



**HAL**  
open science

# Model Theory of XPath on Data Trees. Part I: Bisimulation and Characterization

Diego Figueira, Santiago Figueira, Carlos Areces

► **To cite this version:**

Diego Figueira, Santiago Figueira, Carlos Areces. Model Theory of XPath on Data Trees. Part I: Bisimulation and Characterization. Journal of Artificial Intelligence Research, 2015. hal-01778416

**HAL Id: hal-01778416**

**<https://hal.science/hal-01778416>**

Submitted on 25 Apr 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.

# Model Theory of XPath on Data Trees.

## Part I: Bisimulation and Characterization

**Diego Figueira**  
*CNRS, LaBRI, France*

DIEGO.FIGUEIRA@LABRI.FR

**Santiago Figueira**  
*Universidad de Buenos Aires and CONICET, Argentina*

SANTIAGO@DC.UBA.AR

**Carlos Areces**  
*Universidad Nacional de Córdoba and CONICET, Argentina*

CARLOS.ARECES@GMAIL.COM

### 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`, `parent`, `ancestor` and `descendant` 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 both 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.

## 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 (Clark & DeRose, 1999).

Core-XPath (term coined in Gottlob, Koch, & Pichler, 2005) 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 reasoning 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 Document

Type Definitions (DTDs) (Marx, 2004; Benedikt, Fan, & Geerts, 2008). Moreover, it is known that it is equivalent to  $\text{FO}^2$  (first-order logic with two variables) over an appropriate signature on trees in terms of expressive power (Marx & de Rijke, 2005), and that it is strictly less expressive than PDL with converse over trees (Benedikt & Koch, 2008). 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 the work of Bojańczyk, Muscholl, Schwentick, and Segoufin (2009). Here, we will call this logic  $\text{XPath}_=$ . Models of  $\text{XPath}_=$  are *data trees* which can be seen as XML documents. A data tree is a tree whose nodes contains a *label* from a finite alphabet and a *data value* from an infinite domain (see Figure 1 for an example). We will relax the condition on finiteness and consider also infinite data trees, although all our results hold also on finite structures.

The main characteristic of  $\text{XPath}_=$  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 the works of Figueira (2010) and Bojańczyk and Parys (2011). In particular, all the logics studied in this article have a decidable satisfiability problem (Figueira & Segoufin, 2011; Figueira, 2012); 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.

In this article we focus on the basic model theory tool of *bisimulations*, which defines the structural conditions necessary for ensuring that two models coincide in all the properties expressible by a logic. Whereas the basic notion of bisimulation was introduced for the basic modal logic, one can find adequate notions of bisimulations for different logics, in the sense that they capture the notion of indistinguishability of the logic. The challenge here is to find adequate notions of bisimulation for logics such as XPath, whose navigation is akin to modal logics such as PDL, but which can also test for equality of data values in the data tree.

**Contribution:**  $\text{XPath}_=$  can navigate the data tree by means of axes like **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  $\downarrow$ );  $\text{XPath}_=(\uparrow\downarrow)$  ( $\text{XPath}_=$  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 we introduce bisimulation notions for  $\text{XPath}_=(\downarrow)$ ,  $\text{XPath}_=(\downarrow\downarrow^*)$ ,  $\text{XPath}_=(\downarrow^*)$ ,  $\text{XPath}_=(\uparrow\downarrow)$ ,  $\text{XPath}_=(\downarrow\uparrow\downarrow^*\uparrow^*)$  and show that they precisely characterize the logical equivalence relation of the corresponding logic. We also consider fine grained versions of these bisimulations indexed by two measures of formula complexity. The first measure of formula complexity consists on the maximum number of nested axes in a formula, which we call *downward depth* in the case of  $\text{XPath}_=(\downarrow)$  and *vertical depth* in the case of  $\text{XPath}_=(\uparrow\downarrow)$ . The second one is the number of nested subformulas, called *nesting depth*.

The notion of bisimulation for  $\text{XPath}_=(\uparrow\downarrow)$  relies on a normal form which we also introduce. Basically, this normal form restricts the navigation of the expressions to be very simple: either going downward or going upward and then downward. Similar normal forms for languages on trees are folklore, and our work here consists mainly in adapting them to the setup of tests for data values.

We also 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 §4 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  $\text{XPath}_=(\uparrow\downarrow)$  instead the characterization fails as we show in a counter-example. However, a weaker result can still be established, namely that if a first-order formula is bisimulation-invariant, for the bisimulation notion corresponding to any fixed number of nested axes, then it is equivalent to an  $\text{XPath}_=(\uparrow\downarrow)$  formula.
- Using bisimulations we show some (non)expressivity results about  $\text{XPath}_=$  in §5. We show, for example, that formulas of  $\text{XPath}_=(\downarrow)$  with nesting depth  $n+1$  and downward depth  $d$  have more expressive power than those of nesting and downward depth  $n$  and  $d$  respectively, as long as  $n < d$ .
- All results are proved both over the class of arbitrary (possibly infinite) data trees, and over the class of finite data trees.

## 1.1 Related Work

The notion of bisimulation was introduced independently by van Benthem (1976) in the context of modal correspondence theory, by Milner (1980) and Park (1981) in concurrency theory, and by Forti and Honsell (1983) in non-wellfounded set theory (for a historical outlook see the work of Sangiorgi, 2009). 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. XPath has been known to be closely connected to known modal languages, such as PDL and modal  $\mu$ -calculus, depending on the fragments taken into account (ten Cate, Fontaine, & Litak, 2010). However, the fragments studied hitherto are “data unaware”, that is, allowing to express of the structure of the model as well of the fixed set of labels.

To the best of our knowledge the present is the first work on bisimulations and invariance with logics with data tests.

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 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 (van Benthem, 1976). Van Benthem’s original result over arbitrary structures was proved to hold for finite structures by Rosen (1997). The proof was then simplified and unified by Otto (2004a, 2006), and later expanded by Dawar and Otto (2009) to other classes of structures. Other formalisms of different expressive power have also been considered for querying data trees, such as first-order logic with two variables (Bojańczyk et al., 2009), tree patterns (David, 2008; Figueira & Libkin, 2014), register automata (Neven, Schwentick, & Vianu, 2004),  $\mu$ -calculus with registers (Jurdziński & Lazić, 2011), or datalog programs (Abiteboul, Bourhis, Muscholl, & Wu, 2013). In the absence of data values, logics for semi-structured databases can often be seen as modal logics. In fact, structural characterizations for XPath without equality test were studied in the work of Gyssens, Paredaens, Gucht, and Fletcher (2006), and XPath is known to be captured by PDL (Harel, 1984), whose bisimulation is well-understood (Blackburn, de Rijke, & Venema, 2001). It is then natural to look for an intuitive bisimulation definition for  $\text{XPath}_=$ .

The first significant result concerning an algorithmic solution to the bisimulation problem for the basic modal logic was given by Hopcroft (1971), with a polynomial time algorithm for state minimization in a deterministic finite automaton. This problem is equivalent to determine a coarsest partition of a set that is stable with respect to a finite set of functions. Paige and Tarjan (1987) solved the problem for the more general case, where the restriction of stability concerns a finite set of *relations*. Kanellakis and Smolka (1990) are the first to recognize that the algorithm of Paige and Tarjan can be used to determine the maximum bisimulation for the basic modal logic in an arbitrary graph. Hence it can be decided in polynomial time whether two nodes in finite models are bisimilar for the basic modal logic —though the length of the actual formulas which distinguish two non-bisimilar nodes cannot be polynomially bounded with respect to the size of the models (Figueira & Gorín, 2010). In our case, deciding whether two nodes in finite data trees are bisimilar can also be solved in polynomial time.

A preliminary version of the present paper appeared in a work of Figueira, Figueira, and Areces (2014). Also, there is a continuation (a ‘Part II’): a work by Abriola, Descotte, and Figueira (2015), after the preliminary version by Abriola, Descotte, and Figueira (2014), where they address further model theoretical questions such as definability and separation, and bisimulation notions for pairs of nodes (instead of single nodes) to capture the idea of “indistinguishability by means of path expressions” (instead of node expressions).

## 2. Preliminaries

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  $\lambda$  for the empty string.

Let  $Trees(A)$  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  $Trees(\mathbb{A} \times \mathbb{D})$  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.

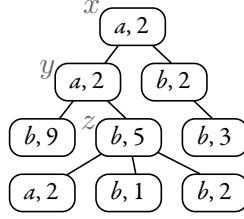


Figure 1: A data tree of  $Trees(\mathbb{A} \times \mathbb{D})$  with  $\mathbb{A} = \{a, b\}$  and  $\mathbb{D} = \mathbb{N}$ .

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  $label(x) \in \mathbb{A}$  to denote the node's label, and  $data(x) \in \mathbb{D}$  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. We write  $x \xrightarrow{*} y$  to denote that  $(x, y)$  is in the reflexive transitive closure of  $\rightarrow$ .  $(x \xrightarrow{n})$  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).

Let  $P$  be a property on nodes of data trees. When property  $P$  is true at node  $u$  of data tree  $\mathcal{T}$ , we say that  $(\mathcal{T}, u)$  satisfies  $P$ . For any binary relation  $R$  over nodes 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  $P$ .

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.  $\text{XPath}_=$  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  $\text{XPath}_=(\mathcal{O})$ , with  $\mathcal{O} \subseteq \{\downarrow, \downarrow^*, \uparrow, \uparrow^*\}$ , is defined by mutual recursion as follows:

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

A **formula** of  $\text{XPath}_=(\mathcal{O})$  is either a node expression or a path expression.

