



**HAL**  
open science

## Basic Model Theory of XPath on Data Trees

Diego Figueira, Santiago Figueira, Carlos Areces

► **To cite this version:**

Diego Figueira, Santiago Figueira, Carlos Areces. Basic Model Theory of XPath on Data Trees. 17th International Conference on Database Theory (ICDT), Mar 2014, Athens, Greece. hal-01793608

**HAL Id: hal-01793608**

**<https://hal.science/hal-01793608>**

Submitted on 16 May 2018

**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.

# 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”.

## 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 FO<sup>2</sup> (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<sub>=</sub>. Models of XPath<sub>=</sub> 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<sub>=</sub> is to allow formulas of the form  $\langle \alpha = \beta \rangle$ , where  $\alpha, \beta$  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  $\alpha, \beta$ , respectively, and such that the data value of  $y$  is equal to the data value of  $z$ .

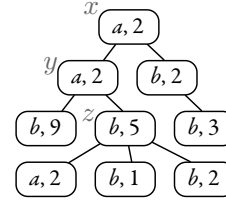
Recent articles investigate several algorithmic problems

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  $\text{XPath}_=^{\downarrow}$  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  $\text{XPath}$  without equality test were studied in [14], and  $\text{XPath}$  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  $\lambda$  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

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

**Table 1: Semantics of XPath<sub>=</sub> 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<sub>=</sub> is a two-sorted language, with **path expressions** (that we write  $\alpha, \beta, \gamma$ ) and **node expressions** (that we write  $\varphi, \psi, \eta$ ). The fragment XPath<sub>=</sub>( $\mathcal{O}$ ), with  $\mathcal{O} \subseteq \{\downarrow, \downarrow_*, \uparrow, \uparrow_*\}$ , is defined by mutual recursion as follows:

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

A **formula** of XPath<sub>=</sub>( $\mathcal{O}$ ) is either a node expression or a path expression. To save space, we use XPath<sub>=</sub> $^\downarrow$  for XPath<sub>=</sub>( $\downarrow$ ); XPath<sub>=</sub> $^\uparrow$  for XPath<sub>=</sub>( $\uparrow$ ); XPath<sub>=</sub> $^{\downarrow*}$  for XPath<sub>=</sub>( $\downarrow, \downarrow_*$ ); and XPath<sub>=</sub> $^{\uparrow*}$  for XPath<sub>=</sub>( $\uparrow, \uparrow_*$ ).

We formally define the semantics of XPath<sub>=</sub> in Table 1. As an example, if  $\mathcal{T}$  is the data tree shown in Figure 1, then  $\llbracket \langle \downarrow_*[b \wedge \downarrow[b] \neq \downarrow[b]] \rangle \rrbracket^T = \{x, y, z\}$ , 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  $u \in \llbracket [\varphi] \rrbracket^T$ , and we say that  $\mathcal{T}, u$  **satisfies**  $\varphi$ . We say that the formulas  $\varphi, \psi$  of XPath<sub>=</sub> are **equivalent** (notation:  $\varphi \equiv \psi$ ) iff  $\llbracket [\varphi] \rrbracket^T = \llbracket [\psi] \rrbracket^T$  for all data trees  $\mathcal{T}$ . Similarly, path expressions  $\alpha, \beta$  of XPath<sub>=</sub> are **equivalent** (notation:  $\alpha \equiv \beta$ ) iff  $\llbracket [\alpha] \rrbracket^T = \llbracket [\beta] \rrbracket^T$  for all data trees  $\mathcal{T}$ .

We call **downward XPath** to XPath<sub>=</sub> $^\downarrow$  and **vertical XPath** to XPath<sub>=</sub> $^\uparrow$ .

In terms of expressive power, it is easy to see that  $\cup$  is unessential: every XPath<sub>=</sub> node expression  $\varphi$  has an equivalent  $\varphi'$  with no  $\cup$  in its path expressions.  $\varphi'$  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  $\cup$

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

where  $\odot \in \{=, \neq\}$ . 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  $\varphi$ , defined in Table 2. Let  $\ell\text{-XPath}_\perp^\downarrow$  be the fragment of XPath<sub>=</sub> $^\downarrow$  consisting of all formulas  $\varphi$  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<sub>=</sub> $^\downarrow$**  (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  **$\ell$ -equivalent for XPath<sub>=</sub> $^\downarrow$**  (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  $\ell$ , there are finitely many different formulas  $\varphi$  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 $\ell$ -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<sub>=</sub> $^\downarrow$**  (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.  $(\xrightarrow{i} v)Z(\xrightarrow{i} v')$  for all  $0 \leq i < n$ , and
  3.  $(\xrightarrow{i} w)Z(\xrightarrow{i} w')$  for all  $0 \leq i < m$ .
- **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:

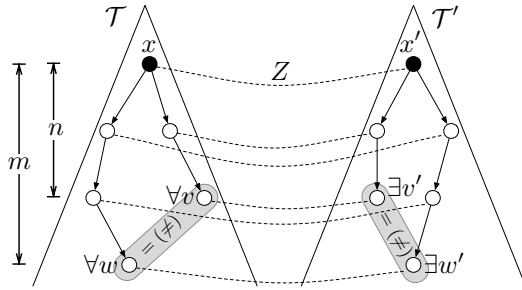
PROPOSITION 3.3.  $\mathcal{T}, u \leftrightarrow^\downarrow (\mathcal{T}|u), u$ .

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  **$\ell$ -bisimilar for XPath<sub>=</sub> $^\downarrow$**  (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

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

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



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

- **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  $\mathcal{T}, u \Leftrightarrow^{\downarrow} \mathcal{T}', u'$  then  $\mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'$  for all  $\ell$ .

**PROPOSITION 3.5.** *Suppose  $\mathcal{T}$  and  $\mathcal{T}'$  have height at most  $\ell$ ,  $u \in T$ , and  $u' \in T'$ . Then  $\mathcal{T}, u \Leftrightarrow_{\ell}^{\downarrow} \mathcal{T}', u'$  iff  $\mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'$ .*

For a data tree  $\mathcal{T}$  and  $u \in T$ , let  $\mathcal{T}|_{\ell}u$  denote the subtree of  $\mathcal{T}$  induced by  $\{v \in T \mid (\exists n \leq \ell) u \xrightarrow{n} v\}$ .

**PROPOSITION 3.6.**  $\mathcal{T}, u \Leftrightarrow_{\ell}^{\downarrow} (\mathcal{T}|_{\ell}u), u$ .

### 3.1.2 Equivalence and bisimulation

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

**THEOREM 3.7.**

1.  $\mathcal{T}, u \Leftrightarrow^{\downarrow} \mathcal{T}', u'$  implies  $\mathcal{T}, u \equiv^{\downarrow} \mathcal{T}', u'$ . 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 (see Appendix for details) is a consequence of the next two propositions:

**PROPOSITION 3.8.**  $\mathcal{T}, u \Leftrightarrow_{\ell}^{\downarrow} \mathcal{T}', u'$  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  $\varphi$  with  $dd(\varphi) \leq j$ , and for all  $\alpha$  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  $\varphi$  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

$$x = x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_n = v \quad \text{in } \mathcal{T},$$

$$x' = x'_0 \rightarrow x'_1 \rightarrow x'_2 \rightarrow \dots \rightarrow x'_n = v' \quad \text{in } \mathcal{T}',$$

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  $\varphi$  of  $\alpha$  and  $k \in \{0, \dots, n\}$  such that  $\mathcal{T}, x_k \models \varphi$  and  $\mathcal{T}', x'_k \not\models \varphi$  (this is shown in Lemma A.1 in the Appendix). This contradicts the inductive hypothesis 1.  $\square$

PROPOSITION 3.9.  $\mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'$  implies  $\mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'$ .

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

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

We show that  $Z$  is an  $\ell$ -bisimulation between  $\mathcal{T}, u$  and  $\mathcal{T}', u'$ . By hypothesis,  $u Z_{\ell} u'$ . 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  $x Z_h x'$ ,

$$\begin{aligned} x &= v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_n = v & \text{in } \mathcal{T}, \\ x &= w_0 \rightarrow w_1 \rightarrow \dots \rightarrow w_m = w & \text{in } \mathcal{T}, \end{aligned}$$

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

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

Since  $\mathcal{T}, x \equiv_h^{\downarrow} \mathcal{T}', x'$ ,  $\text{dd}(\langle \downarrow^n = \downarrow^m \rangle) \leq h$  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  $\text{dd}(\top) = 0$ . For each  $(v', w') \in P$  we define two families of formulas,

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

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

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

Since  $\text{dd}(\varphi_{v', w'}^i) \leq h - i$ , 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  $\Phi^i$  and  $\Psi^j$  are well-formed formulas. Finally, let

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

By construction,  $\text{dd}(\alpha), \text{dd}(\beta) \leq h$  and so  $\text{dd}(\langle \alpha = \beta \rangle) \leq h$ . 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  $\text{XPath}_{=}^{\downarrow}$ . Interestingly, the notion we give is simpler than the one for  $\text{XPath}_{=}^{\downarrow}$  due to a normal form enjoyed by the logic.

In the downward fragment of  $\text{XPath}_{=}^{\downarrow}$  we used  $\text{dd}(\varphi)$  to measure the maximum depth from the current point of evaluation that the formula can access. For the vertical fragment of  $\text{XPath}_{=}^{\downarrow}$ , 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  $\varphi$  (notation:  $\text{vd}(\varphi)$ ) is the pair  $\text{vd}(\varphi) \in \mathbb{Z}_{\geq 0}^2$  defined in Table 2.

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

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

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  $\text{XPath}_{=}^{\downarrow}$**  (notation:  $\mathcal{T}, u \equiv^{\downarrow} \mathcal{T}', u'$ ) iff for all  $\varphi \in \text{XPath}_{=}^{\downarrow}$ , 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  $\text{XPath}_{=}^{\downarrow}$ , 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  $\text{XPath}_{=}^{\downarrow}$  formulas  $\varphi$  so that  $\text{vd}(\varphi) \leq (r, s)$  [resp.  $\text{vd}(\varphi) \leq (r, s)$  and  $\text{nd}(\varphi) \leq k$ ].

### 3.2.1 Normal form

We define a useful normal form for  $\text{XPath}_{=}^{\downarrow}$  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  $\lambda$ ,  $\downarrow^1 = \downarrow$ , and  $\downarrow^{n+1} = \downarrow \downarrow^n$  (similarly for  $\uparrow^n$ ).

A path expression  $\alpha$  of  $\text{XPath}_{=}^{\downarrow}$  is **downward** [resp. **upward**] if it is of the form  $\downarrow^n[\varphi]$  [resp.  $[\varphi]\uparrow^n$ ] for some  $n \geq 0$  with  $\varphi \in \text{XPath}_{=}^{\downarrow}$ . For example,  $\downarrow([\uparrow])$  is a downward expression whereas  $\downarrow([\downarrow])\downarrow$  is not. An **up-down** expression is any expression of the form  $\varepsilon, \alpha^{\uparrow}, \alpha^{\downarrow}$  or  $\alpha^{\uparrow}\alpha^{\downarrow}$  where  $\alpha^{\uparrow}$  is upward and  $\alpha^{\downarrow}$  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  $\alpha^{\updownarrow}, \beta^{\updownarrow}, \gamma^{\updownarrow}$  to denote up-down expressions. Note that in particular any downward or upward expression is an up-down expression. An  $\text{XPath}_{=}^{\downarrow}$  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  $\langle \varepsilon \odot \alpha^{\updownarrow} \rangle$  with  $\odot \in \{=, \neq\}$ .

PROPOSITION 3.10. Let  $\varphi \in (r, s, k)\text{-XPath}_{\perp}^{\downarrow}$ . There is  $\varphi^{\uparrow\downarrow} \in \text{XPath}_{\perp}^{\downarrow}$  in up-down normal form such that

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

### 3.2.2 Finite index

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

PROPOSITION 3.11. If  $r + s \geq 2$  then  $\equiv_{r,s}^{\uparrow\downarrow}$  has infinite index.

In the proof of the above proposition (see Appendix) 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 3.12.  $\equiv_{r,s,k}^{\uparrow\downarrow}$  has finite index.

COROLLARY 3.13.  $\{\mathcal{T}', u' \mid \mathcal{T}, u \equiv_{r,s,k}^{\uparrow\downarrow} \mathcal{T}', u'\}$  is definable by an  $(r, s, k)\text{-XPath}_{\perp}^{\downarrow}$ -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,  $\leftrightarrow_{r,s,k}^{\uparrow\downarrow}$  does not correspond precisely to  $\equiv_{r,s,k}^{\uparrow\downarrow}$ , although  $\leftrightarrow^{\uparrow\downarrow}$  corresponds precisely to  $\equiv^{\uparrow\downarrow}$ . Nonetheless, we obtain, for all  $r, s, k$ ,

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

Let  $\mathcal{T}$  and  $\mathcal{T}'$  be two data-trees. We say that  $u \in T$  and  $u' \in T'$  are **bisimilar for  $\text{XPath}_{\perp}^{\downarrow}$**  (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 3):** 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  $zZz'$ .
- **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 3.10).

We say that  $u \in T$  and  $u' \in T'$  are  $(r, s, k)$ -**bisimilar for  $\text{XPath}_{\perp}^{\downarrow}$**  (notation:  $\mathcal{T}, u \leftrightarrow_{r,s,k}^{\downarrow} \mathcal{T}', u'$ ) 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  $uZ_{\hat{r}, \hat{s}}^{\hat{k}}u'$  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.

