite] data-trees to any $XPath_{\leq}^{\downarrow E}$ -formula.

7. APPLICATIONS

We devote this section to exemplify how the model theoretic tools we developed can be used to show expressiveness results for $XPath_{\leq}$. We do not intend to be comprehensive; rather we will exhibit a number of different results that show possible uses of the notions of bisimulation we introduced.

7.1 Expressiveness hierarchies

Define $\equiv_{\ell, k}^{\downarrow}$ as the equivalence $\equiv_{\ell}^{\downarrow}$ restricted to formulas of nesting depth at most k , that is, $\mathcal{T}, u \equiv_{\ell, k}^{\downarrow} \mathcal{T}', u'$ iff for all $\varphi \in XPath_{\leq}^{\downarrow}$ such that $dd(\varphi) \leq \ell$ and $nd(\varphi) \leq k$ we have $\mathcal{T}, u \models \varphi$ iff $\mathcal{T}', u' \models \varphi$. Define a more fine-grained notion of bisimulation in a similar way. We say that $u \in T$ and $u' \in T'$ are (ℓ, k) -bisimilar for $XPath_{\leq}^{\downarrow}$ (notation: $\mathcal{T}, u \leftrightarrow_{\ell, k}^{\downarrow} \mathcal{T}', u'$) if there is a family of relations $(Z_{j, t})_{j \leq \ell, t \leq k}$ in $T \times T'$ such that $u Z_{\ell, k} u'$ and for all $j \leq \ell, t \leq k, x \in T$ and $x' \in T'$ we have

- **Harmony:** If $x Z_{j, t} x'$ then $label(x) = label(x')$.
- **Zig:** If $x Z_{j, t} x', x \xrightarrow{n} v$ and $x \xrightarrow{m} w$ with $n, m \leq j$ then there are $v', w' \in T'$ such that $x' \xrightarrow{n} v', x' \xrightarrow{m} w'$ and
 1. $data(v) = data(w) \Leftrightarrow data(v') = data(w')$,
 2. if $t > 0, (\xrightarrow{i} v) Z_{j-n+i, t-1} (\xrightarrow{i} v')$ for all $0 \leq i < n$, and
 3. if $t > 0, (\xrightarrow{i} w) Z_{j-m+i, t-1} (\xrightarrow{i} w')$ for all $0 \leq i < m$.
- **Zag:** If $x Z_{j, t} x', x' \xrightarrow{n} v'$ and $x' \xrightarrow{m} w'$ with $n, m \leq j$ then there are $v, w \in T$ such that $x \xrightarrow{n} v, x \xrightarrow{m} w$ and items 1, 2 and 3 above are verified.

Following the same ideas used in Propositions 3.6 and 3.7, it is easy to show that (ℓ, k) -bisimulations characterize (ℓ, k) -equivalence.

PROPOSITION 7.1. $\mathcal{T}, u \leftrightarrow_{\ell, k}^{\downarrow} \mathcal{T}', u'$ iff $\mathcal{T}, u \equiv_{\ell, k}^{\downarrow} \mathcal{T}', u'$.

The following theorem characterizes when an increase in nesting depth results in an increase in expressive power (see Figure 5). We conjecture that a similar hierarchy holds in the absence of data values, but this is not a direct consequence of our result.

THEOREM 7.2. *For all $\ell, k \geq 0, i \geq 1$,*

$$\begin{aligned} \equiv_{\ell, 0}^{\downarrow} \supsetneq \equiv_{\ell, 1}^{\downarrow} \supsetneq \cdots \supsetneq \equiv_{\ell, \ell}^{\downarrow} = \equiv_{\ell, \ell+i}^{\downarrow}, \text{ and} \\ \equiv_{\ell, k}^{\downarrow} \supsetneq \equiv_{\ell+i, k}^{\downarrow}. \end{aligned}$$

7.2 Safe operations on models

Bisimulations can also be used to show that certain operations on models preserve truth. Such operations are usually called *safe* for a given logic, as they can be applied to a model without changing the truth values of any formula in