We formally define the semantics of  $\text{XPath}_=$  as follows, for  $\mathcal{T}$  a data tree:

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

As an example, if  $\mathcal{T}$  is the data tree shown in Figure 1, then

$$\llbracket \langle \downarrow_* [b \wedge \langle \downarrow [b] \neq \downarrow [b] \rangle] \rangle \rrbracket^{\mathcal{T}} = \{x, y, z\},$$

and 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 say that  $\mathcal{T}, u$  is a **pointed data tree**, we write  $\mathcal{T}, u \models \varphi$  to denote  $u \in \llbracket [\varphi] \rrbracket^{\mathcal{T}}$ , and we say that  $\mathcal{T}, u$  **satisfies**  $\varphi$ . We say that the node expressions  $\varphi, \psi$  of XPath<sub>=</sub> are **equivalent** (notation:  $\varphi \equiv \psi$ ) iff  $\llbracket [\varphi] \rrbracket^{\mathcal{T}} = \llbracket [\psi] \rrbracket^{\mathcal{T}}$  for all data trees  $\mathcal{T}$ . Similarly, path expressions  $\alpha, \beta$  of XPath<sub>=</sub> are **equivalent** (notation:  $\alpha \equiv \beta$ ) iff  $\llbracket \langle \alpha \rangle \rrbracket^{\mathcal{T}} = \llbracket \langle \beta \rangle \rrbracket^{\mathcal{T}}$  for all data trees  $\mathcal{T}$ .

The fragment of **downward XPath** is denoted XPath<sub>=</sub>( $\downarrow$ ) and **vertical XPath** is denoted XPath<sub>=</sub>( $\uparrow\downarrow$ ).

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 maximum number of nested axes or the maximum 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. In the sequel we will see that in some situations also expressions of the form  $[\psi]$  can be sometimes avoided (§3.2.1), although not when we only have downward axes (Lemma 10).

## 2.1 Translating to First-Order Logic

In this section we show that there is a truth-preserving translation from  $XPath_{=}(↑↓)$  to first-order logic over an appropriate signature. To do so, we first must interpret data trees into relational structures, and we do this in the most standard way: using a binary ‘child’ relation, an equivalence relation for testing data equality, and monadic relations to test for labels. 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

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

We can now give the translation from XPath to first-order logic over  $\sigma$ . The translation function is indexed by the free variables of the formula produced—one for node expressions, and two for path expressions.

$$\begin{aligned} \text{Tr}_x(a) &= P_a(x) && (a \in \mathbb{A}) \\ \text{Tr}_x(\varphi \dagger \psi) &= \text{Tr}_x(\varphi) \dagger \text{Tr}_x(\psi) && (\dagger \in \{\wedge, \vee\}) \\ \text{Tr}_x(\neg\varphi) &= \neg\text{Tr}_x(\varphi) \\ \text{Tr}_x(\langle\alpha\rangle) &= (\exists y)\text{Tr}_{x,y}(\alpha) && (y \text{ a new variable}) \\ \text{Tr}_x(\langle\alpha = \beta\rangle) &= (\exists y)(\exists z)(y \approx z \wedge \text{Tr}_{x,y}(\alpha) \wedge \text{Tr}_{x,z}(\beta)) && (y, z \text{ new variables}) \\ \text{Tr}_x(\langle\alpha \neq \beta\rangle) &= (\exists y)(\exists z)(y \not\approx z \wedge \text{Tr}_{x,y}(\alpha) \wedge \text{Tr}_{x,z}(\beta)) && (y, z \text{ new variables}) \\ \text{Tr}_{x,y}(\epsilon) &= (x = y) \\ \text{Tr}_{x,y}(\downarrow) &= (x \rightsquigarrow y) \\ \text{Tr}_{x,y}(\uparrow) &= (y \rightsquigarrow x) \\ \text{Tr}_{x,y}(\alpha\beta) &= (\exists z)(\text{Tr}_{x,z}(\alpha) \wedge \text{Tr}_{z,y}(\beta)) && (z \text{ a new variable}) \\ \text{Tr}_{x,y}(\alpha \cup \beta) &= \text{Tr}_{x,y}(\alpha) \vee \text{Tr}_{x,y}(\beta) \\ \text{Tr}_{x,y}([\varphi]) &= \text{Tr}_x(\varphi) \wedge (x = y). \end{aligned}$$

### Proposition 1.

1. If  $\varphi$  is an  $XPath_{=}(↑↓)$  node expression, then  $u \in \llbracket \varphi \rrbracket^{\mathcal{T}}$  iff  $\mathcal{T} \models \text{Tr}_x(\varphi)(u)$ .
2. If  $\alpha$  is an  $XPath_{=}(↑↓)$  path expression, then  $(u, v) \in \llbracket \alpha \rrbracket^{\mathcal{T}}$  iff  $\mathcal{T} \models \text{Tr}_{x,y}(\alpha)(u, v)$ .

*Proof.* The proof is by structural induction on  $\varphi$  and  $\alpha$ . □

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 for  $\odot \in \{=, \neq\}$  we have

$$\begin{aligned} \text{qr}(\text{Tr}_x(\langle\alpha\rangle)) &= 1 + \text{qr}(\text{Tr}_x(\alpha)) \\ \text{qr}(\text{Tr}_x(\langle\alpha \odot \beta\rangle)) &= 2 + \max(\text{qr}(\text{Tr}_x(\alpha)), \text{qr}(\text{Tr}_x(\beta))) \\ \text{qr}(\text{Tr}_{x,y}(\alpha\beta)) &= 1 + \max(\text{qr}(\text{Tr}_{x,y}(\alpha)), \text{qr}(\text{Tr}_{x,y}(\beta))) \\ \text{qr}(\text{Tr}_{x,y}([\phi])) &= \text{qr}(\text{Tr}_x(\varphi)). \end{aligned}$$



### 3. Bisimulation

In this section we define notions of bisimulation for the downward and vertical fragments of XPath, and we show that they coincide with the corresponding logical equivalence relation. For the case of vertical XPath the bisimulation notion relies on a normal form which we introduce for this purpose.

#### 3.1 Downward XPath

We write  $\text{dd}(\varphi)$  to denote the **downward depth** of  $\varphi$  defined as follows:

$$\begin{array}{ll}
 \text{dd}(a) &= 0 & \text{dd}(\lambda) &= 0 \\
 \text{dd}(\varphi \wedge \psi) &= \max\{\text{dd}(\varphi), \text{dd}(\psi)\} & \text{dd}(\varepsilon\alpha) &= \text{dd}(\alpha) \\
 \text{dd}(\neg\varphi) &= \text{dd}(\varphi) & \text{dd}([\varphi]\alpha) &= \max\{\text{dd}(\varphi), \text{dd}(\alpha)\} \\
 \text{dd}(\langle\alpha\rangle) &= \text{dd}(\alpha) & \text{dd}(\downarrow\alpha) &= 1 + \text{dd}(\alpha) \\
 \text{dd}(\langle\alpha \odot \beta\rangle) &= \max\{\text{dd}(\alpha), \text{dd}(\beta)\} & &
 \end{array}$$

where  $a \in \mathbb{A}$ ,  $\odot \in \{=, \neq\}$ , and  $\alpha$  is any path expression or the empty string  $\lambda$ . Let  $\ell\text{-XPath}_{=, \downarrow}$  be the fragment of  $\text{XPath}_{=, \downarrow}$  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>=, \downarrow</sub>** (notation:  $\mathcal{T}, u \equiv \mathcal{T}', u'$ ) iff for all node expression  $\varphi \in \text{XPath}_{=, \downarrow}$ , 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>=, \downarrow</sub>** (notation:  $\mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'$ ) iff for all node expression  $\varphi \in \ell\text{-XPath}_{=, \downarrow}$ , we have  $\mathcal{T}, u \models \varphi$  iff  $\mathcal{T}', u' \models \varphi$ .<sup>1</sup>

We will first show that, for every  $\ell$ , there are finitely many different formulas  $\varphi$  of  $\text{dd}(\varphi) \leq \ell$  up to logical equivalence. Formally, this is usually referred by saying that  $\equiv_{\ell}^{\downarrow}$  has **finite index**.

**Proposition 2.**  $\equiv_{\ell}^{\downarrow}$  has finite index.

*Proof.* It can be easily shown by induction that for any node expression  $\varphi \in \ell\text{-XPath}_{=, \downarrow}$  with unnecessary uses of  $\varepsilon$  (recall that  $\alpha\varepsilon\beta \equiv \alpha\beta$ ) we have that  $\text{qr}(\text{Tr}_x(\varphi))$  is bounded. It is a well-known result of first order that there are finitely many nonequivalent formulas of bounded quantifier rank. Hence there are finitely many nonequivalent node expressions of bounded downward depth.  $\square$

**Corollary 3.**  $\{\mathcal{T}', u' \mid \mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'\}$  is definable by a node expression  $\chi_{\ell, \mathcal{T}, u}$  of  $\ell\text{-XPath}_{=, \downarrow}$ .

*Proof.* Consider the conjunction of all  $\ell\text{-XPath}_{=, \downarrow}$  formulas  $\varphi$  such that  $\mathcal{T}, u \models \varphi$ . By Proposition 2, 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.  $\square$

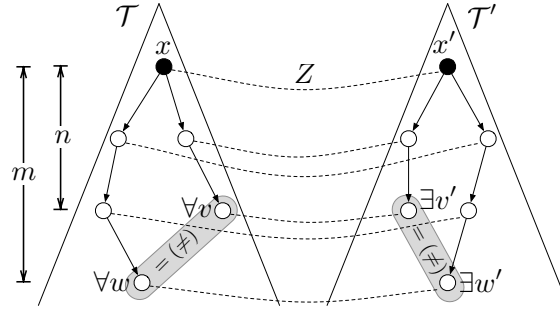
1. Two pointed data trees are ‘equivalent’ when they are ‘indistinguishable’ by the formulas of a given logic. We adopted the terminology of the literature, and so we used the first word, and not the second. The reader should not confuse this notion with that of ‘equivalent’ formulas.