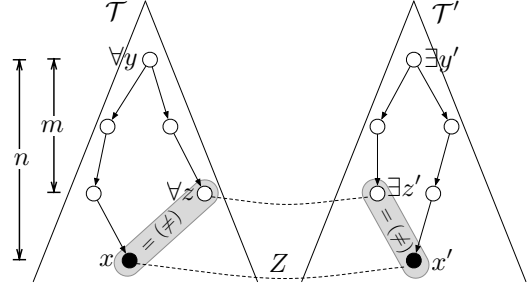


Figure 3: Zig clause of bisimulation for  $\text{XPath}_{\perp}^{\downarrow}$

- **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  $y \xrightarrow{m} z$  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 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 3.14. 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

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

PROPOSITION 3.15.  $\mathcal{T}, u \leftrightarrow_{r,s,k}^{\uparrow\downarrow} (\mathcal{T}|_r^s u), u$ .

### 3.2.4 Equivalence and bisimulation

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

THEOREM 3.16.

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

COROLLARY 3.17.  $\leftrightarrow_{r,s,k}^{\downarrow}$  has finite index.

## 4. SIMULATION

In this section we define notions of directed (non-symmetric) simulations for  $\text{XPath}_{\perp}^{\downarrow}$  and  $\text{XPath}_{\perp}^{\downarrow}$ , as it is done, e.g., in [16] for some modal logics. We obtain results similar to Theorems 3.7 and 3.16 but relating each simulation notion with the corresponding logical implication.

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

A **simulation for  $\text{XPath}_=^\downarrow$**  [resp. for  $\text{XPath}_=^\uparrow$ ] 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  $\text{XPath}_=^\downarrow$  (notation:  $\mathcal{T}, u \xrightarrow{\downarrow} T', u'$ ) 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:** 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.  $\text{data}(v) = \text{data}(w) \Rightarrow \text{data}(v') = \text{data}(w')$ ,
  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  $\text{XPath}_=^\uparrow$  (notation:  $\mathcal{T}, u \xrightarrow{\uparrow} T', u'$ ) 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:** 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  $\text{data}(z) = \text{data}(x)$  then  $\text{data}(z') = \text{data}(x')$ .

Relations  $\xrightarrow{\downarrow}$  and  $\xrightarrow{\uparrow}_{r,s,k}$  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

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

Define  $\Rightarrow^\downarrow$ ,  $\Rightarrow^\uparrow$ , and  $\Rightarrow^\uparrow_{r,s,k}$  in an analogous way for  $\ell\text{-XPath}_=^{\downarrow+}$ ,  $\text{XPath}_=^{\uparrow+}$ ,  $(r, s, k)\text{-XPath}_=^{\uparrow+}$ , respectively. As for bisimulation, we have that  $\Rightarrow$  coincides with  $\Rightarrow$ .

THEOREM 4.1.

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

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

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

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

THEOREM 4.3.

1.  $\varphi \in \text{XPath}_=^\downarrow$  is  $\xrightarrow{\downarrow}$ -invariant [resp.  $\xrightarrow{\downarrow}$ ] iff it is equivalent to a formula of  $\text{XPath}_=^{\downarrow+}$  [resp.  $\ell\text{-XPath}_=^{\downarrow+}$ ].
2.  $\varphi \in \text{XPath}_=^\uparrow$  is  $\xrightarrow{\uparrow}$ -invariant iff it is equivalent to a formula of  $\text{XPath}_=^{\uparrow+}$ .
3. If  $\varphi \in \text{XPath}_=^\uparrow$  is  $\xrightarrow{\uparrow}_{r,s,k}$ -invariant then it is equivalent to a formula of  $(r, s, k)\text{-XPath}_=^{\uparrow+}$ .
4. If  $\varphi \in \text{XPath}_=^\uparrow$  is equivalent to a formula of  $(r, s, k)\text{-XPath}_=^{\uparrow+}$  then  $\varphi$  is  $\xrightarrow{\uparrow}_{r,s,k}$ -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  $\text{XPath}_=^{\downarrow*}$  and  $\text{XPath}_=^{\uparrow*}$  respectively, and let  $\Rightarrow^{\downarrow*}$  and  $\Rightarrow^{\uparrow*}$  be the logical implication for  $\text{XPath}_=^{\downarrow*}$  and  $\text{XPath}_=^{\uparrow*}$  respectively.

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

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

## 6. CHARACTERIZATION

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

### 6.1 Translating to first-order logic

We say that an  $\text{XPath}_=^\downarrow$ -path expression  $\alpha$  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  $n \geq 0$ ,  $\varphi_i \in \text{XPath}_=^\downarrow$ , and  $o_i \in \{\downarrow, \uparrow\}$ .

PROPOSITION 6.1. For any  $\text{XPath}_=^\downarrow$ - [resp.  $\text{XPath}_=^\uparrow$ -] path expression  $\alpha$  there is an equivalent  $\text{XPath}_=^\downarrow$ - [resp.  $\text{XPath}_=^\uparrow$ -] path expression  $\alpha'$  in simple normal form. Further,  $\alpha'$  can be computed in polynomial time from  $\alpha$ .<sup>1</sup>

We say that an  $\text{XPath}_=^\downarrow$ -formula  $\varphi$  is in **simple normal form** if each path expression  $\alpha$  occurring in  $\varphi$  is in simple normal form.

Fix the signature  $\sigma$  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  $\sigma$ -structure such that

$$\rightsquigarrow^\mathcal{T} = \{(x, y) \in T^2 \mid y \text{ is a child of } x\};$$

<sup>1</sup>Note that this proposition holds only for paths expressions without union.



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

We define the following translation  $\text{Tr}$  mapping  $\text{XPath}_{\perp}^{\dagger}$  formulas in simple normal form to first-order  $\sigma$ -formulas:

$$\begin{aligned}\text{Tr}_x(a) &= P_a(x) & (a \in \mathbb{A}) \\ \text{Tr}_x(\varphi \uparrow \psi) &= \text{Tr}_x(\varphi) \uparrow \text{Tr}_x(\psi) & (\uparrow \in \{\wedge, \vee\}) \\ \text{Tr}_x(\neg\varphi) &= \neg\text{Tr}_x(\varphi) \\ \text{Tr}_x(\langle\alpha\rangle) &= (\exists \bar{y})(x = y_0 \wedge \text{Tr}_{\bar{y}}(\alpha)) \\ \text{Tr}_x(\langle\alpha = \beta\rangle) &= (\exists \bar{y})(\exists \bar{z})(x = y_0 \wedge x = z_0 \wedge y_n \approx z_m \wedge \\ &\quad \text{Tr}_{\bar{y}}(\alpha) \wedge \text{Tr}_{\bar{z}}(\beta)) \\ \text{Tr}_x(\langle\alpha \neq \beta\rangle) &= (\exists \bar{y})(\exists \bar{z})(x = y_0 \wedge x = z_0 \wedge y_n \not\approx z_m \wedge \\ &\quad \text{Tr}_{\bar{y}}(\alpha) \wedge \text{Tr}_{\bar{z}}(\beta)) \\ \text{Tr}_{\bar{y}}(\alpha) &= \bigwedge_{i=0}^{n-1} o_{i+1}(y_i, y_{i+1}) \wedge \bigwedge_{i=0}^n \text{Tr}_{y_i}(\varphi_i),\end{aligned}$$

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;

$$\begin{aligned}\alpha &= [\varphi_0]o_1[\varphi_1]o_2[\varphi_2]o_3 \cdots o_n[\varphi_n]; \\ \beta &= [\psi_0]o'_1[\psi_1]o'_2[\psi_2]o'_3 \cdots o'_m[\psi_m];\end{aligned}$$

$o_i, o'_i \in \{\downarrow, \uparrow\}$ ;  $o_j(u, v)$  represents  $u \rightsquigarrow v$  if  $o_j = \downarrow$ , and  $v \rightsquigarrow u$  otherwise.

**PROPOSITION 6.2.** *For  $\varphi \in \text{XPath}_{\perp}^{\dagger}$  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  $\sigma$ , and let  $\mathcal{C}$  be a class of  $\sigma$ -models. An  $\text{FO}(\sigma)$ -formula  $\varphi(x)$  is  $\ell$ -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.

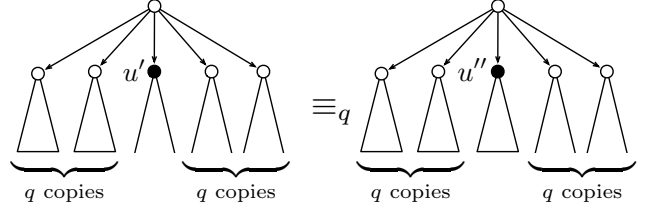
**THEOREM 6.3 (CHARACTERIZATION).** *Let  $\varphi(x) \in \text{FO}(\sigma)$ . The following are equivalent:*

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

**PROOF.** The implication (ii)  $\Rightarrow$  (i) follows straightforwardly from Theorem 3.7. The proof of (i)  $\Rightarrow$  (ii) goes as follows: First, we show that any  $\Leftrightarrow^{\downarrow}$ -invariant  $\varphi(x) \in \text{FO}(\sigma)$  is  $\ell$ -local for  $\ell = 2^{\text{qr}(\varphi)} - 1$  (Proposition 6.4). Then, we prove that any  $\Leftrightarrow^{\downarrow}$ -invariant  $\varphi(x) \in \text{FO}(\sigma)$  that is  $\ell$ -local is  $\Leftrightarrow^{\downarrow}_{\ell}$ -invariant (Proposition B.2 in the Appendix). Finally, we show that any  $\text{FO}(\sigma)$ -definable property which is  $\Leftrightarrow^{\downarrow}_{\ell}$ -invariant is definable in  $\ell$ - $\text{XPath}_{\perp}^{\dagger}$  (Proposition B.3 in the Appendix).  $\square$

**PROPOSITION 6.4.** *Any  $\Leftrightarrow^{\downarrow}$ -invariant  $\varphi(x) \in \text{FO}(\sigma)$  over [finite] data-trees is  $\ell$ -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



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

- (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

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

and hence  $\varphi$  is  $\ell$ -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 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 6.3 fails for  $\text{XPath}_{\perp}^{\dagger}$ :

**LEMMA 6.5.** *The  $\text{FO}(\sigma)$ -formula*

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

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

Hence  $\text{XPath}_{\perp}^{\dagger}$  is not the fragment of  $\text{FO}(\sigma)$  which is  $\Leftrightarrow^{\downarrow}$ -invariant over [finite] data-trees. However, the following analog of Proposition B.3 (needed for the proof of Theorem 6.3) still holds for the case of  $\text{XPath}_{\perp}^{\dagger}$ :

**PROPOSITION 6.6.** *Let  $k' = k \cdot (r+s+2)$ . If  $\varphi(x) \in \text{FO}(\sigma)$  is  $\Leftrightarrow^{\downarrow}_{r,s,k'}$ -invariant over [finite] data-trees, then there is  $\psi \in (r, s, k)$ - $\text{XPath}_{\perp}^{\dagger}$  such that  $\text{Tr}_x(\psi)$  is logically equivalent to  $\varphi$  over [finite] data-trees.*

Notice that the counterexample in Lemma 6.5 is an unrestricted, existential formula. One may wonder if it might be possible to extend the expressive power of  $\text{XPath}_{\perp}^{\downarrow}$  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 6.3 fails. Formally, consider the logic  $\text{XPath}_{\perp}^{\downarrow E}$ , which results from adding the operator  $E$  to  $\text{XPath}_{\perp}^{\downarrow}$  with the following semantics:  $\llbracket E\varphi \rrbracket^T = T$  if  $\llbracket \varphi \rrbracket^T \neq \emptyset$ , and  $\llbracket E\varphi \rrbracket^T = \emptyset$  otherwise.

The following lemma shows a counterexample to the analog of Theorem 6.3, showing that  $\text{XPath}_{\perp}^{\downarrow E}$  is not the fragment of  $\text{FO}(\sigma) \leftarrow^{\downarrow}$ -invariant over [finite] data-trees.

LEMMA 6.7. *The  $\text{FO}(\sigma)$ -formula*

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

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

## 7. APPLICATIONS

We devote this section to exemplify how the model theoretic tools we developed can be used to show expressiveness results for  $\text{XPath}_{\perp}$ . 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  $\varphi \in \text{XPath}_{\perp}^{\downarrow}$  such that  $\text{dd}(\varphi) \leq \ell$  and  $\text{nd}(\varphi) \leq k$  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  $(\ell, k)$ -bisimilar for  $\text{XPath}_{\perp}^{\downarrow}$  (notation:  $\mathcal{T}, u \leftarrow_{\ell, k}^{\downarrow} \mathcal{T}', u'$ ) if there is a family of relations  $(Z_{j,t})_{j \leq \ell, t \leq k}$  in  $T \times T'$  such that  $uZ_{\ell, k}u'$  and for all  $j \leq \ell, t \leq k, x \in T$  and  $x' \in T'$  we have

- **Harmony:** If  $xZ_{j,t}x'$  then  $\text{label}(x) = \text{label}(x')$ .
- **Zig:** If  $xZ_{j,t}x', x' \xrightarrow{n} v$  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.  $\text{data}(v) = \text{data}(w) \Leftrightarrow \text{data}(v') = \text{data}(w')$ ,
  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  $xZ_{j,t}x', x' \xrightarrow{n} v'$  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 3.8 and 3.9, it is easy to show that  $(\ell, k)$ -bisimulations characterize  $(\ell, k)$ -equivalence.