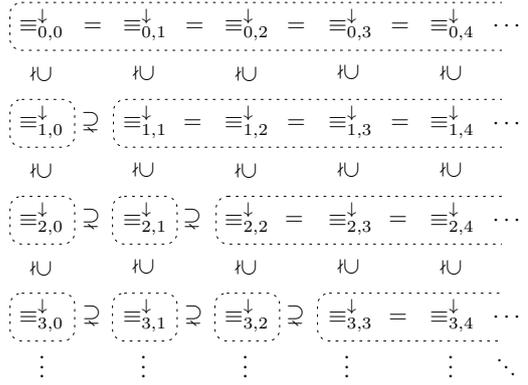


Figure 5: Hierarchy of $XPath_{\leq}^{\downarrow}$.

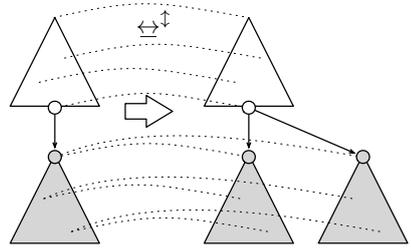


Figure 6: Closure under subtree replication.

the language. Proposition 3.3, for example, is already an example of this kind of results showing that the class of models of a formula is closed under sub-model generation. We will now show a more elaborate example.

We say that \mathcal{T}' is a **subtree replication** of \mathcal{T} , if \mathcal{T}' is the result of inserting $\mathcal{T}|x$ into \mathcal{T} as a sibling of x , where x is any node of \mathcal{T} different from the root. Figure 6 gives a schematic representation of this operation.

PROPOSITION 7.3. $XPath_{\leq}^{\downarrow \uparrow \uparrow^*}$ is closed under subtree replication, i.e. if \mathcal{T}' is a subtree replication of \mathcal{T} , and $u \in T$ then $\mathcal{T}', u \equiv_{\uparrow \uparrow^*}^{\downarrow} \mathcal{T}, u$.

PROOF. Suppose that $x \in T$ is not the root of \mathcal{T} , and that \mathcal{T}' is the result of inserting $\mathcal{T}|x$ into \mathcal{T} as a sibling of x . Let us call \mathcal{T}_x to the new copy of $\mathcal{T}|x$ inserted into \mathcal{T}' , and let X be the set of nodes of $\mathcal{T}|x$. Furthermore, if $v \in X$ then v_x is the corresponding node of \mathcal{T}_x . Nodes v and v_x have the same label and data value, and the position of v in $\mathcal{T}|x$ coincides with the position of v_x in \mathcal{T}_x .

By Theorem 5.1, it suffices to verify that $\mathcal{T}, u \leftrightarrow^{\downarrow \uparrow \uparrow^*} \mathcal{T}', u$ via $Z \subseteq T \times T'$ defined by:

$$Z = \{(y, y) \mid y \in T\} \cup \{(v, v_x) \mid v \in X\}$$

(Z is depicted as dotted lines in Figure 6). \square

7.3 Non-expressivity results

Finally, we will use bisimulation to show the expressivity limits of different fragments of $XPath$. Let $key(a)$ be the property stating that every node with label a has a different data value. Let $fk(a, b)$ (for *foreign key*) be the property $(\forall x)[P_a(x) \Rightarrow (\exists y)[P_b(y) \wedge x \sim y]]$.

PROPOSITION 7.4.

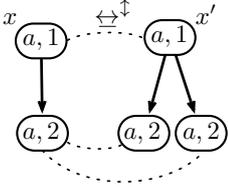


Figure 7: $key(a)$ not in $XPath_{\perp}^{\downarrow\downarrow*}$.

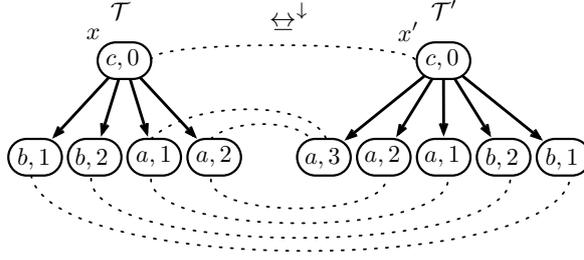


Figure 8: $fk(a,b)$ not in $XPath_{\perp}^{\downarrow\downarrow*}$.