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  $\mathbf{XPath}_{=}$**  ( $\downarrow$ ) (notation:  $\mathcal{T}, u \stackrel{\downarrow}{\leftrightarrow} \mathcal{T}', u'$ ) if 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.  $\text{data}(v) = \text{data}(w) \Leftrightarrow \text{data}(v') = \text{data}(w')$ ,
  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  $0 \leq i < m$ .

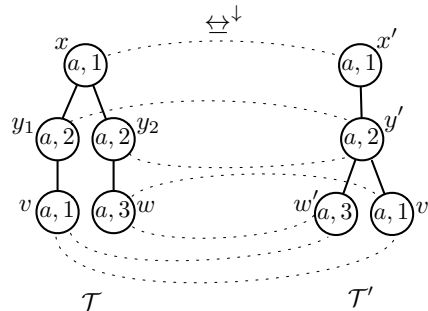
The following picture illustrates the intended requirements.



- **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.

This bisimulation generalizes the classical bisimulation relation (Sangiorgi, 2009), in which the Zig is simply: If  $xZx'$ ,  $x \rightarrow v$ , there is  $v'$  such that  $x' \rightarrow v'$  and  $vZv'$  (the Zag being symmetrical). In fact, if we restrict our Zig condition to having  $n = m = 1$  and  $v = w$ , and the Zag condition to  $n = m = 1$  and  $v' = w'$ , we obtain a relation that corresponds, precisely, to the classical bisimulation relation. Thus, being bisimilar for the notion we define implies being bisimilar for the classical notion.

For example, the dotted lines in the following two data trees represent a bisimulation for  $\mathbf{XPath}_{=}$  ( $\downarrow$ ).



This notion of bisimulation can be seen, as usual, as an Ehrenfeucht-Fraïssé game, where Spoiler tries to find a difference (through the logic’s glasses) between nodes  $u$  and  $u'$ , while Duplicator tries to copy him, showing that  $u$  and  $u'$  are indistinguishable. To gain intuition on our notion of bisimulation, we will briefly explain its associated game (without going into the technical details). The board consists of the data trees  $\mathcal{T}$  and  $\mathcal{T}'$ , and there are two pebbles  $p$  and  $p'$ , so that all along the game,  $p$  will always be over a node of  $\mathcal{T}$ , and  $p'$  will always be over a node of  $\mathcal{T}'$ . Initially,  $p$  is over  $u$ , and  $p'$  is over  $u'$ . If  $u$  and  $u'$  do not satisfy the same label then Spoiler is declared winner and the game finishes. The game proceeds by rounds. Suppose that at some round, the pebbles  $p$  and  $p'$  are in positions  $x$  and  $x'$  satisfying the same label. A *round* consists of one move of Spoiler, followed by an answer by Duplicator, and a final decision made by spoiler.

- Step 1. Spoiler’s first move: he chooses a pebble, two integers  $n$  and  $m$ , and two paths, one of length  $m$  and the other of length  $n$ . These paths start in  $x$  if he chose  $p$  or in  $x'$  if he chose  $p'$ . Suppose he chose  $p$  (the case for  $p'$  is analogous but over data tree  $\mathcal{T}'$  instead of  $\mathcal{T}$ ). His first path is represented by some  $w \in T$  such that  $x \xrightarrow{m} w$ , and his second path is represented by some  $v \in T$  such that  $x \xrightarrow{n} v$ .
- Step 2. Duplicator’s answer: he shows two paths in  $\mathcal{T}'$  (or in  $\mathcal{T}$  in case duplicator had chosen  $p'$  instead of  $p$ ), one of length  $m$  and the other of length  $n$ , both starting on  $x'$ , such that  $data(u) = data(v)$  iff  $data(u') = data(v')$ . If there are no such paths, then Spoiler is declared the winner and the game finishes.
- Step 3. Final move by Spoiler: he chooses

- either  $i \in \{0 \dots n - 1\}$ , and places pebble  $p$  over  $(\xrightarrow{i}v)$  and pebble  $p'$  over  $(\xrightarrow{i}v')$
- or  $i \in \{0 \dots m - 1\}$ , and places pebble  $p$  on  $(\xrightarrow{i}w)$  and pebble  $p'$  over  $(\xrightarrow{i}w')$

If now the pebbles are over two nodes which satisfy the same label, the game proceeds. Else, Spoiler is declared the winner and the game finishes.

Duplicator wins if spoiler is never declared winner in a game at infinitely many rounds. The resemblance of the game rules with **Harmony**, **Zig** and **Zag** is evident. One can see that spoiler has a winning strategy in the game whose initial pebbles are placed over  $x$  and  $x'$  if and only if  $\mathcal{T}, u \leftrightarrow^{\downarrow} \mathcal{T}', u'$ .

It is interesting to compare this game to the one for capturing *bisimulation* for the basic modal logic. In the latter, Spoiler just chooses a successor of  $x$  (or  $x'$ ) following the accessibility relation. This is also the case for Core-XPath or PDL. But in our case, because of adding comparison of data values, Spoiler has to choose a *whole* path (in fact, two paths), making the game ‘less local’ than the one for the basic modal logic (or Core-XPath or PDL). An analogous view for the game of XPath<sub>=</sub>( $\downarrow$ ) would be to allow Spoiler to build his paths in a step-by-step fashion, that is, extending the paths constructed so far with a new *single* node, thus giving Duplicator less certainty. Is it possible to change the rules of the game so that in the second step of the round Spoiler builds his paths in a step-by-step fashion? No. Even if Spoiler has the freedom to extend step-by-step *only one* of his paths, this would be too unfair for Duplicator. Indeed, consider our last example of bisimilar  $\mathcal{T}, x$  and  $\mathcal{T}', x'$ . Spoiler would initially declare that his first path will be  $x'$  of length 0, and so

Duplicator can only define his first path to be  $x$  of length 0. Then Spoiler will construct, step-by-step, his second path. Initially he shows the path  $x' \rightarrow y'$  to Duplicator. Duplicator has two possible answers: either the path  $x \rightarrow y_1$  or  $x \rightarrow y_2$ . Suppose he chooses  $x \rightarrow y_1$  (the case for  $x \rightarrow y_2$  is analogous). Now Spoiler extends his path  $x' \rightarrow y'$  to  $x' \rightarrow y' \rightarrow w'$ . Since  $y_1$  has only one child, Spoiler has one possible answer for extending his own path:  $x \rightarrow y_1 \rightarrow v$ . Next, Spoiler declares that he is done: he has constructed two paths starting in  $x'$ , one is of length 0 and the other one is  $x' \rightarrow y' \rightarrow w'$ , of length 2. Duplicator has constructed one path of length 0 and the other as  $x \rightarrow y_1 \rightarrow v$ . Duplicator loses, as  $\text{data}(x') \neq \text{data}(w')$  but  $\text{data}(x) = \text{data}(v)$ . This example shows that the game fairness is established when spoiler tells duplicator *in advance* which are the two paths he chooses.

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:

**Observation 4.**  $\mathcal{T}, u \leftrightarrow^{\downarrow} (\mathcal{T}|u), u$ .

**Observation 5.** 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

- **Harmony:** If  $xZ_jx'$  then  $\text{label}(x) = \text{label}(x')$ .
- **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.  $\text{data}(v) = \text{data}(w) \Leftrightarrow \text{data}(v') = \text{data}(w')$ ,
  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.

This notion of  $\ell$ -bisimulation corresponds to the game between spoiler and duplicator at  $\ell$  rounds.

Clearly if  $\mathcal{T}, u \leftrightarrow^{\downarrow} \mathcal{T}', u'$  then  $\mathcal{T}, u \leftrightarrow_{\ell}^{\downarrow} \mathcal{T}', u'$  for all  $\ell$ .

**Observation 6.** 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  $\mathcal{T}, u \leftrightarrow^{\downarrow} \mathcal{T}', u'$ .

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

**Observation 7.**  $\mathcal{T}, u \leftrightarrow_{\ell}^{\downarrow} (\mathcal{T}|_{\ell}u), u$ .

## 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 8.**

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 above theorem will be shown as a consequence of Propositions 9 and 11:

**Proposition 9.**  $\mathcal{T}, u \Leftrightarrow_\ell^\downarrow \mathcal{T}', u'$  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  $\text{dd}(\varphi) \leq j$ , and for all  $\alpha$  with  $\text{dd}(\alpha) \leq j$ :

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,  $\text{label}(x) = \text{label}(x')$  and then  $\mathcal{T}, x \models \varphi$  iff  $\mathcal{T}', x' \models \varphi$ . The Boolean cases for  $\varphi$  are straightforward.

Suppose  $\varphi = \langle \alpha = \beta \rangle$ . 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  $\text{data}(v) = \text{data}(w)$ . 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  $\text{data}(v') = \text{data}(w')$ . 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  $x_0, \dots, x_n \in T$  and  $x'_0, \dots, x'_n \in T'$  be such that

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

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

**Lemma 10.** *Let  $\alpha$  be a path expression of  $X\text{Path}_=(\downarrow \downarrow_*)$ . Let  $x \xrightarrow{n} v$  and  $x' \xrightarrow{n} v'$  such that  $(x, v) \in \llbracket \alpha \rrbracket^{\mathcal{T}}$  and  $(x', v') \notin \llbracket \alpha \rrbracket^{\mathcal{T}'}$ . Then there is a subformula  $\varphi$  of  $\alpha$  and  $k \in \{0, \dots, n\}$  such that  $\mathcal{T}, (\xrightarrow{k}v) \models \varphi$  and  $\mathcal{T}', (\xrightarrow{k}v') \not\models \varphi$ .*

*Proof of Lemma.* Let  $x = v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_n = v$  and  $x' = v'_0 \rightarrow v'_1 \rightarrow \dots \rightarrow v'_n = v'$ . 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.  $\blacksquare$

But this contradicts the inductive hypothesis 1.  $\square$

**Proposition 11.**  $\mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'$  implies  $\mathcal{T}, u \Leftrightarrow_{\ell}^{\downarrow} \mathcal{T}', u'$ .