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

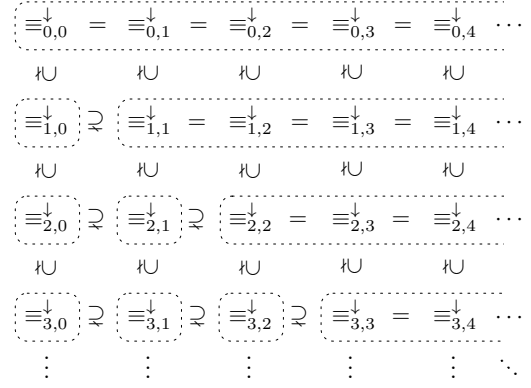


Figure 5: Hierarchy of  $\text{XPath}_{\perp}^{\downarrow}$ .

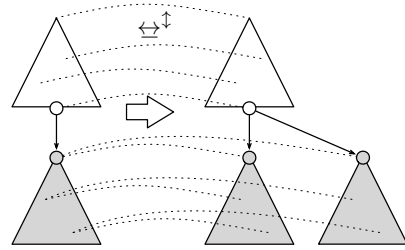


Figure 6: Closure under subtree replication.

The following theorem —proved in the Appendix using the bisimulation notion introduced above— 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  $\ell, k \geq 0, i \geq 1$ ,*

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

### 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 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.  $\text{XPath}_{\perp}^{\downarrow \uparrow \uparrow}$  is closed under subtree replication, i.e. if  $\mathcal{T}'$  is a subtree replication of  $\mathcal{T}$ , and  $u \in T$  then  $\mathcal{T}', u \equiv^{\uparrow \uparrow \uparrow} \mathcal{T}, u$ .

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  $\mathcal{T}, u \leftrightarrow^\downarrow \mathcal{T}', u$  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.

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

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  $XPath_{\downarrow}^{\uparrow\downarrow*}$  because the models  $\mathcal{T}$  and  $\mathcal{T}'$  in Figure 8 are bisimilar for  $XPath_{\downarrow}$  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  $XPath_{\downarrow}^{\uparrow\downarrow*}$ .

Finally, suppose there exists  $\psi \in XPath_{\downarrow}^{\uparrow\downarrow*+}$  expressing  $fk(a, b)$ . Since  $\mathcal{T}$  is a substructure of  $\mathcal{T}'$  we have  $\mathcal{T}, x \rightarrow^\downarrow \mathcal{T}', x$  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  $XPath_{\downarrow}^{\uparrow}$ ;
2.  $dist_3$  is not expressible in  $XPath_{\downarrow}^{\uparrow\downarrow*}$ ;
3. neither  $dist_3$  nor its complement can be expressed in  $XPath_{\downarrow}^{\uparrow\downarrow*+}$ .

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

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  $\mathcal{T}, x \leftrightarrow^\downarrow \mathcal{T}', x'$  via  $Z \subseteq T \times T'$  defined by

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

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  $XPath_{\downarrow}^{\uparrow\downarrow*}$ .

For 3, one can verify that  $\mathcal{T}, x \rightarrow^\downarrow \mathcal{T}', x'$  via  $Z$  as defined above. If  $dist_3$  were definable in  $XPath_{\downarrow}^{\uparrow\downarrow*+}$  via  $\psi$  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  $XPath_{\downarrow}^{\uparrow\downarrow*+}$ .  $\square$

## 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. *J. of the ACM*, 55(2):1–79, 2008.

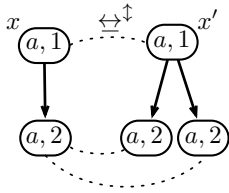


Figure 7:  $key(a)$  not in  $XPath_{\downarrow}^{\uparrow\downarrow*}$ .

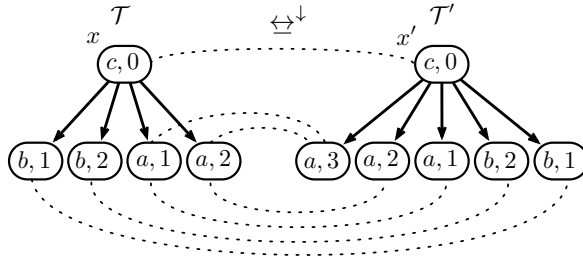


Figure 8:  $fk(a, b)$  not in  $XPath_{\downarrow}^{\uparrow\downarrow*}$ .

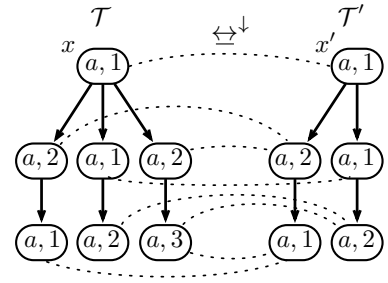


Figure 9:  $dist_3$  not in  $XPath_{\downarrow}^{\uparrow\downarrow*}$ .

- [2] M. Benedikt and C. Koch. XPath leashed. *ACM Comput. Surv.*, 41(1), 2008.
- [3] P. Blackburn, M. de Rijke, and Y. Venema. *Modal Logic*, volume 53 of *Cambridge Tracts Theoret. Comput. Sci.* Cambridge University Press, 2001.
- [4] M. Bojańczyk, A. Muscholl, T. Schwentick, and L. Segoufin. Two-variable logic on data trees and XML reasoning. *J. of the ACM*, 56(3):1–48, 2009.
- [5] M. Bojańczyk and P. Parys. XPath evaluation in linear time. *J. 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. *Ann. Pure Appl. 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 Trans. Comput. Log.*, 13(4), 2012.
- [10] D. Figueira and L. Segoufin. Bottom-up automata on data trees and vertical XPath. In *Int. Symp. on Theoretical Aspects of Computer Science (STACS'11)*, volume 9 of *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 Trans. 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. *J. Logic Comput.*, 7:503–524, 1997.
- [17] M. Marx. XPath with conditional axis relations. In *Int. Conf. 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. *Ann. Pure Appl. Logic*, 130(1-3):173–205, 2004.
- [22] M. Otto. Bisimulation invariance and finite models. In *Logic Colloquium'02*, volume 27 of *Lect. Notes Log.*, pages 276–298, 2006.
- [23] D. Park. Concurrency and automata on infinite sequences. In *Theoret. Comput. Sci.*, volume 104 of *LNCS*, pages 167–183. Springer, 1981.
- [24] E. Rosen. Modal logic over finite structures. *J. Logic Lang. Inform.*, 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.

## APPENDIX

### A. PROOFS OF SECTION 3

Given a path expression  $\alpha$ , the **navigation** of  $\alpha$  (notation:  $\text{nav}(\alpha)$ ) is the string of  $\{\uparrow, \downarrow\}^*$  that results from removing all node expressions  $[\psi]$  and  $\varepsilon$  from  $\alpha$ . For example,  $\text{nav}(\downarrow[\uparrow]\downarrow[\downarrow = \uparrow]\uparrow[b]) = \downarrow\uparrow$ .

LEMMA A.1. *Let  $\alpha$  be a path expression of  $XPath_{\downarrow}^{\uparrow\downarrow*}$ . Let  $x \xrightarrow{n} v$  and  $x' \xrightarrow{n} v'$  such that  $(x, v) \in \llbracket \alpha \rrbracket^T$  and  $(x', v') \notin \llbracket \alpha \rrbracket^{T'}$ .*

Then there is a subformula  $\varphi$  of  $\alpha$  and  $k \in \{0, \dots, n\}$  such that  $\mathcal{T}, (\overset{k}{\rightarrow}v) \models \varphi$  and  $\mathcal{T}', (\overset{k}{\rightarrow}v') \not\models \varphi$ .

PROOF. Let  $x = v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_n = v$  and  $x' = v'_0 \rightarrow v'_1 \rightarrow \dots \rightarrow v'_n \neq \varphi^{\uparrow\downarrow} \equiv \varphi$ ;  $v'$ . We proceed by induction on  $|\alpha|$ . If  $\alpha = \varepsilon$  then  $x = v$  and so  $n = 0$ . Hence  $x' = v'$  and  $(x', v') \in \llbracket \alpha \rrbracket^{\mathcal{T}'}$ , which contradicts the hypothesis, and thus the statement is trivially true. If  $\alpha = \downarrow$  then  $x \rightarrow v$  and so  $n = 1$ . Hence  $x' \rightarrow v'$  and  $(x', v') \in \llbracket \alpha \rrbracket^{\mathcal{T}'}$ . This case is also trivial. The case  $\alpha = \downarrow_*$  is similar.

Suppose  $\alpha = [\psi]$ . Since  $(x', v') \notin \llbracket \alpha \rrbracket^{\mathcal{T}'}$ , we have  $x' = v'$  and  $\mathcal{T}', v' \not\models \psi$ . Taking  $k = 0$  and  $\varphi = \psi$  the statement holds. Observe that  $\psi$  is a subformula of  $\alpha$ .

Suppose  $\alpha = \beta\gamma$ . Then there is  $k$  such that  $(x, v_k) \in \llbracket \beta \rrbracket^{\mathcal{T}}$  and  $(v_k, v) \in \llbracket \gamma \rrbracket^{\mathcal{T}}$ . Since  $(x', v') \notin \llbracket \alpha \rrbracket^{\mathcal{T}'}$ , we have  $(x', v'_k) \notin \llbracket \beta \rrbracket^{\mathcal{T}'}$  or  $(v'_k, v') \notin \llbracket \gamma \rrbracket^{\mathcal{T}'}$ . In either case, apply inductive hypothesis straightforwardly.  $\square$

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

PROOF. We show by induction on  $\ell$  that there are finitely many non-equivalent formulas of downward depth at most  $\ell$ , and finitely many non-equivalent path expressions of downward depth at most  $\ell$ .

For the base case, any formula of downward depth 0 is a Boolean combination of labels, and hence there are finitely many non-equivalent of them. Any path expression of downward depth 0 is equivalent to  $[\varphi]$  for  $\text{dd}(\varphi) = 0$ , and hence there are finitely many non-equivalent of them.

For the induction, any formula of downward depth  $\ell + 1$  is a boolean combination of labels or formulas of the form  $\langle \alpha \rangle$ ,  $\langle \alpha = \beta \rangle$  or  $\langle \alpha \neq \beta \rangle$ , where  $\text{dd}(\alpha), \text{dd}(\beta) \leq \ell + 1$ , so it suffices to show that there are finitely many non-equivalent path expressions of downward depth at most  $\ell + 1$ . Let  $\alpha$  be such that  $\text{dd}(\alpha) \leq \ell + 1$ . By Proposition 6.1,  $\alpha$  is either equivalent to a path expression of the form  $[\psi]$  or of the form  $[\psi] \downarrow \beta$ , where  $\text{dd}(\psi), \text{dd}(\beta) \leq \ell$ . By inductive hypothesis there are finitely many non-equivalent  $\psi$ 's and  $\beta_i$ 's, and hence finitely many non-equivalent  $\alpha$ 's.  $\square$

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

PROOF. Consider the conjunction of all  $\ell$ -XPath $_{\downarrow}^{\downarrow}$  formulas  $\varphi$  such that  $\mathcal{T}, u \models \varphi$ . By Proposition 3.1, up to logical equivalence, there are finitely many such  $\varphi$ 's, and hence the conjunction is equivalent to a finite one. Define  $\chi_{\ell, \mathcal{T}, u}$  as this finite conjunction.  $\square$

THEOREM 3.7.

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

PROOF. Item 2 is a direct consequence of Propositions 3.8 and 3.9. The argument for item 1 is similar to the one of the aforementioned propositions, but working with a single  $Z$  instead of  $(Z_i)_{i \leq \ell}$ . For the converse implication, define  $Z$  by  $xZx'$  iff  $\mathcal{T}, x \equiv^{\downarrow} \mathcal{T}', x'$ . The conjunctions in (1) are then finite because  $\mathcal{T}'$  is finitely branching, and so  $P$  is finite (the fact that  $\mathcal{T}$  is finitely branching is used to show that  $Z$  is satisfied).  $\square$

PROPOSITION 3.10. Let  $\varphi \in (r, s, k)$ -XPath $_{\downarrow}^{\downarrow}$ . There is  $\varphi^{\uparrow\downarrow} \in \text{XPath}_{\downarrow}^{\downarrow}$  in up-down normal form such that

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

PROOF. The idea is that we can factorize any path in the tree going down and up as a node tests in the expression. Consider for instance the expression  $\alpha = \uparrow\downarrow[a]\uparrow\downarrow$ . It is immediate that  $\alpha$  is equivalent to the up-down expression  $[(\uparrow[\downarrow[a]])]\uparrow\downarrow$ , which is in up-down normal form.