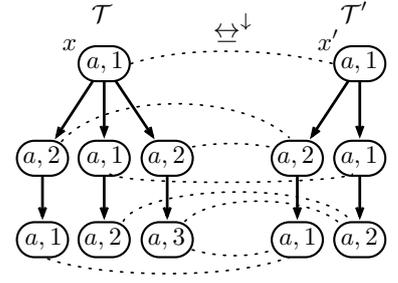


Figure 9: $dist_3$ not in $XPath_{\perp}^{\downarrow\downarrow*}$.

1. $key(a)$ is not expressible in $XPath_{\perp}^{\downarrow\downarrow*}$.
2. $fk(a,b)$ is expressible in $XPath_{\perp}^{\downarrow\downarrow*}$ but it is not expressible in $XPath_{\perp}^{\downarrow\downarrow*}$ or $XPath_{\perp}^{\downarrow\downarrow*+}$.

PROOF. The first item follows from Proposition 7.3. Since the logic is closed under subtree replication, the trees of Figure 7 are equivalent. As $key(a)$ holds in one and not in the other, the statement follows.

For the second item, it is easy to see that $fk(a,b)$ is expressible with the formula $\neg(\uparrow^* \downarrow_* [a \wedge \neg(\varepsilon = \uparrow^* \downarrow_* [b])])$. However, this property cannot be expressed in $XPath_{\perp}^{\downarrow\downarrow*}$ because the models \mathcal{T} and \mathcal{T}' in Figure 8 are bisimilar for $XPath_{\perp}$ via Z , depicted as dotted lines. Since \mathcal{T}, x satisfies $fk(a,b)$ but \mathcal{T}', x' does not, from Theorem 5.1 it follows that $fk(a,b)$ is not expressible in $XPath_{\perp}^{\downarrow\downarrow*}$.

Finally, suppose there exists $\psi \in XPath_{\perp}^{\downarrow\downarrow*+}$ expressing $fk(a,b)$. Since \mathcal{T} is a substructure of \mathcal{T}' we have $\mathcal{T}, x \xrightarrow{\downarrow} \mathcal{T}', x$ by Lemma 4.2. By Theorem 5.1(2) and the fact that $\mathcal{T}, x \models \psi$, we have $\mathcal{T}', x \models \psi$, which is a contradiction. \square

Let $dist_3(x)$ be the property stating that there are nodes y, z so that $x \rightarrow y \rightarrow z$ and x, y, z have pairwise distinct data values.

PROPOSITION 7.5.

1. $dist_3$ is expressible in $XPath_{\perp}^{\downarrow}$;
2. $dist_3$ is not expressible in $XPath_{\perp}^{\downarrow\downarrow*}$;
3. neither $dist_3$ nor its complement can be expressed in $XPath_{\perp}^{\downarrow\downarrow*+}$.

PROOF. For 1, one can check that $\mathcal{T}, x \models \varphi$ iff \mathcal{T}, x satisfies $dist_3$, for $\varphi = \langle \varepsilon \neq \downarrow\downarrow[(\varepsilon \neq \uparrow)] \rangle$.

Let us see 2. Consider the data trees \mathcal{T}, x and \mathcal{T}', x' depicted in Figure 9. It is straightforward that \mathcal{T}, x satisfies $dist_3$ and \mathcal{T}', x' does not.

Let v'_1 and v'_2 be the leaves of \mathcal{T}' and let v be the only node of \mathcal{T} with data value 3. One can check that $\mathcal{T}, x \xrightarrow{\downarrow} \mathcal{T}', x'$ via $Z \subseteq T \times T'$ defined by

$$Z = \{ \langle u, u' \rangle \mid h(u) = h(u') \wedge data(u) = data(u') \} \cup \{ \langle v, v'_1 \rangle, \langle v, v'_2 \rangle \},$$

where $h(y)$ denotes the height of y , i.e., the distance from y to the root of the corresponding tree (Z is depicted as dotted lines in Figure 9). Since \mathcal{T}, x satisfies $dist_3$ but \mathcal{T}', x' does not, from Theorem 5.1 it follows that $dist_3$ is not expressible in $XPath_{\perp}^{\downarrow\downarrow*}$.