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

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

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

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

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 node expressions,

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

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:

- Assume (a) 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 (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

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

However these conjunctions could be potentially infinite if trees are infinitely branching. Since  $\text{dd}(\varphi_{v',w'}^i) \leq h - i$ , by Proposition 2 there are finitely many non-equivalent node expressions  $\varphi_{v',w'}^i$  and 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  $\alpha = [\Phi^0] \downarrow [\Phi^1] \downarrow \dots \downarrow [\Phi^n]$  and  $\beta = [\Psi^0] \downarrow [\Psi^1] \downarrow \dots \downarrow [\Psi^m]$ . By construction,  $\text{dd}(\alpha), \text{dd}(\beta) \leq h$  and so  $\text{dd}(\langle \alpha = \beta \rangle) \leq h$ .

It is clear that by construction  $(x, v) \in \llbracket \alpha \rrbracket^{\mathcal{T}}$  and  $(x, w) \in \llbracket \beta \rrbracket^{\mathcal{T}}$ , and therefore  $\mathcal{T}, x \models \langle \alpha = \beta \rangle$ . We see next that  $\mathcal{T}', x' \not\models \langle \alpha = \beta \rangle$ . This contradicts  $\mathcal{T}, x \equiv_h^{\downarrow} \mathcal{T}', x'$ , and hence we are done. Suppose that  $\mathcal{T}', x' \models \langle \alpha = \beta \rangle$ . Then there is  $(v', w') \in P$  such that  $(x', v') \in \llbracket \alpha \rrbracket^{\mathcal{T}'}$  and  $(x', w') \in \llbracket \beta \rrbracket^{\mathcal{T}'}$ . In particular, i and ii are true, and then either (a) or (b) hold. In the first case, we have, by construction, that  $(x', v') \notin \llbracket \alpha \rrbracket^{\mathcal{T}'}$ , and in the second it is clear that  $(x', w') \notin \llbracket \beta \rrbracket^{\mathcal{T}'}$ . In either case, we arrive to a contradiction.  $\square$

*Proof of Theorem 8.* Item 2 is a direct consequence of Propositions 9 and 11.

The left-to-right argument for item 1 can be seen as a consequence of item 2. Indeed,  $\mathcal{T}, u \leftrightarrow_{\ell}^{\downarrow} \mathcal{T}', u'$  implies  $\mathcal{T}, u \leftrightarrow_{\ell}^{\downarrow} \mathcal{T}', u'$  for all  $\ell$ , which by item 2 implies  $\mathcal{T}, u \equiv_{\ell}^{\downarrow} \mathcal{T}', u'$  for all  $\ell$ , which in turn entails  $\mathcal{T}, u \equiv^{\downarrow} \mathcal{T}', u'$ .

The right-to-left argument for item 1 is similar to that of Proposition 11, but defining  $Z$  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 Zag is satisfied).  $\square$

### 3.2 Vertical XPath

We now study bisimulation for  $\text{XPath}_=(\uparrow\downarrow)$ . 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  $\text{XPath}_=$  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  $\text{XPath}_=$ , 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 (notation:  $\text{vd}(\varphi)$ ). The **nesting depth** of a formula  $\varphi$  (notation:  $\text{nd}(\varphi)$ ) is the maximum number of nested  $[ ]$  appearing in  $\varphi$ .

$$\begin{array}{ll}
 \text{vd}(a) = (0, 0) & \text{vd}(\lambda) = (0, 0) \\
 \text{vd}(\varphi \wedge \psi) = \max\{\text{vd}(\varphi), \text{vd}(\psi)\} & \text{vd}(\varepsilon\alpha) = \text{vd}(\alpha) \\
 \text{vd}(\neg\varphi) = \text{vd}(\varphi) & \text{vd}([\varphi]\alpha) = \max\{\text{vd}(\varphi), \text{vd}(\alpha)\} \\
 \text{vd}(\langle\alpha\rangle) = \text{vd}(\alpha) & \text{vd}(\downarrow\alpha) = \max\{(0, 0), \text{vd}(\alpha) + (1, -1)\} \\
 \text{vd}(\langle\alpha \odot \beta\rangle) = \max\{\text{vd}(\alpha), \text{vd}(\beta)\} & \text{vd}(\uparrow\alpha) = \max\{(0, 0), \text{vd}(\alpha) + (-1, 1)\} \\
 \\
 \text{nd}(a) = 0 & \text{nd}(\alpha\beta) = \max\{\text{nd}(\alpha), \text{nd}(\beta)\} \\
 \text{nd}(\varphi \wedge \psi) = \max\{\text{nd}(\varphi), \text{nd}(\psi)\} & \text{nd}(\varepsilon) = 0 \\
 \text{nd}(\neg\varphi) = \text{nd}(\varphi) & \text{nd}([\varphi]) = 1 + \text{nd}(\varphi) \\
 \text{nd}(\langle\alpha\rangle) = \text{nd}(\alpha) & \text{nd}(\downarrow) = 0 \\
 \text{nd}(\langle\alpha \odot \beta\rangle) = \max\{\text{nd}(\alpha), \text{nd}(\beta)\} & \text{nd}(\uparrow) = 0
 \end{array}$$

where,  $a \in \mathbb{A}$ ,  $\odot \in \{=, \neq\}$ , ‘+’ and ‘max’ are performed component-wise, and  $\alpha$  is any path expression or the empty string  $\lambda$ .

Let  $(r, s, k)$ - $\text{XPath}_=(\uparrow\downarrow)$  be the set of formulas  $\varphi$  in  $\text{XPath}_=(\uparrow\downarrow)$  with  $\text{vd}(\varphi) \leq (r, s)$  and  $\text{nd}(\varphi) \leq k$ . Let  $\mathcal{T}, u$  and  $\mathcal{T}', u'$  be pointed data trees. We say that  $\mathcal{T}, u$  and  $\mathcal{T}', u'$  are **equivalent for  $\text{XPath}_=(\uparrow\downarrow)$**  (notation:  $\mathcal{T}, u \equiv^{\downarrow\uparrow} \mathcal{T}', u'$ ) iff for all  $\varphi \in \text{XPath}_=(\uparrow\downarrow)$ , 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  $\text{XPath}_=(\uparrow\downarrow)$  (notation:  $\mathcal{T}, x \equiv_{r,s}^{\downarrow\uparrow} \mathcal{T}', x'$  [resp.  $\mathcal{T}, x \equiv_{r,s,k}^{\downarrow\uparrow} \mathcal{T}', x'$ ]) if they satisfy the same node expressions  $\varphi$  of  $\text{XPath}_=(\uparrow\downarrow)$  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 normal form for  $\text{XPath}_=(\uparrow\downarrow)$  that will be implicitly used in the definition of bisimulation in this 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  $\text{XPath}_=(\uparrow\downarrow)$  is **downward** [resp. **upward**] if it is of the form  $\downarrow^n[\varphi]$  [resp.  $[\varphi]\uparrow^n$ ] for some  $n > 0$  with  $\varphi \in \text{XPath}_=(\uparrow\downarrow)$ . For example,  $\downarrow[\langle\uparrow\rangle]$  is a downward expression whereas  $\downarrow[\langle\downarrow\rangle]\downarrow$  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  $\text{XPath}_=(\uparrow\downarrow)$  formula 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  $\odot \in \{=, \neq\}$ . Next, we show that for every  $\text{XPath}_=(\uparrow\downarrow)$  formula there is an equivalent



one in up-down normal form. Of course, the idea of replacing child axes with parent axes is by no means novel, and there have been a number of works on rewriting expressions into equivalent ones to improve performance or streamability, such as the works of Olteanu (2007) and Olteanu, Meuss, Furche, and Bry (2002). However, these rewrite systems aim at removing backward axes (parent, ancestor, etc) maintaining equivalence at the root, different to our end, which is to maintain equivalence of path expressions over any pair of nodes. Further, we doubt our normal form can be useful in such scenarios since from a computational or streamability perspective the resulting formula in up-down normal form seems more complex. Although being more ‘regular’, it has more nested modalities and introduces parent axes, as it will be shown later. Our motivation for this normal form is to simplify the definition of bisimulation by means of rendering the formula as simple as possible in form.

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,  $\text{nav}(\downarrow[\langle \uparrow \rangle] \downarrow[\langle \downarrow = \uparrow \rangle] \uparrow[b]) = \downarrow \downarrow \uparrow$ .

**Proposition 12.** *Let  $\varphi \in (r, s, k)\text{-XPath}_=(\uparrow \downarrow)$ , then there is  $\varphi^{\uparrow \downarrow} \in \text{XPath}_=(\uparrow \downarrow)$  in up-down normal form such that the following hold*

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 test 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  $[\langle \uparrow[\langle \downarrow[a] \rangle] \rangle] \uparrow \downarrow$ , which is in up-down normal form.

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

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

In the expressions above, each  $\xi_i$  is the empty string, or it is 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,  $\downarrow[-a] \downarrow[-b] \uparrow[-c] \uparrow \equiv^{\downarrow \uparrow} [\langle \downarrow[-a] \downarrow[-b] \downarrow[-c] \rangle]$ . 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  $[a] \downarrow[b] \downarrow \equiv^{\downarrow \uparrow} \downarrow \downarrow[\langle \uparrow[b] \uparrow[a] \rangle]$  and  $\uparrow[a] \uparrow[b] \equiv^{\downarrow \uparrow} [\langle \uparrow[a] \uparrow[b] \rangle] \uparrow \uparrow$ . It is thus clear that the left and right expressions above are semantically equivalent.