We use the following directed equivalences to translate any path expression into an equivalent up-down expression.

$$\begin{aligned} \varepsilon\gamma &\equiv^{\uparrow} \gamma && (\varepsilon) \\ \alpha[\psi_1][\psi_2]\beta &\equiv^{\uparrow} \alpha[\psi_1 \wedge \psi_2]\beta && (\text{merge}) \\ \alpha \xi_{-n} \downarrow \dots \downarrow \xi_{-1} \downarrow \xi_0 \uparrow \xi_1 \uparrow \dots \uparrow \xi_n \beta &\equiv^{\uparrow} \\ &\alpha [(\xi_{-n} \downarrow \dots \downarrow \xi_{-1} \downarrow \xi_0)] \beta && (\text{factor}) \\ \alpha \xi_n \downarrow \xi_{n-1} \downarrow \dots \downarrow \xi_0 &\equiv^{\uparrow} \alpha \downarrow^n [(\xi_0 \uparrow \xi_1 \uparrow \dots \uparrow \xi_n)] && (\text{shift-r}) \\ \xi_0 \uparrow \xi_1 \uparrow \dots \uparrow \xi_n \beta &\equiv^{\uparrow} [(\xi_0 \uparrow \xi_1 \uparrow \dots \uparrow \xi_n)] \uparrow^n \beta && (\text{shift-l}) \end{aligned}$$

In the expressions above, each  $\xi_i$  is the empty string, or of the form  $\varepsilon$  or  $[\varphi_1][\varphi_2] \dots [\varphi_n]$ ,  $\alpha$  and  $\beta$  can be any path expression, or the empty string, and  $\gamma$  is any path expression. The idea is that (*factor*) converts an expression that goes down  $n$  times and then up  $n$  times into a node expression, and when doing this, any test done in the  $i$ -th node when going down is merged with the  $(n - i)$ -th test when going up. For example,  $\downarrow[\neg a]\downarrow[c]\uparrow[\neg b]\uparrow \equiv^{\uparrow} [(\downarrow[\neg a][\neg b]\downarrow[c])]$ . On the other hand, (*shift-r*) and (*shift-l*) group all the node tests in the lowest node in the expression, making use of the fact that the parent relation is functional. Thus, for example  $[a]\downarrow[b]\downarrow \equiv^{\uparrow} \downarrow\downarrow[(\uparrow[b]\uparrow[a])]$  and  $\uparrow[a]\uparrow[b] \equiv^{\uparrow} [(\uparrow[a]\uparrow[b])]\uparrow\uparrow$ . It is thus clear that the left and right expressions above are semantically equivalent.

LEMMA A.2. Let  $\alpha$  be an XPath $_{\downarrow}^{\downarrow}$ -path expression with  $\text{vd}(\alpha) = (r, s)$  and  $\text{nd}(\alpha) = k$ , Then there is an up-down path expression  $\alpha^{\uparrow\downarrow}$  such that:

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

PROOF. We first apply rule (*factor*) as many times as possible. It is clear that if  $\text{nav}(\alpha)$  is of the form  $\uparrow^n \downarrow^m$  for some  $n, m \geq 0$  then rule (*factor*) cannot be applied and we are done. Hence, suppose  $\text{nav}(\alpha)$  contains the pattern  $\uparrow\downarrow$ . Let

$$\begin{aligned} \alpha &= \gamma \uparrow \alpha_1 \gamma \downarrow \\ \alpha_1 &= \gamma_1 \underbrace{\xi_{-n_1}^1 \downarrow \dots \downarrow \xi_0^1 \uparrow \dots \uparrow \xi_{n_1}^1}_{\text{matches (factor)}} \\ &\quad \gamma_2 \underbrace{\xi_{-n_2}^2 \downarrow \dots \downarrow \xi_0^2 \uparrow \dots \uparrow \xi_{n_2}^2}_{\text{matches (factor)}} \\ &\quad \vdots \end{aligned}$$

$$\gamma_{m-1} \underbrace{\xi_{-n_m}^m \downarrow \dots \downarrow \xi_0^m \uparrow \dots \uparrow \xi_{n_m}^m}_{\text{matches (factor)}} \gamma_m,$$

where  $\text{nav}(\gamma_\uparrow), \text{nav}(\gamma_m) \in \uparrow^*$ ,  $\text{nav}(\gamma_\downarrow), \text{nav}(\gamma_1) \in \downarrow^*$ , and  $\xi_j^i$  are the empty string,  $\varepsilon$  or  $[\varphi_1^{i,j}][\varphi_2^{i,j}] \dots [\varphi_{h_{i,j}}^{i,j}]$ . Furthermore, assume that  $m$  is maximal (i.e., it is impossible to apply (factor) in any of the  $\gamma_i$ 's) and that the length of each  $\gamma_i$  is minimal (i.e., it is not the case that  $\text{nav}(\gamma_i)$  ends with  $\downarrow$  and that  $\text{nav}(\gamma_{i+1})$  begins with  $\uparrow$ ). Observe that  $\text{nav}(\gamma_i) \in \uparrow^* \downarrow^*$ . We apply rule (factor) in the  $m-1$  marked places and obtain

$$\begin{aligned} \alpha_2 &= \gamma_1 \underbrace{[\xi_{-n_1}^1 \xi_{n_1}^1 \downarrow \dots \downarrow \xi_{-1}^1 \xi_1^1 \downarrow \xi_0^1]}_{\text{(factor) applied}} \\ &\quad \gamma_2 \underbrace{[\xi_{-n_2}^2 \xi_{n_2}^2 \downarrow \dots \downarrow \xi_{-1}^2 \xi_1^2 \downarrow \xi_0^2]}_{\text{(factor) applied}} \\ &\quad \vdots \\ &\quad \gamma_{m-1} \underbrace{[\xi_{-n_m}^m \xi_{n_m}^m \downarrow \dots \downarrow \xi_{-1}^m \xi_1^m \downarrow \xi_0^m]}_{\text{(factor) applied}} \gamma_m, \end{aligned}$$

Let  $\text{vd}(\text{nav}(\alpha_1)) = (r_1, s_1)$ . Since  $\text{nav}(\alpha) = \text{nav}(\gamma_\uparrow \alpha_1 \gamma_\downarrow)$  contains the pattern  $\downarrow \uparrow$ , we have that  $r_1 > 0$ . It can be shown that  $\text{vd}(\gamma_\uparrow \alpha_2 \gamma_\uparrow) = (r, s)$ ,  $\text{nd}(\alpha_2) \leq \text{nd}(\alpha_1) + 1$ , and  $\text{vd}(\text{nav}(\alpha_2)) \leq (r_1 - 1, s_1)$ . If we repeat this procedure with  $\alpha_2$  and so on until we can no longer apply rule (factor), we end up with an up-down path expression  $\alpha_f$  so that

1.  $\alpha_f \equiv^\dagger \alpha_1$ ,
2.  $\text{vd}(\gamma_\uparrow \alpha_f \gamma_\downarrow) = (r, s)$ , and
3.  $\text{nd}(\alpha_f) \leq \text{nd}(\alpha_1) + r_1$ .

We can now apply (shift-r), (shift-l), ( $\varepsilon$ ) and (merge) to  $\gamma_\uparrow \alpha_f \gamma_\downarrow$  in order to obtain an up-down path expression  $\alpha^\updownarrow$  satisfying the desired conditions:

1.  $\alpha^\updownarrow \equiv^\dagger \alpha$
2.  $\text{vd}(\alpha^\updownarrow) = (r, s)$ , and
3.  $\text{nd}(\alpha^\updownarrow) \leq k + r_1 + 1 \leq r + s + 1$ .

This concludes the proof of Lemma A.2.  $\square$

LEMMA A.3. Let  $\alpha^\updownarrow, \beta^\updownarrow$  be up-down path expressions and let  $\varphi = \langle \alpha^\updownarrow \odot \beta^\updownarrow \rangle$  (for  $\odot \in \{=, \neq\}$ ) with  $\text{vd}(\varphi) = (r, s)$  and  $\text{nd}(\varphi) = k$ . Then there is an up-down path expression  $\gamma^\updownarrow$  such that:

1.  $\langle \varepsilon \odot \gamma^\updownarrow \rangle \equiv^\dagger \varphi$ ,
2.  $\text{vd}(\gamma^\updownarrow) = (r, s)$ , and
3.  $\text{nd}(\gamma^\updownarrow) \leq k + 1$ .

PROOF. Let us analyse the case where

$$\begin{aligned} \alpha^\updownarrow &= [\psi_\alpha] \uparrow^{n_\alpha} \downarrow^{m_\alpha} [\tau_\alpha] \\ \beta^\updownarrow &= [\psi_\beta] \uparrow^{n_\beta} \downarrow^{m_\beta} [\tau_\beta] \end{aligned}$$

(the remaining cases being only simpler), where  $\psi_\alpha, \psi_\beta, \tau_\alpha, \tau_\beta$  are in up-down normal form. Suppose, without loss of generality, that  $n_\alpha \leq n_\beta$ . Hence, we have  $\langle \alpha^\updownarrow \odot \beta^\updownarrow \rangle \equiv^\dagger$

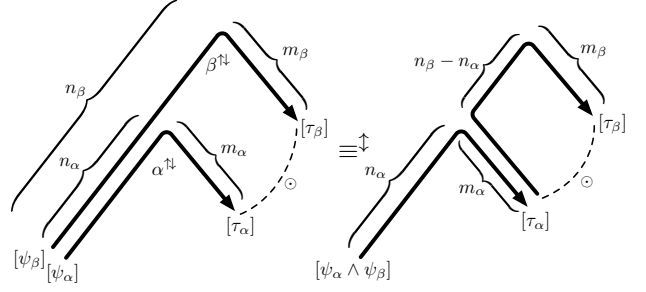


Figure 10: Normal form for data tests.

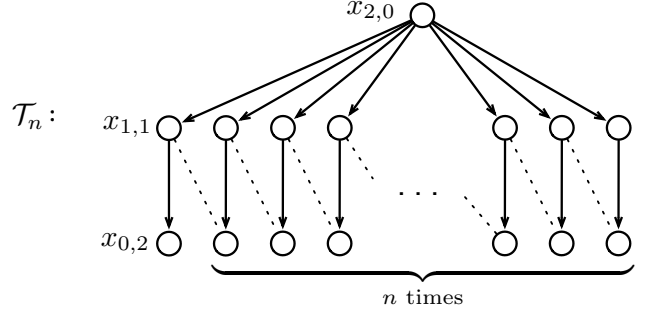


Figure 11: Model verifying  $\psi_i^j$  for all  $i \geq n$  and not verifying  $\varphi_l$  for no  $l < n$ . Dotted lines represent equal data values.

$\langle \varepsilon \odot \gamma^\updownarrow \rangle$ , where

$$\gamma^\updownarrow = [\psi_\alpha \wedge \psi_\beta] \uparrow^{n_\alpha} \downarrow^{m_\alpha} [\tau_\alpha \wedge \langle \varepsilon \odot \uparrow^{m_\alpha} \uparrow^{n_\beta - n_\alpha} \downarrow^{m_\beta} [\tau_\beta] \rangle].$$

It is clear that the formulas are equivalent (cf. Figure 10). Moreover, the right-hand formula has at most one more nesting than the left-hand formula, and its vertical depth is at most  $(r, s)$ . This concludes the proof of Lemma A.3.  $\square$

By induction on  $\varphi$ , and using lemmas A.2 and A.3, one can show that there is  $\varphi^\updownarrow$  as desired.

PROPOSITION 3.11. If  $r + s \geq 2$  then  $\equiv_{r,s}^\dagger$  has infinite index.

PROOF. We show that for every  $r, s$  so that  $r + s = 2$  there is an infinite set of non-equivalent formulas  $\{\psi_{r,s}^i\}_{i \geq 0}$  of vertical depth  $(r, s)$ . It thus follows that for every  $r, s$  so that  $r + s \geq 2$ ,  $\equiv_{r,s}^\dagger$  has infinite index.

Consider the following formulas.

$$\begin{aligned} \psi_{1,1}^0 &= \langle \varepsilon = \uparrow \downarrow \downarrow \rangle & \psi_{1,1}^{i+1} &= \langle \varepsilon = \uparrow \downarrow [\psi_{1,1}^i] \downarrow \rangle \\ \psi_{0,2}^0 &= \langle \uparrow = \uparrow \uparrow \downarrow \downarrow \rangle & \psi_{0,2}^{i+1} &= \langle \uparrow = \uparrow \uparrow \downarrow [\psi_{1,1}^i] \downarrow \rangle \\ \psi_{2,0}^0 &= \langle \downarrow = \downarrow \downarrow \rangle & \psi_{2,0}^{i+1} &= \langle \downarrow = \downarrow [\psi_{1,1}^i] \downarrow \rangle \end{aligned}$$

Note that  $\text{vd}(\psi_{r,s}^n) = (r, s)$  and  $\text{nd}(\psi_{r,s}^n) = n$  for every  $n$ . The formula  $\psi_{r,s}^n$  intuitively says that there is a chain of length  $n$  as depicted in Figure 11.

In the data tree  $\mathcal{T}_n$  of the figure, we have that  $\mathcal{T}_n, x_{r,s} \models \psi_{r,s}^n$  but  $\mathcal{T}_n, x_{r,s} \not\models \psi_{r,s}^{n'}$  for any  $n' > n$ . Therefore,  $\{\psi_{r,s}^i\}_{i \geq 0}$  is an infinite set of non-equivalent formulas of vertical depth  $(r, s)$ .  $\square$

PROPOSITION 3.12.  $\equiv_{r,s,k}^{\uparrow}$  has finite index.

PROOF. For any  $\varphi$  with  $\text{nd}(\varphi) = k$  and  $\text{vd}(\varphi) = (r, s)$ , let  $F(\varphi) = (k, r + s)$ . Define  $F$  in a similar way for path expressions  $\alpha$ . In this proof “finite” means finite up to logical equivalence. By Proposition 3.10 we can consider only formulas in up-down normal form.