For 3, one can verify that $\mathcal{T}, x \xrightarrow{\downarrow} \mathcal{T}', x'$ via Z as defined above. If $dist_3$ were definable in $XPath_{\perp}^{\downarrow\downarrow*+}$ via ψ and the fact that $\mathcal{T}, x \models \psi$, by Theorem 5.1(2) we would have $\mathcal{T}', x' \models \psi$, and this is a contradiction.

Let $\overline{dist_3}$ denote the complement of $dist_3$, i.e., $\overline{dist_3}(x)$ iff for all y, z so that $x \rightarrow y \rightarrow z$, we have that x, y, z do not have pairwise distinct data values. Now \mathcal{T}', x' satisfies $\overline{dist_3}$ and \mathcal{T}, x does not. Since \mathcal{T}' is a substructure of \mathcal{T} , by an argument analog to the one used in the proof of Proposition 7.4-2, we conclude that $\overline{dist_3}$ is not expressible in $XPath_{\perp}^{\downarrow\downarrow*+}$. \square

8. DISCUSSION

In this article we studied model theoretic properties of XPath over both finite and arbitrary data trees using bisimulations. One of the main results we discuss is the characterization of the downward and vertical fragments of XPath as the fragments of first-order logic which are invariant under suitable notions of bisimulation. This can be seen as a first step in the larger program of studying the model theory and expressiveness of XPath with data values and, more generally, of logics on data trees. It would be interesting to study notions of bisimulation with only descendant; or characterizations of XPath with child and descendant, as a fragment of FO with the descendant relation on data trees. We did not consider XPath with horizontal navigation between siblings, such as the axes `next-sibling` and `previous-sibling`. In fact, adding these axes results in a fragment that is somewhat less interesting since the adequate bisimulation notion on finite data trees corresponds precisely to data tree isomorphism modulo renaming of data values.

In Section 7 we show a number of concrete application of the model theoretic tools we developed, discussing both expressivity and non-expressivity results. We also show examples of operations which are safe for a given XPath fragment. It would be worthwhile to devise other model operations that preserve truth of XPath formulas as we show is the case for *subtree replication*.

An important application of bisimulation is as a minimization method: given a data tree \mathcal{T}_1 we want to find a data tree \mathcal{T}_2 , as small as possible, so that \mathcal{T}_1 and \mathcal{T}_2 are bisimilar for some fragment \mathcal{L} of XPath. Since \mathcal{L} cannot distinguish between \mathcal{T}_1 and \mathcal{T}_2 , we can use \mathcal{T}_2 as representative of \mathcal{T}_1 while the expressive power of \mathcal{L} is all that is required by a given application. The complexity of several inference tasks (e.g., model checking) depends directly on the model size. This is why in some cases it may be profitable to first apply

a minimization step. The existence of efficient minimization algorithms is intimately related to bisimulations: we can minimize a data tree \mathcal{T} by partitioning it in terms of its coarsest auto-bisimulation. We plan to design and implement algorithms for data tree minimization using bisimulation and investigate their computational complexity.