The following lemma treats the case of path expressions:

**Lemma 13.** *Let  $\alpha$  be an  $XPath_{=}(↑↓)$ -path expression with  $\text{vd}(\alpha) = (r, s)$  and  $\text{nd}(\alpha) = k$ , then there is an up-down path expression  $\alpha^{\updownarrow}$  such that:*

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

*Proof of Lemma.* 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  $\downarrow \uparrow$ . Let

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

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

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

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  $\text{vd}(\gamma_{\uparrow} \alpha_2 \gamma_{\downarrow}) = (r, s)$ ,  $\text{nd}(\alpha_2) \leq \text{nd}(\alpha_1) + 1$ , 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.  $\alpha_f \equiv \downarrow \uparrow \alpha_1$ ,
2.  $\text{vd}(\gamma \uparrow \alpha_f \gamma \downarrow) = (r, s)$ , and
3.  $\text{nd}(\alpha_f) \leq \text{nd}(\alpha_1) + r_1$ .

After applying  $(\varepsilon)$  and  $(merge)$  to  $\gamma \uparrow \alpha_f \gamma \downarrow$  as many times as possible, we obtain an equivalent  $\alpha'$ , of the same vertical and nesting depth as  $\gamma \uparrow \alpha_f \gamma \downarrow$ , of the form

$$\alpha' = \chi_1 \uparrow \chi_2 \uparrow \dots \uparrow \chi_n \downarrow \chi_{n+1} \downarrow \chi_{n+2} \downarrow \dots \downarrow \chi_{n+m},$$

where  $\chi_i$  is a (possibly empty) string of the form  $[\varphi_i^1] \dots [\varphi_i^{n_i}]$ . We apply  $(shift-l)$  and  $(shift-r)$  to  $\alpha'$  to obtain an equivalent  $\alpha''$  of the same vertical depth (i.e.  $\text{vd}(\alpha'') = \text{vd}(\alpha') = (r, s)$ ) and of nesting depth equal to  $\text{nd}(\alpha') + 1$  of the form

$$\alpha^{\updownarrow} = \chi' \uparrow^n \downarrow^m \chi'',$$

where  $\chi'$  and  $\chi''$  are the empty string or of the form  $[\varphi]$ . Observe that  $\text{nd}(\alpha^{\updownarrow}) = \text{nd}(\alpha') + 1 \leq k + r_1 + 1 \leq r + s + 1$ , and so  $\alpha^{\updownarrow}$  satisfies all requirements of the lemma. This concludes the proof of Lemma 13.  $\blacksquare$

The following lemma treats the case of data tests:

**Lemma 14.** *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:*

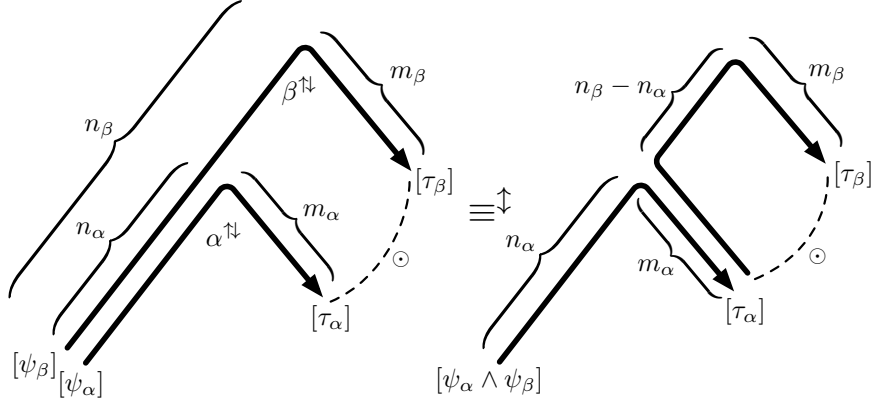
1.  $\langle \gamma^{\updownarrow} \rangle \equiv \downarrow \uparrow \varphi$ ,
2.  $\text{vd}(\gamma^{\updownarrow}) = (r, s)$ , and
3.  $\text{nd}(\gamma^{\updownarrow}) \leq k + 1$ .

*Proof of Lemma.* Let us analyse the case where  $\alpha^{\updownarrow} = [\psi_\alpha] \uparrow^{n_\alpha} \downarrow^{m_\alpha} [\tau_\alpha]$  and  $\beta^{\updownarrow} = [\psi_\beta] \uparrow^{n_\beta} \downarrow^{m_\beta} [\tau_\beta]$ , where  $n_\alpha + m_\alpha > 0$ ,  $n_\beta + m_\beta > 0$ , and  $\psi_\alpha, \psi_\beta, \tau_\alpha, \tau_\beta$  are in up-down normal form (the remaining cases where  $\alpha^{\updownarrow} = \varepsilon$  or  $\beta^{\updownarrow} = \varepsilon$  being simpler). Suppose, without loss of generality, that  $n_\alpha \leq n_\beta$ .

Hence, we have  $\langle \alpha^{\updownarrow} \odot \beta^{\updownarrow} \rangle \equiv \downarrow \uparrow \langle \gamma^{\updownarrow} \rangle$ , 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.* the picture below).



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 14.  $\blacksquare$

By induction on  $\varphi$ , and using Lemmas 13 and 14, one can show that there is  $\varphi^{\updownarrow}$  as desired.  $\square$

### 3.2.2 FINITE INDEX

Contrary to the case of  $\text{XPath}_{=}\langle\downarrow\rangle$  (*cf.*, Proposition 2), the logical equivalence relation restricted to  $\text{XPath}_{=}\langle\updownarrow\rangle$ -formulas of bounded vertical depth has infinitely many equivalence classes.

**Proposition 15.** *If  $r + s \geq 2$  then  $\equiv_{r,s}^{\updownarrow}$  has infinite index.*

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

Consider the following formulas.

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

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 2.

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  $(r, s)$ .  $\square$

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

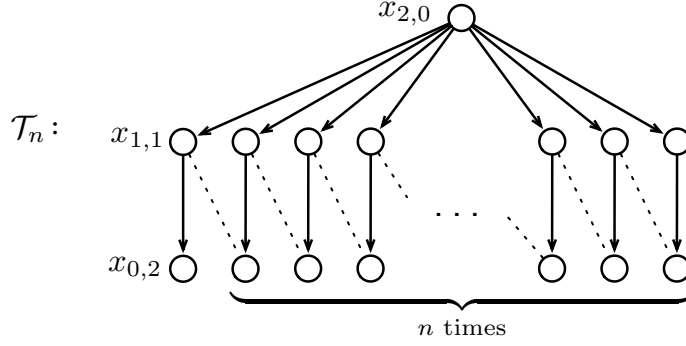


Figure 2: Model verifying  $\psi_i^j$  for all  $i \geq n$  and not verifying  $\varphi_l$  for no  $l < n$ . Dotted lines represent equal data values.

**Proposition 16.**  $\equiv_{r,s,k}^{\downarrow\uparrow}$  has finite index.

*Proof.* Same argument as in the proof of Proposition 2.  $\square$

**Corollary 17.**  $\{\mathcal{T}', u' \mid \mathcal{T}, u \equiv_{r,s,k}^{\downarrow\uparrow} \mathcal{T}', u'\}$  is definable by a node expression of  $(r, s, k)$ - $XPath_{=}$ ( $\uparrow\downarrow$ ).

*Proof.* Similar to the proof of Corollary 3.  $\square$

### 3.2.3 BISIMULATION AND $(r, s, k)$ -BISIMULATION

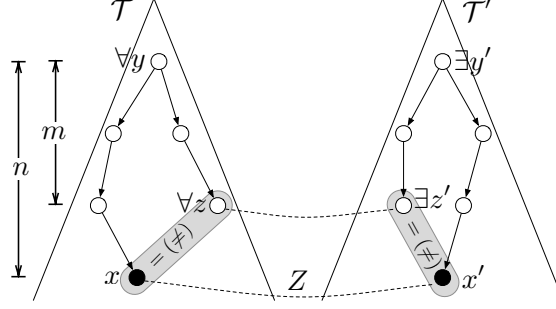
The advantage of the normal form presented in Section 3.2.1 is that it makes it possible to use a very simple notion of bisimulation. The disadvantage is that, since it does not preserve nesting depth,  $\leftrightarrow_{r,s,k}^{\downarrow\uparrow}$  does not correspond precisely to  $\equiv_{r,s,k}^{\downarrow\uparrow}$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

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

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

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

The following picture illustrates the intended requirements



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

The above definitions heavily rely on the normal form of Proposition 12. In fact, the normal form is not strictly necessary for giving a notion of bisimulation for the vertical fragment. For the case of the downward fragment, no normal form was used, as every path expression is essentially a repetition of node test and ‘child’ relation. On the contrary, for the case of the vertical fragment, the use of a normal form would be very beneficial. A notion of bisimulation not taking into account the existence of normal form would have a rule **Zig** of the form: If  $xZx'$  and  $n_1, \dots, n_k, m_1, \dots, m_k, \tilde{n}_1, \dots, \tilde{n}_{\tilde{k}}$  and  $\tilde{m}_1, \dots, \tilde{m}_{\tilde{k}}$  are such that

- $v_1^i \rightarrow \dots \rightarrow v_{n_i}^i$  in  $\mathcal{T}$  and  $w_1^i \rightarrow \dots \rightarrow w_{m_i}^i$  in  $\mathcal{T}$  for  $i \in \{1, \dots, k\}$ ,
- $\tilde{v}_1^i \rightarrow \dots \rightarrow \tilde{v}_{\tilde{n}_i}^i$  in  $\mathcal{T}$  and  $\tilde{w}_1^i \rightarrow \dots \rightarrow \tilde{w}_{\tilde{m}_i}^i$  in  $\mathcal{T}$  for  $i \in \{1, \dots, \tilde{k}\}$ ,
- $v_1^i = w_1^i$  and  $\tilde{v}_1^i = \tilde{w}_1^i$
- $w_{m_i}^i = v_{n_i+1}^{i+1}$  and  $\tilde{w}_{\tilde{m}_i}^i = \tilde{v}_{\tilde{n}_i+1}^{i+1}$
- $x = v_{n_1}^1$  and  $x = \tilde{v}_{\tilde{n}_1}^1$

then there are  $v_1^i, \dots, v_{n_i}^i$  in  $\mathcal{T}'$ ,  $w_1^i, \dots, w_{m_i}^i$  for  $i \in \{1, \dots, k\}$ , and  $\tilde{v}_1^i, \dots, \tilde{v}_{\tilde{n}_i}^i, \tilde{w}_1^i, \dots, \tilde{w}_{\tilde{m}_i}^i$  in  $\mathcal{T}'$  for  $i \in \{1, \dots, \tilde{k}\}$  such that