We show that there are finitely many formulas  $\varphi$  in up-down normal form such that  $F(\varphi) \leq (k, t)$ , and that there are finitely many path expressions  $\alpha$  in up-down normal form such that  $F(\alpha) \leq (k, t)$ . We use induction on the lexicographic ordering of  $(k, t)$ . Observe that if  $F(\varphi) = (k, t)$  then  $\varphi$  is a boolean combination of labels and formulas of the form  $\langle \varepsilon = \alpha \rangle$ ,  $\langle \varepsilon \neq \alpha \rangle$  or  $\langle \alpha \rangle$ , where  $F(\alpha) \leq (k, t)$ . Hence it suffices to show the statement for path expressions. If  $F(\alpha) = (0, t)$  then  $\alpha$  is either  $\varepsilon$  or  $\uparrow^n \downarrow^m$ , where  $n, m \leq 2t$ , so there are finitely many of them. If  $F(\alpha) = (k + 1, t)$ , then  $\alpha$  is  $[\varphi_1] \uparrow^n \downarrow^m [\varphi_2]$ , where  $n, m \leq 2t$  and  $\text{nd}(\varphi_i) \leq k$  for  $i = 1, 2$ . Since  $F(\varphi_i) <_{\text{lex}} (k + 1, t)$ , by inductive hypothesis there are finitely many such  $\varphi_i$ 's, and therefore  $\alpha$ 's.  $\square$

COROLLARY 3.13.  $\{\mathcal{T}', u' \mid \mathcal{T}, u \equiv_{r,s,k}^{\uparrow} \mathcal{T}', u'\}$  is definable by an  $(r, s, k)$ -XPath $_{\equiv}^{\uparrow}$ -formula.

PROOF. Similar to the proof of Corollary 3.2.  $\square$

THEOREM 3.16.

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

PROOF. Items 2 and 3 are shown in Propositions A.4 and A.5.

The argument for item 1 is similar to the one of the aforementioned propositions, but working with a single  $Z$  instead of  $(Z_{\hat{r}, \hat{s}}^{\hat{k}})_{\hat{r}, \hat{s}, \hat{k}}$ . For the converse implication, define  $Z$  by  $xZx'$  iff  $\mathcal{T}, x \equiv^{\uparrow} \mathcal{T}', x'$ . The conjunctions in (2) are then finite because  $\mathcal{T}'$  is finitely branching, and so  $P$  is finite (the fact that  $\mathcal{T}'$  is finitely branching is used for showing that Zag is satisfied).  $\square$

PROPOSITION A.4.  $\mathcal{T}, u \leftrightarrow_{r,s,k \cdot (r+s+2)}^{\uparrow} \mathcal{T}', u'$  implies  $\mathcal{T}, u \equiv_{r,s,k}^{\uparrow} \mathcal{T}', u'$ .

PROOF. We show that if  $\mathcal{T}, u \leftrightarrow_{r,s,k}^{\uparrow} \mathcal{T}', u'$  via

$$(Z_{\hat{r}, \hat{s}}^{\hat{k}})_{\hat{r} + \hat{s} \leq r + s, \hat{k} \leq k}$$

then for all  $n \leq \hat{s}$  and  $m \leq \hat{r} + n$ , for all  $\varphi$  in up-down normal form with  $\text{vd}(\varphi) \leq (\hat{r}, \hat{s})$ ,  $\text{nd}(\varphi) \leq \hat{k}$ , for all upward expression  $\alpha^{\uparrow}$  in up-down normal form, and for all downward expression  $\alpha^{\downarrow}$  in up-down normal form with  $\text{vd}(\alpha^{\uparrow}), \text{vd}(\alpha^{\downarrow}) \leq (\hat{r}, \hat{s})$ ,  $\text{nd}(\alpha^{\uparrow}), \text{nd}(\alpha^{\downarrow}) \leq \hat{k}$ :

1. If  $xZ_{\hat{r}, \hat{s}}^{\hat{k}}x'$  then  $\mathcal{T}, x \models \varphi$  iff  $\mathcal{T}', x' \models \varphi$ .
2. If  $y \xrightarrow{n} x, y' \xrightarrow{n} x', xZ_{\hat{r}, \hat{s}}^{\hat{k}-1}x'$ , then  $(x, y) \in \llbracket \alpha^{\uparrow} \rrbracket^{\mathcal{T}}$  iff  $(x', y') \in \llbracket \alpha^{\uparrow} \rrbracket^{\mathcal{T}'}$ .
3. If  $y \xrightarrow{m} z, y' \xrightarrow{m} z', zZ_{\hat{r}', \hat{s}'}^{\hat{k}-1}z'$  for  $\hat{r}' = \hat{r} + n - m, \hat{s}' = \hat{s} - n + m$ , then  $(y, z) \in \llbracket \alpha^{\downarrow} \rrbracket^{\mathcal{T}}$  iff  $(y', z') \in \llbracket \alpha^{\downarrow} \rrbracket^{\mathcal{T}'}$ .

Hence, by Proposition 3.10, the main statement follows. We simultaneously show 1, 2 and 3 by induction on  $|\varphi| + |\alpha^{\downarrow}|$ .

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

Suppose  $\varphi = \langle \varepsilon = \alpha^{\uparrow} \alpha^{\downarrow} \rangle$ . We show  $\mathcal{T}, x \models \varphi \Rightarrow \mathcal{T}', x' \models \varphi$ , so assume  $\mathcal{T}, x \models \varphi$ . Suppose there are  $y, z \in \mathcal{T}$  and  $n \leq \hat{s}, m \leq \hat{r} + n$  such that  $y \xrightarrow{n} x, y \xrightarrow{m} z, (x, y) \in \llbracket \alpha^{\uparrow} \rrbracket^{\mathcal{T}}, (y, z) \in \llbracket \alpha^{\downarrow} \rrbracket^{\mathcal{T}}$  and  $\text{data}(x) = \text{data}(z)$ . By Zig, there are  $y', z' \in \mathcal{T}'$  such that  $zZ_{\hat{r}', \hat{s}'}^{\hat{k}-1}z'$  for  $\hat{r}' = \hat{r} + n - m, \hat{s}' = \hat{s} - n + m$ , and  $\text{data}(x') = \text{data}(z')$ . By inductive hypothesis 2 and 3,  $(x', y') \in \llbracket \alpha^{\uparrow} \rrbracket^{\mathcal{T}'}$  and  $(y', z') \in \llbracket \alpha^{\downarrow} \rrbracket^{\mathcal{T}'}$ . Hence  $\mathcal{T}', x' \models \varphi$ . The implication  $\mathcal{T}', x' \models \varphi \Rightarrow \mathcal{T}, x \models \varphi$  is analogous. The cases  $\varphi = \langle \varepsilon \neq \alpha^{\uparrow} \rangle$ , and  $\varphi = \langle \varepsilon \odot \alpha^{\uparrow} \rangle, \varphi = \langle \varepsilon \odot \alpha^{\downarrow} \rangle$  ( $\odot \in \{=, \neq\}$ ) and  $\varphi = \langle \alpha \rangle$  (for  $\alpha$  in up-down normal form) are shown in a similar way. The cases  $\varphi = \langle \varepsilon \odot \varepsilon \rangle$  ( $\odot \in \{=, \neq\}$ ) are trivial.

Let us now analyze item 2. Let  $\alpha^{\uparrow} = [\psi] \uparrow^n$  ( $n \geq 0$ ), and let

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

be such that

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

and  $xZ_{\hat{r}, \hat{s}}^{\hat{k}-1}x'$ . By Observation 3.14, we have  $x_0Z_{\hat{r}', \hat{s}'}^{\hat{k}-1}x'_0$ , for  $\hat{r}' = \hat{r} + n, \hat{s}' = \hat{s} - n$ . Assume by contradiction that  $(x', y') \notin \llbracket \alpha^{\uparrow} \rrbracket^{\mathcal{T}'}$ . This necessarily means that  $\mathcal{T}, x_0 \models \psi$  but  $\mathcal{T}', x'_0 \not\models \psi$ . But  $\psi$  is a subformula of  $\alpha^{\uparrow}$ ,  $\text{nd}(\psi) \leq \hat{k} - 1$  and  $\text{nd}(\psi) \leq (\hat{r}', \hat{s}')$  and this contradicts inductive hypothesis 1.

Item 3 is shown in a similar way. Let  $\alpha^{\downarrow} = \downarrow^m[\psi]$  ( $m \geq 0$ ), and let

$$z_0, \dots, z_m \in T \quad \text{and} \quad z'_0, \dots, z'_m \in T'$$

be such that

$$\begin{aligned} y &= z_0 \rightarrow z_1 \rightarrow \dots \rightarrow z_m = z && \text{in } \mathcal{T}, \\ y' &= z'_0 \rightarrow z'_1 \rightarrow \dots \rightarrow z'_m = z' && \text{in } \mathcal{T}', \end{aligned}$$

and  $zZ_{\hat{r}', \hat{s}'}^{\hat{k}-1}z'$ . Assume by contradiction that  $(y', z') \notin \llbracket \alpha^{\downarrow} \rrbracket^{\mathcal{T}'}$ . This necessarily means that  $\mathcal{T}, z_m \models \psi$  but  $\mathcal{T}', z'_m \not\models \psi$ . But  $\psi$  is a subformula of  $\alpha^{\downarrow}$ ,  $\text{nd}(\psi) \leq \hat{k} - 1$  and  $\text{nd}(\psi) \leq (\hat{r}', \hat{s}')$  and this contradicts inductive hypothesis 1.  $\square$

PROPOSITION A.5.  $\mathcal{T}, u \equiv_{r,s,k}^{\uparrow} \mathcal{T}', u'$  implies  $\mathcal{T}, u \leftrightarrow_{r,s,k}^{\uparrow} \mathcal{T}', u'$ .

PROOF. Fix  $u \in T$  and  $u' \in T'$  such that  $\mathcal{T}, u \equiv_{r,s,k}^{\uparrow} \mathcal{T}', u'$ . Define  $(Z_{\hat{r}, \hat{s}}^{\hat{k}})_{\hat{r} + \hat{s} \leq r + s, \hat{k} \leq k}$  by

$$xZ_{\hat{r}, \hat{s}}^{\hat{k}}x' \quad \text{iff} \quad \mathcal{T}, x \equiv_{\hat{r}, \hat{s}, \hat{k}}^{\uparrow} \mathcal{T}', x'.$$

We show that  $Z_{r,s}^k$  is a  $(r, s, k)$ -bisimulation between  $\mathcal{T}, u$  and  $\mathcal{T}', u'$ . By hypothesis,  $uZ_{r,s}^ku'$ . Now fix  $\hat{r} + \hat{s} \leq r + s, \hat{k} \leq k$ . By construction,  $Z_{\hat{r}, \hat{s}}^{\hat{k}}$  satisfies Harmony. Let us see that  $Z_{\hat{r}, \hat{s}}^{\hat{k}}$  satisfies Zig (the case for Zag is analogous). Suppose  $xZ_{\hat{r}, \hat{s}}^{\hat{k}}x'$ ,

$$\begin{aligned} y &= x_0 \rightarrow x_1 \rightarrow \dots \rightarrow v_n = x && \text{in } \mathcal{T}, \\ y &= z_0 \rightarrow z_1 \rightarrow \dots \rightarrow z_m = z && \text{in } \mathcal{T}, \end{aligned}$$

and  $data(x) = data(z)$  (the case  $data(x) \neq data(z)$  is shown in a similar way), where  $m \leq \hat{r} + n$ . Let  $P \subseteq T'^2$  be defined by

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

Since  $\mathcal{T}, x \equiv_{r,s,k}^{\dagger} \mathcal{T}', x'$ ,  $\text{vd}(\langle \varepsilon = \uparrow^n \downarrow^m \rangle) \leq (r, s)$ ,  $\text{nd}(\langle \varepsilon = \uparrow^n \downarrow^m \rangle) = 0$ , and  $\mathcal{T}, x \models \langle \varepsilon = \uparrow^n \downarrow^m \rangle$ , we conclude that  $P \neq \emptyset$ . We next show that there exists  $(y', z') \in P$  such that

- i.  $y' = x'_0 \rightarrow x'_1 \rightarrow \dots \rightarrow x'_n = x'$  in  $\mathcal{T}'$
- ii.  $y' = z'_0 \rightarrow z'_1 \rightarrow \dots \rightarrow z'_m = z'$  in  $\mathcal{T}'$ ,
- iii.  $\mathcal{T}, x \equiv_{\hat{r}, \hat{s}, \hat{k}-1}^{\dagger} \mathcal{T}', x'$ , and
- iv.  $\mathcal{T}, z \equiv_{\hat{r}', \hat{s}', \hat{k}-1}^{\dagger} \mathcal{T}', z'$ , where  $\hat{r}' = \hat{r} + n - m$ ,  $\hat{s}' = \hat{s} - n + m$ ,

and hence, by inductive hypothesis, Zig is satisfied by  $Z_{\hat{r}, \hat{s}}^{\hat{k}}$ . By way of contradiction, assume that for all  $(y', z') \in P$  satisfying i and ii we have either

- (a)  $\mathcal{T}, x \not\equiv_{\hat{r}, \hat{s}, \hat{k}-1}^{\dagger} \mathcal{T}', x'$ ; or
- (b)  $\mathcal{T}, z \not\equiv_{\hat{r}', \hat{s}', \hat{k}-1}^{\dagger} \mathcal{T}', z'$  for  $\hat{r}' = \hat{r} + n - m$ ,  $\hat{s}' = \hat{s} - n + m$ .