References

- [1] M. Benedikt, W. Fan, and F. Geerts. XPath satisfiability in the presence of DTDs. *Journal of the ACM*, 55(2):1–79, 2008.
- [2] M. Benedikt and C. Koch. XPath leashed. *ACM Computing Surveys*, 41(1), 2008.
- [3] P. Blackburn, M. de Rijke, and Y. Venema. *Modal Logic*, volume 53 of *Cambridge Tracts in Theoretical Computer Science*. Cambridge University Press, 2001.
- [4] M. Bojańczyk, A. Muscholl, T. Schwentick, and L. Segoufin. Two-variable logic on data trees and XML reasoning. *Journal of the ACM*, 56(3):1–48, 2009.
- [5] M. Bojańczyk and P. Parys. XPath evaluation in linear time. *Journal of the ACM*, 58(4):17, 2011.
- [6] J. Clark and S. DeRose. XML path language (XPath). Website, 1999. W3C Recommendation. <http://www.w3.org/TR/xpath>.
- [7] A. Dawar and M. Otto. Modal characterisation theorems over special classes of frames. *Annals of Pure and Applied Logic*, 161(1):1–42, 2009.
- [8] D. Figueira. *Reasoning on Words and Trees with Data*. PhD thesis, Laboratoire Spécification et Vérification, ENS Cachan, France, 2010.
- [9] D. Figueira. Decidability of downward XPath. *ACM Transactions on Computational Logic*, 13(4), 2012.
- [10] D. Figueira and L. Segoufin. Bottom-up automata on data trees and vertical XPath. In *International Symposium on Theoretical Aspects of Computer Science (STACS'11)*, volume 9 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 93–104. Leibniz-Zentrum für Informatik, 2011.
- [11] M. Forti and F. Honsell. Set theory with free construction principles. *Annali Scuola Normale Superiore, Pisa*, X(3):493–522, 1983.
- [12] V. Goranko and M. Otto. Model theory of modal logic. In J. Van Benthem P. Blackburn and F. Wolter, editors, *Handbook of Modal Logic*, volume 3 of *Studies in Logic and Practical Reasoning*, chapter 5, pages 249–329. Elsevier, 2007.
- [13] G. Gottlob, C. Koch, and R. Pichler. Efficient algorithms for processing XPath queries. *ACM Transactions on Database Systems*, 30(2):444–491, 2005.
- [14] Marc Gyssens, Jan Paredaens, Dirk Van Gucht, and George H. L. Fletcher. Structural characterizations of the semantics of xpath as navigation tool on a document. In *PODS*, pages 318–327. ACM, 2006.
- [15] D. Harel. Dynamic logic. In D. Gabbay and F. Guenther, editors, *Handbook of Philosophical Logic. Vol. II*, volume 165 of *Synthese Library*, pages 497–604. D. Reidel Publishing Co., Dordrecht, 1984. Extensions of classical logic.
- [16] N. Kurtonina and M. de Rijke. Simulating without negation. *Journal of Logic and Computation*, 7:503–524, 1997.
- [17] M. Marx. XPath with conditional axis relations. In *International Conference on Extending Database Technology (EDBT'04)*, volume 2992 of *LNCS*, pages 477–494. Springer, 2004.
- [18] M. Marx and M. de Rijke. Semantic characterizations of navigational XPath. *SIGMOD Record*, 34(2):41–46, 2005.
- [19] R. Milner. *A Calculus of Communicating Systems*, volume 92 of *LNCS*. Springer, 1980.
- [20] M. Otto. Elementary proof of the Van Benthem-Rosen characterisation theorem. Technical Report 2342, Fachbereich Mathematik, Technische Universität Darmstadt, 2004.
- [21] M. Otto. Modal and guarded characterisation theorems over finite transition systems. *Annals of Pure and Applied Logic*, 130(1-3):173–205, 2004.
- [22] M. Otto. Bisimulation invariance and finite models. In *Logic Colloquium'02*, volume 27 of *Lecture Notes in Logic*, pages 276–298, 2006.
- [23] D. Park. Concurrency and automata on infinite sequences. In *Theoretical Computer Science*, volume 104 of *LNCS*, pages 167–183. Springer, 1981.
- [24] E. Rosen. Modal logic over finite structures. *Journal of Logic, Language and Information*, 6(4):427–439, 1997.
- [25] Davide Sangiorgi. On the origins of bisimulation and coinduction. *ACM Transactions on Programming Languages and Systems*, 31(4), 2009.
- [26] J. van Benthem. *Modal Correspondence Theory*. PhD thesis, Universiteit van Amsterdam, 1976.