- $v_1^i \rightarrow \dots \rightarrow v_{n_i}^i$  in  $\mathcal{T}'$  and  $w_1^i \rightarrow \dots \rightarrow w_{m_i}^i$  in  $\mathcal{T}'$ ,
- $\tilde{v}_1^i \rightarrow \dots \rightarrow \tilde{v}_{\tilde{n}_i}^i$  in  $\mathcal{T}'$  and  $\tilde{w}_1^i \rightarrow \dots \rightarrow \tilde{w}_{\tilde{m}_i}^i$  in  $\mathcal{T}'$ ,
- $v_1^i = w_1^i$  and  $\tilde{v}_1^i = \tilde{w}_1^i$
- $w_{m_i}^i = v_{n_i+1}^{i+1}$  and  $\tilde{w}_{\tilde{m}_i}^i = \tilde{v}_{\tilde{n}_i+1}^{i+1}$
- $x = v_{n_1}^1$  and  $x = \tilde{v}_{\tilde{n}_1}^1$
- $v_j^i Z v_j^i, w_j^i Z w_j^i, \tilde{v}_j^i Z \tilde{v}_j^i, \tilde{w}_j^i Z \tilde{w}_j^i$  and
- $\text{data}(w_{n_k}^k) = \text{data}(\tilde{w}_{\tilde{n}_k}^k)$  iff  $\text{data}(w_{\tilde{n}_k}^k) = \text{data}(\tilde{w}_{n_k}^k)$

Intuitively, this definition establishes that for every paths  $p$  and  $\tilde{p}$  going up and down, many times in  $\mathcal{T}$ , there has to be a similar paths  $p'$  and  $\tilde{p}'$  in  $\mathcal{T}$ , going up and down many times, and respecting the ‘shape’ of  $p$  and  $\tilde{p}$  respectively such that the  $j$ -th node of  $p$  is connected to the  $j$ -th node of  $p'$  via  $Z$ , and the same for the nodes along  $\tilde{p}$ , and such that the data values of the last node of  $p$  and  $\tilde{p}$  in  $\mathcal{T}$  are equal if and only if the data values of the last nodes of  $p'$  and  $\tilde{p}'$  are so.

This definition is the ‘natural’ extension of the downward bisimulation, but instead of considering downward paths, it considers more general ones. As it happens with the downward bisimulation, every intermediate node of  $\mathcal{T}$  in both paths have to be related with the corresponding node in  $\mathcal{T}'$ . As one can immediately see, this definition is quite long and checks too many conditions. Working with normal forms allow us to get a much simpler definition of bisimulation, which has two main advantages: *a*) there is only one path in each data tree, and *b*) we do not require intermediate nodes to be related by  $Z$ . These two features are direct consequences of the up-down normal form, as the normal form *a*) only compares values against the root (so there is only one non-empty path expression in each ‘diamond’ node expression) and *b*) it only make tests at the beginning and at the end of each path expressions (but not in intermediate nodes). However, notice that due to the normal form both definitions denote *the same* bisimulation relation.

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

- **Harmony:** If  $xZ_{\hat{r},\hat{s}}^{\hat{k}} x'$  then  $label(x) = label(x')$ .
- **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.  $data(z) = data(x) \Leftrightarrow data(z') = data(x')$ ,
  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 18.** 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

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

**Observation 19.**  $\mathcal{T}, u \stackrel{\uparrow\downarrow}{\leftrightarrow}_{r,s,k} (\mathcal{T}|_r^s u), u$ .

## 3.2.4 EQUIVALENCE AND BISIMULATION

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

**Theorem 20.**

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

The above theorem will be shown as a consequence of the following Propositions 21 and 22.

**Proposition 21.**  $\mathcal{T}, u \Leftrightarrow_{r,s,k \cdot (r+s+2)}^{\downarrow\uparrow} \mathcal{T}', u'$  implies  $\mathcal{T}, u \equiv_{r,s,k}^{\downarrow\uparrow} \mathcal{T}', u'$ .

*Proof.* We show that if  $\mathcal{T}, u \Leftrightarrow_{r,s,k}^{\downarrow\uparrow} \mathcal{T}', u'$  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  $x Z_{\hat{r}, \hat{s}}^{\hat{k}} x'$  then  $\mathcal{T}, x \models \varphi$  iff  $\mathcal{T}', x' \models \varphi$ .
2. If  $y \xrightarrow{n} x, y' \xrightarrow{n} x', x Z_{\hat{r}, \hat{s}}^{k-1} x'$ , then  $(x, y) \in \llbracket \alpha^\uparrow \rrbracket^{\mathcal{T}}$  iff  $(x', y') \in \llbracket \alpha^\uparrow \rrbracket^{\mathcal{T}'}$ .
3. If  $y \xrightarrow{m} z, y' \xrightarrow{m} z', z Z_{\hat{r}', \hat{s}'}^{k-1} z'$  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 12, the main statement follows. We simultaneously show 1, 2 and 3 by induction on  $|\varphi| + |\alpha^\uparrow| + |\alpha^\downarrow|$ .

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  $z Z_{\hat{r}', \hat{s}'}^{k-1} z'$  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  $x_0, \dots, x_n \in T$  and  $x'_0, \dots, x'_n \in T'$  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 18, 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  $z_0, \dots, z_m \in T$  and  $z'_0, \dots, z'_m \in T'$  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 22.**  $\mathcal{T}, u \equiv_{r,s,k}^{\downarrow\uparrow} \mathcal{T}', u'$  implies  $\mathcal{T}, u \Leftrightarrow_{r,s,k}^{\downarrow\uparrow} \mathcal{T}', u'$ .

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

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

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

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

Since  $\mathcal{T}, x \equiv_{r,s,k}^{\downarrow\uparrow} \mathcal{T}', x'$ ,  $\text{vd}(\langle \varepsilon = \uparrow^n \downarrow^m \rangle) \leq (r, s)$ ,  $\text{nd}(\langle \varepsilon = \uparrow^n \downarrow^m \rangle) = 0$ , 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.  $\mathcal{T}, x \equiv_{\hat{r}, \hat{s}, \hat{k}-1}^{\downarrow\uparrow} \mathcal{T}', x'$ , and
- iv.  $\mathcal{T}, z \equiv_{\hat{r}', \hat{s}', \hat{k}-1}^{\downarrow\uparrow} \mathcal{T}', z'$ , 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)  $\mathcal{T}, x \not\equiv_{\hat{r}, \hat{s}, \hat{k}-1}^{\downarrow\uparrow} \mathcal{T}', x'$ ; or
- (b)  $\mathcal{T}, z \not\equiv_{\hat{r}', \hat{s}', \hat{k}-1}^{\downarrow\uparrow} \mathcal{T}', z'$  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 node expressions,  $\varphi_{y', z'}$  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  $\text{vd}(\varphi_{y', z'}) \leq (\hat{r}, \hat{s})$ ,  $\text{nd}(\varphi_{y', z'}) < \hat{k}$ , and such that  $\mathcal{T}, x \models \varphi_{y', z'}$  but  $\mathcal{T}', x' \not\models \varphi_{y', z'}$ ; and let  $\psi_{y', z'} = \top$ .
- 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 16, 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  $\mathcal{T}, x \equiv_{\hat{r}, \hat{s}, \hat{k}}^{\downarrow\uparrow} \mathcal{T}', x'$ .  $\square$

*Proof of Theorem 20.* Items 2 and 3 are shown in Propositions 21 and 22.

The left-to-right argument for item 1 can be seen as a consequence of item 2. Indeed,  $\mathcal{T}, u \leftrightarrow^{\downarrow\uparrow} \mathcal{T}', u'$  implies  $\mathcal{T}, u \leftrightarrow_{(r,s)}^{\downarrow\uparrow} \mathcal{T}', u'$  for all  $r, s$ , which by item 2 implies  $\mathcal{T}, u \equiv_{(r,s)}^{\downarrow\uparrow} \mathcal{T}', u'$  for all  $r, s$ , which in turn entails  $\mathcal{T}, u \equiv^{\downarrow\uparrow} \mathcal{T}', u'$ .

The right-to-left argument for item 1 is similar to that of Proposition 22, 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  $\mathcal{T}, x \equiv^{\downarrow\uparrow} \mathcal{T}', x'$ . 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).  $\square$

**Corollary 23.**  $\leftrightarrow_{r,s,k}^{\downarrow\uparrow}$  has finite index.

*Proof.* Immediate from Theorem 20 and Proposition 16.  $\square$

### 3.3 Simulation

In this section we define notions of directed (non-symmetric) simulations for  $\text{XPath}_=(\downarrow)$  and  $\text{XPath}_=(\uparrow\downarrow)$ , as it is done, e.g., in the works of Kurtonina and de Rijke (1997) or Lutz, Piro, and Wolter (2011) for some modal logics. We obtain results similar to Theorems 8 and 20 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  $\text{XPath}_=(\downarrow)$ ,  $\text{XPath}_=(\uparrow\downarrow)$ ,  $\text{XPath}_=(\downarrow\downarrow^*)$ , or  $\text{XPath}_=(\downarrow\uparrow\downarrow^*\uparrow^*)$ , we write  $\mathcal{L}^+$  for the positive fragment of  $\mathcal{L}$ .