Fix  $\top$  as any tautology such that  $\text{vd}(\top) = (0, 0)$ ,  $\text{nd}(\top) = 0$ . For each  $(y', z') \in P$  we define formulas,  $\varphi_{y', z'}$  and  $\psi_{y', z'}$ , satisfying that  $\text{vd}(\varphi_{y', z'}) \leq (\hat{r}, \hat{s})$ ,  $\text{nd}(\varphi_{y', z'}) < \hat{k}$  and  $\text{vd}(\psi_{y', z'}) \leq (\hat{r}', \hat{s}')$ ,  $\text{nd}(\psi_{y', z'}) < \hat{k}$  as follows:

- Suppose (a) holds. Let  $\varphi_{y', z'}$  be such that  $\text{vd}(\varphi_{y', w'}) \leq (\hat{r}, \hat{s})$ ,  $\text{nd}(\varphi_{y', w'}) < \hat{k}$ , and such that  $\mathcal{T}, x \models \varphi_{y', z'}$  but  $\mathcal{T}', x' \not\models \varphi_{y', z'}$ ; and let  $\psi_{y', w'} = \top$ .
- Suppose (a) does not hold. Then (b) holds. Let  $\psi_{y', z'}$  be such that  $\text{vd}(\psi_{y', z'}) \leq (\hat{r}', \hat{s}')$ ,  $\text{nd}(\psi_{y', z'}) < \hat{k}$  and such that  $\mathcal{T}, z \models \psi_{y', z'}$  but  $\mathcal{T}', z' \not\models \psi_{y', z'}$ ; and let  $\varphi_{y', z'} = \top$ .

Let

$$\Phi = \bigwedge_{(y', z') \in P} \varphi_{y', z'} \quad \text{and} \quad \Psi = \bigwedge_{(y', z') \in P} \psi_{y', z'}. \quad (2)$$

Since  $\text{vd}(\varphi_{y', z'}) \leq (\hat{r}, \hat{s})$ ,  $\text{nd}(\varphi_{y', z'}) < \hat{k}$ , by Proposition 3.12, there are finitely many non-equivalent formulas  $\varphi_{y', z'}$ . The same applies to formulas  $\psi_{y', z'}$ . Hence both infinite conjunctions in (2) are equivalent to finite ones, and therefore without loss of generality we may assume that  $\Phi$  and  $\Psi$  are well-formed formulas.

Finally, let

$$\alpha^{\uparrow} = [\Phi]^{\uparrow n} \quad \text{and} \quad \alpha^{\downarrow} = \downarrow^m[\Psi].$$

By construction,  $\text{vd}(\alpha^{\uparrow} \alpha^{\downarrow}) \leq (\hat{r}, \hat{s})$ ,  $\text{nd}(\alpha^{\uparrow} \alpha^{\downarrow}) \leq \hat{k}$ . Furthermore,  $\mathcal{T}, x \models \langle \varepsilon = \alpha^{\uparrow} \alpha^{\downarrow} \rangle$  and  $\mathcal{T}', x' \not\models \langle \varepsilon = \alpha^{\uparrow} \alpha^{\downarrow} \rangle$ , but this contradicts the fact that  $\mathcal{T}, x \equiv_{\hat{r}, \hat{s}, \hat{k}}^{\dagger} \mathcal{T}', x'$ .  $\square$

**COROLLARY 3.17.**  $\Leftrightarrow_{r,s,k}^{\dagger}$  has finite index.

**PROOF.** Immediate from Theorem 3.16 and Proposition 3.12.

## B. PROOFS OF SECTION 4

**THEOREM 4.1.**

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

**PROOF.** The proofs are straightforward adaptations of the proofs of Propositions 3.8 and 3.9 and Propositions A.4 and A.5 respectively, and are omitted here. In particular, for the 'if' part, in the adaptation of the proofs of Propositions 3.9 and A.5, the simulations are defined by

$$\begin{aligned} xZ_i x' &\text{ iff } \mathcal{T}, x \equiv_i^{\downarrow} \mathcal{T}', x \\ xZ_{\hat{r}, \hat{s}}^{\hat{k}} x' &\text{ iff } \mathcal{T}, x \equiv_{\hat{r}, \hat{s}, \hat{k}}^{\dagger} \mathcal{T}', x \end{aligned}$$

respectively, and the conditions (a) and (b) on page 5 become now

- (a)  $[\exists i \in \{0, \dots, n\} \exists \varphi \in \text{XPath}_{\leq}^{\dagger}] \text{dd}(\varphi) \leq h - i \wedge \mathcal{T}, v_i \models \varphi \wedge \mathcal{T}', v'_i \not\models \varphi$ ; or
- (b)  $[\exists j \in \{0, \dots, m\} \exists \varphi \in \text{XPath}_{\leq}^{\dagger}] \text{dd}(\varphi) \leq h - j \wedge \mathcal{T}, w_j \models \varphi \wedge \mathcal{T}', w'_j \not\models \varphi$ ,

and

- (a)  $[\exists i \in \{0, \dots, n\} \exists \varphi \in \text{XPath}_{\leq}^{\dagger}] \text{vd}(\varphi) \leq (\hat{r} + i, \hat{s} - i) \wedge \text{nd}(\varphi) \leq k - 1 \wedge \mathcal{T}, v_i \models \varphi \wedge \mathcal{T}', v'_i \not\models \varphi$ ; or
- (b)  $[\exists j \in \{0, \dots, m\} \exists \varphi \in \text{XPath}_{\leq}^{\dagger}] \text{vd}(\varphi) \leq (\hat{r} + j', \hat{s} - j') \wedge \text{nd}(\varphi) \leq k - 1 \wedge \mathcal{T}, w_j \models \varphi \wedge \mathcal{T}', w'_j \not\models \varphi$

respectively.  $\square$

**LEMMA B.1.**

- (1)  $\{\mathcal{T}', u' \mid \mathcal{T}, u \xrightarrow{\downarrow} \mathcal{T}', u'\}$  is definable by an  $\text{XPath}_{\leq}^{\dagger+}$ -formula  $\chi_{\ell, u, \mathcal{T}}^{\dagger+}$  of downward depth  $\leq \ell$ .
- (2)  $\{\mathcal{T}', u' \mid \mathcal{T}, u \xrightarrow{\dagger}_{r,s,k} \mathcal{T}', u'\}$  is definable by an  $\text{XPath}_{\leq}^{\dagger+}$ -formula  $\chi_{r,s,k,u,\mathcal{T}}^{\dagger+}$  of vertical depth  $\leq (r, s)$  and nesting depth  $\leq k$ .

**PROOF.** For item (2), let  $\text{sim}_{r,s,k}^{\dagger}(\mathcal{T}, u) = \{\mathcal{T}', u' \mid \mathcal{T}, u \xrightarrow{\dagger}_{r,s,k} \mathcal{T}', u'\}$ . Let  $\Phi_{\mathcal{T}', u'}$  be the set of all positive formulas  $\varphi \in \text{XPath}_{\leq}^{\dagger+}$  of vertical depth at most  $(r, s)$  and nesting depth at most  $k$  so that  $\mathcal{T}', u' \models \varphi$ . Let  $\Psi$  be

$$\Psi = \bigvee_{\mathcal{T}', u' \in \text{sim}_{r,s,k}^{\dagger}(\mathcal{T}, u)} \bigwedge \Phi_{\mathcal{T}', u'}.$$

Since every  $\Phi_{\mathcal{T}', u'}$  is finite up to logical equivalence by Proposition 3.12, it follows that  $\Psi$  is a valid formula. We show that it defines  $\text{sim}_{r,s,k}^{\dagger}(\mathcal{T}, u)$ .

Let  $\mathcal{T}', u' \in \text{sim}_{r,s,k}^{\dagger}(\mathcal{T}, u)$ . Then,  $\mathcal{T}', u' \models \bigwedge \Phi_{\mathcal{T}', u'}$  and thus  $\mathcal{T}', u' \models \Psi$ . If on the other hand  $\mathcal{T}', u' \models \Psi$  we have that  $\mathcal{T}', u' \models \bigwedge \Phi_{\mathcal{T}'', u''}$  for some  $\mathcal{T}'', u'' \in \text{sim}_{r,s,k}^{\dagger}(\mathcal{T}, u)$  and then  $\mathcal{T}', u' \equiv_{r,s,k}^{\dagger} \mathcal{T}'', u''$ . By Theorem 3.16-3 we then have that  $\mathcal{T}', u' \Leftrightarrow_{r,s,k}^{\dagger} \mathcal{T}'', u''$ , and in particular  $\mathcal{T}'', u'' \xrightarrow{\dagger}_{r,s,k}$



$\mathcal{T}', u'$ . Since  $\mathcal{T}, u \xrightarrow{\dagger}_{r,s,k} \mathcal{T}'', u''$  and  $\mathcal{T}'', u'' \xrightarrow{\dagger}_{r,s,k} \mathcal{T}', u'$ , then  $\mathcal{T}, u \xrightarrow{\dagger}_{r,s,k} \mathcal{T}', u'$  (by transitivity of  $\xrightarrow{\dagger}_{r,s,k}$ ) and thus  $\mathcal{T}', u' \in \text{sim}_{r,s,k}^{\dagger}(\mathcal{T}, u)$ .

Item (1) is shown in a similar way, making use of Proposition 3.1 and Theorem 3.7-2.  $\square$

THEOREM 4.3.

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

PROOF. We start with item (1), for the case of  $\xrightarrow{\downarrow}$ . The ‘if’ part is straightforward from Theorem 4.1-2, and here we focus on the ‘only if’ part. Let  $\varphi$  be preserved under  $\xrightarrow{\downarrow}$ . Let  $\{(\mathcal{T}_i, u_i)\}_{i \leq n}$  be the set of all pointed models of  $\varphi$  modulo  $\xrightarrow{\downarrow}$  (which is finite due to Theorem 3.7-2 together with Proposition 3.1). We claim that

$$\mathcal{T}, u \models \varphi \text{ iff } \mathcal{T}_i, u_i \xrightarrow{\downarrow} \mathcal{T}, u \text{ for some } i \leq n. \quad (3)$$

On the one hand, if  $\mathcal{T}, u \models \varphi$  then there is  $i \leq n$  such that  $\mathcal{T}_i, u_i \xrightarrow{\downarrow} \mathcal{T}, u$ , and so  $\mathcal{T}_i, u_i \xrightarrow{\downarrow} \mathcal{T}, u$ . On the other hand, suppose  $\mathcal{T}_i, u_i \xrightarrow{\downarrow} \mathcal{T}, u$ . Since  $\varphi$  is preserved under  $\xrightarrow{\downarrow}$  and  $\mathcal{T}_i, u_i \models \varphi$ , we conclude  $\mathcal{T}, u \models \varphi$ .

Let  $\chi_{\ell, u_i, \mathcal{T}_i} \in \text{XPath}_{\leq}^{\downarrow+}$ ,  $\text{dd}(\psi_i) \leq \ell$ , be as in Lemma B.1-(1). Using (3) one shows that  $\bigvee_{i \leq n} \chi_{\ell, u_i, \mathcal{T}_i} \equiv \varphi$ .

For the case of  $\xrightarrow{\downarrow}$  of item (1), the ‘if’ direction follows from Theorem 4.1-1. For the ‘only if’ direction, let  $\varphi$  be preserved under  $\xrightarrow{\downarrow}$ . It is easy to see that  $\varphi$  is preserved under  $\xrightarrow{\downarrow}$  iff it is preserved under  $\xrightarrow{\downarrow}_{\text{dd}(\varphi)}$ . We can then apply the same reasoning as before and the statement follows.

Item (3) follows the same argument as item (1) but this time using Corollary 3.17 and Lemma B.1-(2).

Item (4) is straightforward from Theorem 4.1-3.

Item (2) follows from items (3) and (4) and the observation that  $\varphi$  is preserved under  $\xrightarrow{\dagger}$  iff it is preserved under  $\xrightarrow{\dagger}_{r,s,k \cdot (r+s+2)}$  for  $\text{vd}(\varphi) = (r, s)$  and  $\text{nd}(\varphi) = k$ .  $\square$

PROPOSITION B.2. Any  $\xrightarrow{\downarrow}$ -invariant  $\varphi(x) \in \text{FO}(\sigma)$  over [finite] data-trees that is  $\ell$ -local, is  $\xrightarrow{\downarrow}$ -invariant.

PROOF. Let  $\varphi(x)$  be  $\ell$ -local and  $\xrightarrow{\downarrow}$ -invariant. Suppose  $\mathcal{T}, u \xrightarrow{\downarrow} \mathcal{T}', u'$  and  $\mathcal{T} \models \varphi(u)$ . By  $\ell$ -locality,  $\mathcal{T}|_{eu} \models \varphi(u)$ . Now

$$\mathcal{T}, u \xrightarrow{\downarrow} \mathcal{T}', u' \text{ iff } (\mathcal{T}|_{eu}), u \xrightarrow{\downarrow} (\mathcal{T}'|_{eu'}), u' \quad (\text{Prop. 3.6})$$

$$\text{iff } (\mathcal{T}|_{eu}), u \xrightarrow{\downarrow} (\mathcal{T}'|_{eu'}), u'. \quad (\text{Prop. 3.5})$$

By  $\xrightarrow{\downarrow}$ -invariance,  $\mathcal{T}'|_{eu'} \models \varphi(u')$  and by  $\ell$ -locality again,  $\mathcal{T}' \models \varphi(u')$ .  $\square$

PROPOSITION B.3. If  $\varphi(x) \in \text{FO}(\sigma)$  is  $\xrightarrow{\downarrow}$ -invariant over [finite] data-trees, then there is  $\psi \in \ell\text{-XPath}_{\leq}^{\downarrow}$  such that  $\text{Tr}_x(\psi)$  is logically equivalent to  $\varphi$  over [finite] data-trees.