A **simulation for  $\text{XPath}_=(\downarrow)$**  [resp. for  $\text{XPath}_=(\uparrow\downarrow)$ ] 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:  $\mathcal{T}, u \xrightarrow{\downarrow} \mathcal{T}', u'$ ) 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:** 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  $\text{XPath}_=(\uparrow\downarrow)$  (notation:  $\mathcal{T}, u \xrightarrow{\uparrow\downarrow} \mathcal{T}', u'$ ) 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:** 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  $\xrightarrow{\downarrow}_\ell$  and  $\xrightarrow{\uparrow\downarrow}_{r,s,k}$  are defined accordingly. We define one-way (non-symmetric) logical implication between models as follows. We write  $\mathcal{T}, u \Rightarrow^\downarrow \mathcal{T}', u'$  for

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

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

#### Theorem 24.

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

4.  $\mathcal{T}, u \rightleftharpoons_{r,s,k}^{\downarrow\uparrow} \mathcal{T}', u'$  implies  $\mathcal{T}, u \rightrightarrows_{r,s,k}^{\downarrow\uparrow} \mathcal{T}', u'$ .

*Proof.* The proofs are straightforward adaptations of the proofs of Propositions 9 and 11 and Propositions 21 and 22 respectively, and are omitted here. In particular, for the ‘if’ part, in the adaptation of the proofs of Propositions 11 and 22, the simulations are defined by

$$\begin{aligned} xZ_i x' &\text{ iff } \mathcal{T}, x \rightleftharpoons_i^{\downarrow} \mathcal{T}', x \\ xZ_{\hat{r},\hat{s}}^k x' &\text{ iff } \mathcal{T}, x \rightleftharpoons_{\hat{r},\hat{s},k}^{\downarrow\uparrow} \mathcal{T}', x \end{aligned}$$

respectively, and the conditions (a) and (b) on page 13 become now

(a)  $[\exists i \in \{0, \dots, n\} \exists \varphi \in \text{XPath}_=(\downarrow)^+] \text{dd}(\varphi) \leq h - i \wedge \mathcal{T}, v_i \models \varphi \wedge \mathcal{T}', v'_i \not\models \varphi$ ; or

(b)  $[\exists j \in \{0, \dots, m\} \exists \varphi \in \text{XPath}_=(\downarrow)^+] \text{dd}(\varphi) \leq h - j \wedge \mathcal{T}, w_j \models \varphi \wedge \mathcal{T}', w'_j \not\models \varphi$ ,

and

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

(b)  $[\exists j \in \{0, \dots, m\} \exists \varphi \in \text{XPath}_=(\downarrow\uparrow)^+] \text{vd}(\varphi) \leq (\hat{r} + j', \hat{s} - j')$  for  $j' = n - m + j$   
 $\wedge \text{nd}(\varphi) \leq k - 1 \wedge \mathcal{T}, w_j \models \varphi \wedge \mathcal{T}', w'_j \not\models \varphi$

respectively. □

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

**Lemma 25.** *If  $\mathcal{T}'$  is a substructure of  $\mathcal{T}$  and  $u' \in T'$  then  $\mathcal{T}', u' \rightrightarrows^{\downarrow\uparrow} \mathcal{T}, u'$ .*

**Lemma 26.**

(1)  $\{\mathcal{T}', u' \mid \mathcal{T}, u \rightrightarrows_{\ell}^{\downarrow} \mathcal{T}', u'\}$  is definable by a node expression  $\chi_{\ell, u, \mathcal{T}}^+$  of  $\text{XPath}_=(\downarrow)^+$  with downward depth  $\leq \ell$ .

(2)  $\{\mathcal{T}', u' \mid \mathcal{T}, u \rightrightarrows_{r,s,k}^{\downarrow\uparrow} \mathcal{T}', u'\}$  is definable by a node expression  $\chi_{r,s,k,u,\mathcal{T}}^+$  of  $\text{XPath}_=(\downarrow\uparrow)^+$  with vertical depth  $\leq (r, s)$  and nesting depth  $\leq k$ .

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

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

Since every  $\Phi_{\mathcal{T}', u'}$  is finite up to logical equivalence by Proposition 16, it follows that  $\Psi$  is a valid node expression. We show that it defines  $\text{sim}_{r,s,k}^{\downarrow\uparrow}(\mathcal{T}, u)$ .

Let  $\mathcal{T}', u' \in \text{sim}_{r,s,k}^{\downarrow\uparrow}(\mathcal{T}, u)$ . 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  $\mathcal{T}'', u'' \in \text{sim}_{r,s,k}^{\downarrow\uparrow}(\mathcal{T}, u)$  and then  $\mathcal{T}', u' \equiv_{r,s,k}^{\downarrow\uparrow} \mathcal{T}'', u''$ . By Theorem 20-3 we then have that  $\mathcal{T}', u' \leftrightarrow_{r,s,k}^{\downarrow\uparrow} \mathcal{T}'', u''$ , and in particular  $\mathcal{T}'', u'' \rightrightarrows_{r,s,k}^{\downarrow\uparrow} \mathcal{T}', u'$ . Since  $\mathcal{T}, u \rightrightarrows_{r,s,k}^{\downarrow\uparrow} \mathcal{T}'', u''$  and  $\mathcal{T}'', u'' \rightrightarrows_{r,s,k}^{\downarrow\uparrow} \mathcal{T}', u'$ , then  $\mathcal{T}, u \rightrightarrows_{r,s,k}^{\downarrow\uparrow} \mathcal{T}', u'$  (by transitivity of  $\rightrightarrows_{r,s,k}^{\downarrow\uparrow}$ ) and thus  $\mathcal{T}', u' \in \text{sim}_{r,s,k}^{\downarrow\uparrow}(\mathcal{T}, u)$ .

Item (1) is shown in a similar way, using Proposition 2 and Theorem 8-2.  $\square$

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

**Theorem 27.**

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

*Proof.* We start with item (1), for the case of  $\rightrightarrows_{\ell}^{\downarrow}$ . The ‘if’ part is straightforward from Theorem 24-2, and here we focus on the ‘only if’ part. Let  $\varphi$  be preserved under  $\rightrightarrows_{\ell}^{\downarrow}$ . Let  $\{(\mathcal{T}_i, u_i)\}_{i \leq n}$  be the set of all pointed models of  $\varphi$  modulo  $\leftrightarrow_{\ell}^{\downarrow}$  (which is finite due to Theorem 8-2 together with Proposition 2). We claim that

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

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

Let  $\chi_{\ell, u_i, \mathcal{T}_i} \in \text{XPath}_=(\downarrow)^+$ ,  $\text{dd}(\psi_i) \leq \ell$ , be as in Lemma 26-(1). Using (3) one shows that  $\bigvee_{i \leq n} \chi_{\ell, u_i, \mathcal{T}_i} \equiv \varphi$ .

For the case of  $\rightrightarrows^{\downarrow}$  of item (1), the ‘if’ direction follows from Theorem 24-1. For the ‘only if’ direction, let  $\varphi$  be preserved under  $\rightrightarrows^{\downarrow}$ . It is easy to see that  $\varphi$  is preserved under  $\rightrightarrows^{\downarrow}$  iff it is preserved under  $\rightrightarrows_{\text{dd}(\varphi)}^{\downarrow}$ . 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 23 and Lemma 26-(2).

Item (4) is straightforward from Theorem 24-3.

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

### 3.4 Transitive Axes

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 also 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  $\uparrow$ . Hence the notions of bisimulations for  $\text{XPath}_=(\downarrow\downarrow_*)$  and  $\text{XPath}_=(\downarrow\uparrow\downarrow_*\uparrow^*)$  (denoted  $\Leftrightarrow^{\downarrow\downarrow_*}$  and  $\Leftrightarrow^{\downarrow\uparrow\downarrow_*\uparrow^*}$  respectively) coincide with  $\Leftrightarrow^{\downarrow}$  and  $\Leftrightarrow^{\downarrow\uparrow}$ , respectively.

Let  $\equiv^{\downarrow\downarrow_*}$  and  $\equiv^{\downarrow\uparrow\downarrow_*\uparrow^*}$  be the logical equivalence relation for the fragments  $\text{XPath}_=(\downarrow\downarrow_*)$  and  $\text{XPath}_=(\downarrow\uparrow\downarrow_*\uparrow^*)$  respectively, and let  $\Rightarrow^{\downarrow\downarrow_*}$  and  $\Rightarrow^{\downarrow\uparrow\downarrow_*\uparrow^*}$  be the logical implication for  $\text{XPath}_=(\downarrow\downarrow_*)^+$  and  $\text{XPath}_=(\downarrow\uparrow\downarrow_*\uparrow^*)^+$  respectively.

**Theorem 28.** *Let  $\dagger \in \{\downarrow\downarrow_*, \downarrow\uparrow\downarrow_*\uparrow^*\}$ .*

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

*Proof.* The proof that  $\mathcal{T}, u \Leftrightarrow^\downarrow \mathcal{T}', u' \Rightarrow \mathcal{T}, u \equiv^{\downarrow\downarrow_*} \mathcal{T}', u'$  follows from a simple adaptation of Proposition 9 to the logic  $\text{XPath}_=(\downarrow\downarrow_*)$  and Lemma 10. The fact that for finitely branching,  $\mathcal{T}, u \equiv^{\downarrow\downarrow_*} \mathcal{T}', u' \Rightarrow \mathcal{T}, u \Leftrightarrow^\downarrow \mathcal{T}', u'$  is straightforward from Theorem 8-1 since  $\equiv^{\downarrow\downarrow_*} \subseteq \equiv^\downarrow$ . The cases for  $\text{XPath}_=(\downarrow\uparrow\downarrow_*\uparrow^*)$ ,  $\text{XPath}_=(\downarrow\downarrow_*)$  and  $\text{XPath}_=(\downarrow\uparrow\downarrow_*\uparrow^*)^+$  are analogous.  $\square$

It is not hard to see that the adequate notion of (bi)simulation of any other ‘intermediate’ fragment between  $\text{XPath}_=(\uparrow\downarrow)$  and  $\text{XPath}_=(\downarrow\uparrow\downarrow_*\uparrow^*)$  —such as  $\text{XPath}_=(\downarrow\uparrow\downarrow_*)$  and  $\text{XPath}_=(\downarrow\uparrow\uparrow^*)$ — also corresponds to that of  $\text{XPath}_=(\uparrow\downarrow)$  in the sense of the statement above.