PROOF. By Corollary 3.2, for every data tree  $\mathcal{T}$  and  $u \in \mathcal{T}$  there is an  $\ell\text{-XPath}_{\leq}^{\downarrow}$  formula  $\chi_{\ell, \mathcal{T}, u}$  such that  $\mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'$  iff  $\mathcal{T}', u' \models \chi_{\ell, \mathcal{T}, u}$ . Let

$$\psi = \bigvee_{\mathcal{T} \models \varphi(u)} \chi_{\ell, \mathcal{T}, u}.$$

Since  $\chi_{\ell, \mathcal{T}, u} \in \ell\text{-XPath}_{\leq}^{\downarrow}$  and, by Proposition 3.1,  $\equiv_{\ell}^{\downarrow}$  has finite index, it follows that  $\psi$  is equivalent to a finite disjunction.

We now show that  $\varphi \equiv \text{Tr}_x(\psi)$ . Let us see that  $\varphi \models \text{Tr}_x(\psi)$ . Suppose  $\mathcal{T} \models \varphi(u)$ . Since  $\mathcal{T}, u \models \chi_{\ell, \mathcal{T}, u}$ , we have  $\mathcal{T}, u \models \psi$  and so  $\mathcal{T} \models \text{Tr}_x(\psi)(u)$ . Let us now see that  $\text{Tr}_x(\psi) \models \varphi$ . Assume  $\mathcal{T} \models \text{Tr}_x(\psi)(u)$ , and so  $\mathcal{T}, u \models \psi$ . Then there exists  $\mathcal{T}', u'$  such that  $\mathcal{T}' \models \varphi(u')$  and  $\mathcal{T}, u \models \chi_{\ell, \mathcal{T}', u'}$ . By the property of  $\chi_{\ell, \mathcal{T}', u'}$ , we have  $\mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'$  and since  $\varphi$  is  $\xrightarrow{\downarrow}$ -invariant (and hence  $\equiv_{\ell}^{\downarrow}$ -invariant by Theorem 3.7-2) we conclude  $\mathcal{T} \models \varphi(u)$ .  $\square$

## C. PROOFS OF SECTION 5

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

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

PROOF. The proof that  $\mathcal{T}, u \xleftrightarrow{\dagger} \mathcal{T}', u' \Rightarrow \mathcal{T}, u \equiv^{\dagger} \mathcal{T}', u'$  follows from a simple adaptation of Proposition 3.8 to the logic  $\text{XPath}_{\leq}^{\dagger, \downarrow^*, \uparrow^*}$  and Lemma A.1. The fact that for finitely branching,  $\mathcal{T}, u \equiv^{\dagger} \mathcal{T}', u' \Rightarrow \mathcal{T}, u \xleftrightarrow{\dagger} \mathcal{T}', u'$  is straightforward from Theorem 3.7-1 since  $\equiv^{\dagger} \subseteq \xrightarrow{\dagger}$ .

The cases for  $\text{XPath}_{\leq}^{\dagger, \downarrow^*, \uparrow^*}$ ,  $\text{XPath}_{\leq}^{\dagger, \downarrow^*, \uparrow^*}$  and  $\text{XPath}_{\leq}^{\dagger, \downarrow^*, \uparrow^*+}$  are analogous.  $\square$

## D. PROOFS OF SECTION 6

PROPOSITION 6.1. For any  $\text{XPath}_{\leq}^{\dagger}$ - [resp.  $\text{XPath}_{\leq}^{\downarrow}$ -] path expression  $\alpha$  there is an equivalent  $\text{XPath}_{\leq}^{\dagger}$ - [resp.  $\text{XPath}_{\leq}^{\downarrow}$ -] path expression  $\alpha'$  in simple normal form. Further,  $\alpha'$  can be computed in polynomial time from  $\alpha$ .

PROOF. The translation is straightforward, given the following equivalences:

$$\begin{aligned} \varepsilon &\equiv [\top] \\ \alpha &\equiv [\top]\alpha \equiv \alpha[\top] \\ \alpha[\varphi][\psi]\beta &\equiv \alpha[\varphi \wedge \psi]\beta \end{aligned}$$

where  $\top$  denotes any fixed tautology, for example  $a \vee \neg a$ , for some  $a \in \mathcal{A}$ .  $\square$

LEMMA 6.5. The  $\text{FO}(\sigma)$ -formula

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

is  $\xrightarrow{\dagger}$ -invariant though not logically equivalent over [finite] data-trees to any  $\text{XPath}_{\leq}^{\dagger}$ -formula.

PROOF. Let  $\varphi(x)$  be the  $\text{FO}(\sigma)$ -formula for there is a node labeled  $a$  in the tree, i.e.,

$$\varphi(x) = (\exists y) P_a(y).$$

We prove that  $\varphi$  is  $\leftrightarrow^\downarrow$ -invariant over [finite] data-trees, though it is not logically equivalent over [finite] data-trees to any  $\text{XPath}_{\downarrow}^\downarrow$ -formula.

To see that  $\varphi$  is  $\leftrightarrow^\downarrow$ -invariant over [finite] data-trees, take  $\mathcal{T}, u$  and  $\mathcal{T}', u'$  such that  $\mathcal{T}, u \leftrightarrow^\downarrow \mathcal{T}', u'$  and  $\mathcal{T} \models \varphi(u)$ . Furthermore, suppose that  $\mathcal{T}, u \models \uparrow^m \downarrow^n a$  for adequate  $n$  and  $m$ . By Theorem 3.16,  $\mathcal{T}', u' \models \uparrow^n \downarrow^m a$  and so  $\mathcal{T}' \models \varphi(u')$ .

Assume by contradiction that there is  $\psi \in \text{XPath}_{\downarrow}^\downarrow$  such that  $\mathcal{T}, u \models \psi$  iff  $\mathcal{T} \models \varphi(u)$  for all data-tree  $\mathcal{T}$  and  $u \in T$ . Suppose  $\text{vd}(\psi) = (r, s)$  and  $\text{nd}(\psi) = k$ . Let  $\mathcal{T}$  be a data tree formed by a chain of length  $r + 1$  starting from the root  $u$  with all its nodes containing a label  $b$  except the leaf, which has label  $a$  (the data values are irrelevant). By Proposition 3.15 we have  $\mathcal{T}, u \leftrightarrow_{r,s,k}^\downarrow (\mathcal{T}|_r^s u), u$ . Since  $\mathcal{T}, u \models \psi$ , by Theorem 3.16, we have  $(\mathcal{T}|_r^s u), u \models \psi$ , and so  $\mathcal{T}|_r^s u \models \varphi(u)$ . This last fact is a contradiction because no node of  $\mathcal{T}|_r^s u$  is labeled with  $a$ .  $\square$

**PROPOSITION 6.6.** *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  $\psi \in (r, s, k)\text{-XPath}_{\downarrow}^\downarrow$  such that  $\text{Tr}_x(\psi)$  is logically equivalent to  $\varphi$  over [finite] data-trees.*

**PROOF.** By Corollary 3.13, for every data tree  $\mathcal{T}$  and  $u \in T$  there is an  $(r, s, k)\text{-XPath}_{\downarrow}^\downarrow$  formula  $\chi_{r,s,k,\mathcal{T},u}$  such that  $\mathcal{T}, u \equiv_{r,s,k}^\downarrow \mathcal{T}', u'$  iff  $\mathcal{T}', u' \models \chi_{r,s,k,\mathcal{T},u}$ . Let

$$\psi = \bigvee_{\mathcal{T} \models \varphi(u)} \chi_{r,s,k,\mathcal{T},u}.$$

As  $\chi_{r,s,k,\mathcal{T},u} \in (r, s, k)\text{-XPath}_{\downarrow}^\downarrow$  and, by Proposition 3.12,  $\equiv_{r,s,k}^\downarrow$  has finite index, it follows that  $\psi$  is equivalent to a finite disjunction. The proof that  $\varphi(x) \equiv \text{Tr}_x(\psi)$  is similar to Proposition B.3, as we show next. Let us see that  $\varphi \models \text{Tr}_x(\psi)$ . Suppose  $\mathcal{T} \models \varphi(u)$ . Since  $\mathcal{T}, u \models \chi_{r,s,k,\mathcal{T},a}$ , we have  $\mathcal{T}, u \models \psi$  and so  $\mathcal{T} \models \text{Tr}_x(\psi)(u)$ . Let us see that  $\text{Tr}_x(\psi) \models \varphi$ . Assume  $\mathcal{T} \models \text{Tr}_x(\psi)(u)$ , and so  $\mathcal{T}, u \models \psi$ . Then there exists  $\mathcal{T}', u'$  such that  $\mathcal{T}' \models \varphi(u')$  and  $\mathcal{T}, u \models \chi_{r,s,k,\mathcal{T}',u'}$ . By the property of  $\chi_{r,s,k,\mathcal{T}',u'}$ , we have  $\mathcal{T}, u \equiv_{r,s,k}^\downarrow \mathcal{T}', u'$  and since  $\varphi$  is  $\leftrightarrow_{r,s,k \cdot (r+s+2)}^\downarrow$ -invariant (and hence  $\equiv_{r,s,k}^\downarrow$ -invariant by Theorem 3.16-2) we conclude  $\mathcal{T} \models \varphi(u)$ .  $\square$

**LEMMA 6.7.** *The  $\text{FO}(\sigma)$ -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  $\text{XPath}_{\downarrow}^{\downarrow,E}$ -formula.*

**PROOF.** Let  $\varphi(x)$  be the  $\text{FO}(\sigma)$ -formula for *there are two nodes with same data value and labels  $a$  and  $b$  respectively*, i.e.,

$$\varphi(x) = (\exists y, z) [y \approx z \wedge P_a(y) \wedge P_b(z)].$$

We show that  $\varphi$  cannot be expressed in  $\text{XPath}_{\downarrow}^{\downarrow,\uparrow,E}$ . Suppose, by means of contradiction, that there is a formula  $\psi \in \text{XPath}_{\downarrow}^{\downarrow,\uparrow,E}$  expressing  $\varphi$ , with  $\text{vd}(\psi) = (r, s)$  ( $\text{vd}(\cdot)$  for  $\text{XPath}_{\downarrow}^{\downarrow,\uparrow,E}$  is defined as in Table 2 plus the clause  $\text{vd}(E\varphi) = \text{vd}(\varphi)$ ). Let  $n = r + s$ , and let  $\mathcal{T}$  be the chain-like data-tree

$$u_0 \rightarrow u_1 \rightarrow \dots \rightarrow u_n$$

such that  $\text{label}(u_0) = a$ ,  $\text{label}(u_n) = b$ ,  $\text{label}(u_i) = c$  for  $i \in \{1, \dots, n-1\}$  and  $\text{data}(u_i) = i$  for  $i \in \{0, \dots, n\}$ . Let  $\mathcal{T}'$  be the chain-like data-tree

$$u'_0 \rightarrow u'_1 \rightarrow \dots \rightarrow u'_n$$

such that  $\text{label}(u'_i) = \text{label}(u_i)$  for  $i \in \{0, \dots, n\}$ ,  $\text{data}(u'_i) = \text{data}(u_i)$  for  $i \in \{0, \dots, n-1\}$  and  $\text{data}(u'_n) = 0$ . Note that  $\mathcal{T} \not\models \varphi(u_0)$  and  $\mathcal{T}' \models \varphi(u'_0)$ . However, one can show that for all  $i \in \{0, \dots, n\}$  we have  $\mathcal{T}, u_i \models \psi$  iff  $\mathcal{T}', u'_i \models \psi$ . Hence,  $\psi$  does not express  $\varphi$  and thus  $\varphi$  is not expressible in  $\text{XPath}_{\downarrow}^{\downarrow,\uparrow,E}$ .  $\square$

## E. PROOFS OF SECTION 7

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

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

**PROOF.** Consider the data trees defined in Figure 12 for every  $k$ . Note that  $\equiv_{\ell,k+1}^\downarrow \subseteq \equiv_{\ell,k}^\downarrow$  and  $\equiv_{\ell+1,k}^\downarrow \subseteq \equiv_{\ell,k}^\downarrow$  by definition. We show that  $\equiv_{\ell,k}^\downarrow \neq \equiv_{\ell,k+1}^\downarrow$  for all  $\ell \geq k + 1$ . For this purpose, we show that  $\mathcal{T}_k^1, x_k^1 \equiv_{k+1,k}^\downarrow \mathcal{T}_k^1, x_k^1$  but  $\mathcal{T}_k^1, x_k^1 \not\equiv_{k+1,k+1}^\downarrow \mathcal{T}_k^1, x_k^1$ .

The fact that  $\mathcal{T}_k^1, x_k^1 \not\equiv_{k+1,k+1}^\downarrow \mathcal{T}_k^1, x_k^1$  comes from the fact that the property “there is a path of length  $k + 1$  ending with a label  $a$  whose every pair of consecutive nodes have distinct data value” is definable with the following formula  $\varphi_{k+1}$  of depth  $k + 1$  and nesting depth  $k + 1$ ,

$$\begin{aligned} \varphi_1 &= \langle \varepsilon \neq \downarrow[a] \rangle \\ \varphi_{i+1} &= \langle \varepsilon \neq \downarrow[\varphi_i] \rangle \text{ for } i > 0. \end{aligned}$$

Since  $\mathcal{T}_k^1, x_k^1 \models \varphi_{k+1}$  but  $\mathcal{T}_k^1, x_k^1 \not\models \varphi_{k+1}$ , it follows that  $\mathcal{T}_k^1, x_k^1 \not\equiv_{k+1,k+1}^\downarrow \mathcal{T}_k^1, x_k^1$ .

To show  $\mathcal{T}_k^1, x_k^1 \equiv_{k+1,k}^\downarrow \mathcal{T}_k^1, x_k^1$  we actually use Proposition 7.1 and show  $\mathcal{T}_k^1, x_k^1 \leftrightarrow_{k+1,k}^\downarrow \mathcal{T}_k^1, x_k^1$ . Note that  $\mathcal{T}_k^1$  and  $\mathcal{T}_k^2$  (resp.  $\mathcal{T}_k^1$  and  $\mathcal{T}_k^2$ ) are equal modulo renaming of data values, so we are also showing that the roots of any two data trees with subindex  $k$  are  $(k + 1, k)$ -bisimilar.

**OBSERVATION E.1.** *Note that the set of immediate subtrees of the roots of  $\mathcal{T}_k^1, \mathcal{T}_k^1, \mathcal{T}_k^2, \mathcal{T}_k^2$  are the same as those of  $\mathcal{T}_k^1, \mathcal{T}_k^2, \mathcal{T}_k^2$  (and of  $\mathcal{T}_k^1, \mathcal{T}_k^1, \mathcal{T}_k^2$ ) by construction.*

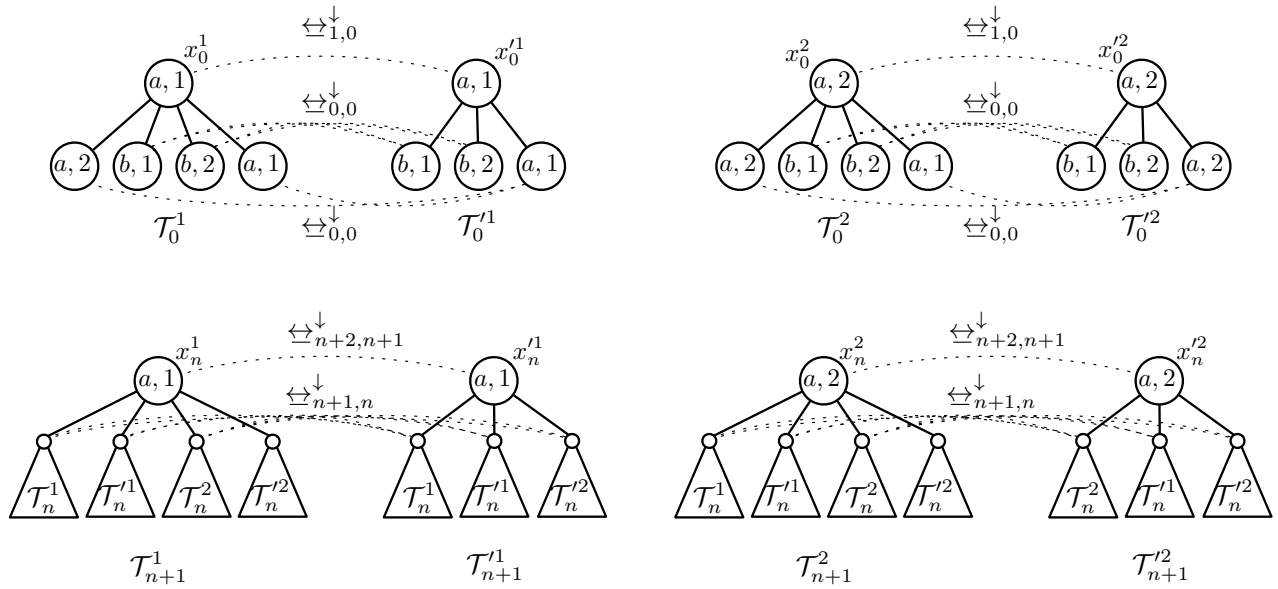
We now show  $\mathcal{T}_k^1, x_k^1 \leftrightarrow_{k+1,k}^\downarrow \mathcal{T}_k^1, x_k^1$ . For every  $j \leq k + 1, t \leq k$ , let  $Z_{j,t}$  be the set of all pairs  $(x, y) \in T_k^1 \times T_k^1$  so that  $x$  and  $y$  are some  $x_{k'}^i$  or  $x_{k'}^i$  for  $i \in \{1, 2\}$  and  $k' \geq t$ .<sup>2</sup> Observe that

$$Z_{j+1,t} \subseteq Z_{j,t} \text{ for all } j, t \leq k. \quad (4)$$

We show that  $(Z_{j,t})_{j \leq k+1, t \leq k}$  verify the bisimulation conditions. We proceed by induction on  $j + t$ . The base case,  $j = t = 0$ , is trivial. The case  $l > 0, t = 0$  is also straightforward.

Suppose then that  $t > 0$ . Let  $(u, u') \in Z_{j,t}$ . Again, Harmony is met since  $Z_{l,t}$  relates only nodes with label  $a$ . Let us suppose that  $u$  is some  $x_{l'}^1$  and  $u'$  is  $x_{l'}^1$  for some  $l' \leq t$ , the other cases being similar or simpler.

<sup>2</sup>Note that  $x_{k'}^i$  or  $x_{k'}^i$  do not necessarily uniquely identify one node, but many possible. The intended meaning is that  $x, y$  can be *any* of them.



**Figure 12: Definition of data trees  $\mathcal{T}_n^i$ ,  $\mathcal{T}_n'^i$  ( $n \geq 0, i \in \{1, 2\}$ ) for proof of Theorem 7.2.**

Let us now show Zig. Let  $v, w$  be so that  $x_{t'}^1 \xrightarrow{n} v$  and  $x_{t'}^1 \xrightarrow{m} w$  with  $n, m \leq j$ .

- If  $v$  is inside the subtree  $\mathcal{T}_{t'-1}^2$  of  $\mathcal{T}_{t'}^1$ , but it is not  $x_{t'-1}^2$ , then we choose  $v'$  as the corresponding<sup>3</sup> node inside the subtree  $\mathcal{T}_{t'-1}^1$  of  $\mathcal{T}_{t'}^1$ . Note that  $\text{data}(v) = \text{data}(v')$  by Observation E.1. Further, since every node of  $\mathcal{T}_{t'-1}^1$  is in a  $Z_{j,t-1}$ -relation with the corresponding node in  $\mathcal{T}_{t'-1}^2$  by construction of  $Z_{j,t-1}$ , it follows that  $(\xrightarrow{i}v)Z_{j,t-1}(\xrightarrow{i}v')$  for all  $i \leq n$ . Thus, by (4),  $(\xrightarrow{i}v)Z_{j-n+i,t-1}(\xrightarrow{i}v')$  for all  $i \leq n$ .
- If, on the other hand,  $v$  is  $x_{t'-1}^2$ , we choose  $v'$  as the root of  $\mathcal{T}_{t'-1}^2$ ,  $x_{t'-1}^2$ . Again, we have that  $\text{data}(v') = \text{data}(v)$  and by construction that  $vZ_{j,t-1}v'$ . Thus, by (4),  $vZ_{j-1,t-1}v'$ .
- Finally, if  $v$  falls outside  $\mathcal{T}_{t'-1}^2$ , we choose  $v'$  as the same node in  $\mathcal{T}_{t'}^1$ , where of course we will have that  $\text{data}(v) = \text{data}(v')$  and that  $(\xrightarrow{i}v)Z_{j,t-1}(\xrightarrow{i}v')$  for all  $i \leq n$ . Thus, by (4),  $(\xrightarrow{i}v)Z_{j-n+i,t-1}(\xrightarrow{i}v')$  for all  $i \leq n$ .

We do the same with  $w$  and  $w'$ . Since in every case we can reach a node with the same data value and so that the corresponding nodes in the path are  $Z_{j,t-1}$ -related, it follows that the Zig condition is satisfied. The Zag condition is only easier, and hence we conclude that  $\mathcal{T}_k^1, x_{k+1,k} \xleftrightarrow{\downarrow} \mathcal{T}_k^1, x'$  for every  $k$ .

We therefore have that  $\equiv_{\ell,k+1}^{\downarrow} \subsetneq \equiv_{\ell,k}^{\downarrow}$  for all  $\ell \geq k+1$ .

The fact that  $\equiv_{\ell+1,k}^{\downarrow} \subsetneq \equiv_{\ell,k}^{\downarrow}$  is of course trivial, formulas of depth  $\ell+1$  can express “the tree has at least depth  $\ell+1$ ”, which cannot be expressed by formulas of depth  $\ell$ .

<sup>3</sup>Remember that  $\mathcal{T}_{t'-1}^1$  and  $\mathcal{T}_{t'-1}^2$  are isomorphic modulo a renaming of data values, so by *corresponding* we mean the node in the same position in the tree

It remains to show that  $\equiv_{\ell,k}^{\downarrow} = \equiv_{\ell,k+1}^{\downarrow}$  for all  $\ell \leq k$ . To show this, we prove  $\mathcal{T}, x \xleftrightarrow{\downarrow} \mathcal{T}', x'$  for every  $\mathcal{T}, \mathcal{T}'$  so that  $\mathcal{T}, x \xleftrightarrow{\downarrow} \mathcal{T}', x'$ . We prove it by induction on  $\ell+k$ . The base case is easy.

For the inductive case, let  $Z_{j,t} = \xleftrightarrow{\downarrow} Z_{j,t}$  for all  $j \leq \ell, t \leq k$ . Hence,  $(Z_{j,t})_{j \leq \ell, t \leq k}$  verify the bisimulation conditions. Let  $Z_{\ell,k+1} = \{(x, x')\}$ . We show that  $Z_{\ell,k+1}$  together with  $(Z_{j,t})_{j \leq \ell, t \leq k}$  verifies the bisimulation conditions. Harmony follows from  $xZ_{\ell,k}x'$ . We show Zig since Zag is equivalent. Suppose  $x \xrightarrow{n} v, x \xrightarrow{m} w$  with  $n, m \leq \ell$ . Then, since  $Z_{\ell,k}$  verifies Zig, there are  $x' \xrightarrow{n} v', x' \xrightarrow{m} w'$  where

- (1)  $\text{data}(v) = \text{data}(w')$  iff  $\text{data}(v') = \text{data}(w')$ ,
- (2)  $(\xrightarrow{i}v)Z_{\ell-n+i,k-1}(\xrightarrow{i}v')$  for all  $i \in \{0, \dots, n-1\}$ , and
- (3)  $(\xrightarrow{i}w)Z_{\ell-m+i,k-1}(\xrightarrow{i}w')$  for all  $i \in \{0, \dots, m-1\}$ .

Since  $\ell \leq k$ , then  $\ell-n+i \leq k-1$ . Further,  $\ell-n+i+k < \ell+k$ , which means that we can apply the inductive hypothesis. Hence, by inductive hypothesis,  $\mathcal{T}, (\xrightarrow{i}v) \xleftrightarrow{\downarrow} \mathcal{T}', (\xrightarrow{i}v')$  and thus  $(\xrightarrow{i}v)Z_{\ell-n+i,k}(\xrightarrow{i}v')$ . By an identical reasoning,  $\mathcal{T}, (\xrightarrow{i}w) \xleftrightarrow{\downarrow} \mathcal{T}', (\xrightarrow{i}w')$  and thus  $(\xrightarrow{i}w)Z_{\ell-n+i,k}(\xrightarrow{i}w')$ . Thus, the Zig condition for  $\xleftrightarrow{\downarrow} Z_{\ell,k+1}$  is verified. The Zag condition holds by symmetry.  $\square$

With respect to vertical XPath, note that since  $\equiv_{r,s,k}^{\downarrow} \subseteq \equiv_{r',s',k'}^{\downarrow}$  for all  $(r, s, k) \leq (r', s', k')$ , as a consequence of Proposition 3.11 we obtain that for every  $r, s, k$  with  $r+s \geq 2$  there is some  $k' > k$  so that  $\equiv_{r,s,k}^{\downarrow} \supseteq \equiv_{r,s,k'}^{\downarrow}$ . In fact, we conjecture that  $\equiv_{r,s,k}^{\downarrow} \supseteq \equiv_{r,s,k+1}^{\downarrow}$  for every  $k$ . We argue that this can be proven through the models  $(\mathcal{T}_n)_n$  in the proof of Proposition 3.11, by showing that  $\mathcal{T}_k, x_{r',s'} \equiv_{r,s,k}^{\downarrow} \mathcal{T}_{k+1}, x_{r',s'}$  but  $\mathcal{T}_k, x_{r',s'} \not\equiv_{r,s,k+1}^{\downarrow} \mathcal{T}_{k+1}, x_{r',s'}$  for every  $(r, s) \geq (r', s')$ . The fact that  $\equiv_{r,s,k}^{\downarrow} \supseteq \equiv_{r+1,s,k}^{\downarrow}$  and  $\equiv_{r,s,k}^{\downarrow} \supseteq \equiv_{r,s+1,k}^{\downarrow}$  are straightforward. We then obtain the following.

CLAIM E.2.  $\equiv_{r,s,k}^{\dagger} \supseteq \equiv_{r',s',k'}^{\dagger}$  for all  $(r, s, k) < (r', s', k')$ ,  
 $r + s \geq 2$ .