On the other hand, if we restrict formulas to have only transitive axes, we obtain that the notion of bisimulation for  $\text{XPath}_=(\downarrow_*)$  has a coarser bisimulation notion, as we define next.

Let  $\mathcal{T}$  and  $\mathcal{T}'$  be two data-trees. We say that  $u \in T$  and  $u' \in T'$  are **bisimilar for  $\text{XPath}_=(\downarrow_*)$**  (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:** If  $xZx'$  and there are  $x \xrightarrow{*} v_1 \xrightarrow{*} \dots \xrightarrow{*} v_n$ , and  $x \xrightarrow{*} w_1 \xrightarrow{*} \dots \xrightarrow{*} w_m$  in  $\mathcal{T}$  then there are  $x' \xrightarrow{*} v'_1 \xrightarrow{*} \dots \xrightarrow{*} v'_n$ , and  $x' \xrightarrow{*} w'_1 \xrightarrow{*} \dots \xrightarrow{*} w'_m$  in  $\mathcal{T}'$  such that
  1.  $\text{data}(v_n) = \text{data}(w_m) \Leftrightarrow \text{data}(v'_n) = \text{data}(w'_m)$ ,
  2.  $v_i Z v'_i$  for all  $1 \leq i \leq n$ , and
  3.  $w_i Z w'_i$  for all  $1 \leq i \leq m$ .
- **Zag:** If  $xZx'$  and there are  $x' \xrightarrow{*} v'_1 \xrightarrow{*} \dots \xrightarrow{*} v'_n$ , and  $x' \xrightarrow{*} w'_1 \xrightarrow{*} \dots \xrightarrow{*} w'_m$  in  $\mathcal{T}'$  then there are  $x \xrightarrow{*} v_1 \xrightarrow{*} \dots \xrightarrow{*} v_n$ , and  $x \xrightarrow{*} w_1 \xrightarrow{*} \dots \xrightarrow{*} w_m$  in  $\mathcal{T}$  such that items 1, 2 and 3 above are verified.

As before, one can define  $\ell$ -bisimilarity for  $\text{XPath}_=(\downarrow_*)$ , notated  $\mathcal{T}, u \Leftrightarrow_{\ell}^{\downarrow_*} \mathcal{T}', u'$ , and also the notions of equivalence  $\equiv^{\downarrow_*}$  and  $\equiv_{\ell}^{\downarrow_*}$  as expected.

These notions of bisimulation coincide with the corresponding logical equivalences:

**Theorem 29.**

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 finite.
2.  $\mathcal{T}, u \Leftrightarrow_{\ell}^{\downarrow_*} \mathcal{T}', u'$  iff  $\mathcal{T}, u \equiv_{\ell}^{\downarrow_*} \mathcal{T}', u'$ .

For the case of  $\text{XPath}_=(\downarrow_*\uparrow^*)$  it is not obvious how to adapt the bisimulation of vertical XPath, since the normal form results does not hold for  $\text{XPath}_=(\downarrow_*\uparrow^*)$ .

## 4. Characterization

In Section 4.1 we characterize  $\text{XPath}_=(\downarrow)$  as the fragment of first-order logic  $\Leftrightarrow^{\downarrow}$ -invariant over data trees. In Section 4.2 we show that this result fails for  $\text{XPath}_=(\uparrow\downarrow)$  in general, but a weaker result holds: Any first-order formula  $\Leftrightarrow_{r,s,k}^{\uparrow\downarrow}$ -invariant for some  $r, s, k$  is equivalent to a  $\text{XPath}_=(\uparrow\downarrow)$  formula.

### 4.1 Downward XPath

Recall that a data tree  $\mathcal{T}$  and  $u \in T$ , let  $\mathcal{T}|_{\ell}u$  denote the subtree of  $\mathcal{T}$  induced by  $\{v \in T \mid (\exists n \leq \ell) u \xrightarrow{n} v\}$ . Any data tree can be regarded as a  $\sigma$ -structure, as explained in §2.1. An  $\text{FO}(\sigma)$ -formula  $\varphi(x)$  is  $\ell$ -**local** if for all data trees  $\mathcal{T}$  and  $u \in T$ , we have  $\mathcal{T} \models \varphi(u) \Leftrightarrow \mathcal{T}|_{\ell}u \models \varphi(u)$ . Recall that  $\varphi \in \text{FO}(\sigma)$ ,  $\text{qr}(\varphi)$  is the quantifier rank of  $\varphi$ .

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

**Theorem 30** (Characterization). *Let  $\varphi(x) \in \text{FO}(\sigma)$ . The following are equivalent:*

1.  $\varphi$  is  $\Leftrightarrow^{\downarrow}$ -invariant over [finite] data-trees;
2.  $\varphi$  is logically equivalent over [finite] data-trees to a node expression of  $\ell$ - $\text{XPath}_=(\downarrow)$ , where  $\ell = 2^{\text{qr}(\varphi)} - 1$ .

The proof of this theorem, whose proof is afterwards, will be a consequence of the following three propositions:

**Proposition 31.** *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 the proof by Otto (2004a). Assume that  $\varphi(x) \in \text{FO}(\sigma)$  is  $\Leftrightarrow^{\downarrow}$ -invariant, let  $q = \text{qr}(\varphi)$ , and put  $\ell = 2^q - 1$ . Given a data tree  $\mathcal{T}$  and  $u \in T$  it suffices to show the existence of data trees  $\mathcal{T}'$  and  $\mathcal{T}''$ , with corresponding elements  $u' \in T'$  and  $u'' \in T''$  such that

- (a)  $\mathcal{T}', u' \Leftrightarrow^{\downarrow} \mathcal{T}, u$ ,

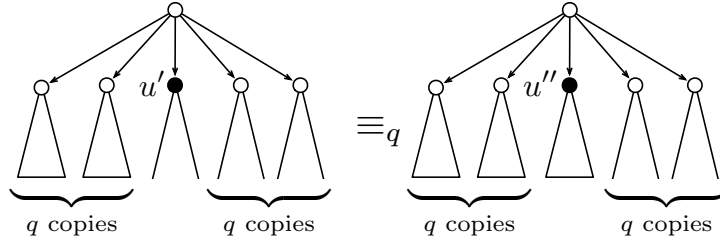
(b)  $\mathcal{T}'', u'' \leftrightarrow^\downarrow (\mathcal{T}|_{\ell u}), u$ , and

(c)  $\mathcal{T}', u' \equiv_q \mathcal{T}'', u''$ .

Indeed, from the above conditions it follows that

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

and hence  $\varphi$  is  $\ell$ -local. By Observation 4 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  $\mathcal{T}|_{\ell u}$ . Consider the structures  $\mathcal{T}'$  and  $\mathcal{T}''$  in the diagram below,



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 Observation 5, 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 as the one used in the paper by Otto (2004a).  $\square$

**Proposition 32.** *Any  $\leftrightarrow^\downarrow$ -invariant  $\varphi(x) \in \text{FO}(\sigma)$  over [finite] data-trees that is  $\ell$ -local, is  $\leftrightarrow_\ell^\downarrow$ -invariant.*

*Proof.* Let  $\varphi(x)$  be  $\ell$ -local and  $\leftrightarrow^\downarrow$ -invariant. Suppose  $\mathcal{T}, u \leftrightarrow_\ell^\downarrow \mathcal{T}', u'$  and  $\mathcal{T} \models \varphi(u)$ . By  $\ell$ -locality,  $\mathcal{T}|_{\ell u} \models \varphi(u)$ . Now

$$\mathcal{T}, u \leftrightarrow_\ell^\downarrow \mathcal{T}', u' \text{ iff } (\mathcal{T}|_{\ell u}), u \leftrightarrow_\ell^\downarrow (\mathcal{T}'|_{\ell u'}), u' \quad \text{(Obs. 7)}$$

$$\text{iff } (\mathcal{T}|_{\ell u}), u \leftrightarrow^\downarrow (\mathcal{T}'|_{\ell u'}), u'. \quad \text{(Obs. 6)}$$

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

**Proposition 33.** *If  $\varphi(x) \in \text{FO}(\sigma)$  is  $\leftrightarrow_\ell^\downarrow$ -invariant over [finite] data-trees, then there is  $\psi \in \ell\text{-XPath}_=(\downarrow)$  such that  $\text{Tr}_x(\psi)$  is logically equivalent to  $\varphi$  over [finite] data-trees.*



*Proof.* By Corollary 3, for every data tree  $\mathcal{T}$  and  $u \in T$  there is a node expression  $\chi_{\ell, \mathcal{T}, u} \in \ell\text{-XPath}_{=}(↓)$  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  $\chi_{\ell, \mathcal{T}, u} \in \ell\text{-XPath}_{=}(↓)$  and, by Proposition 2,  $\equiv_{\ell}^{\downarrow}$  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  $\leftrightarrow_{\ell}^{\downarrow}$ -invariant (and hence  $\equiv_{\ell}^{\downarrow}$ -invariant by Theorem 8-2) we conclude  $\mathcal{T} \models \varphi(u)$ .  $\square$

*Proof of theorem 30.* The implication  $2 \Rightarrow 1$  follows straightforwardly from Theorem 8. The proof of  $1 \Rightarrow 2$  is as follows: First, we show that any  $\leftrightarrow_{\ell}^{\downarrow}$ -invariant  $\varphi(x) \in \text{FO}(\sigma)$  is  $\ell$ -local for  $\ell = 2^{\text{ar}(\varphi)} - 1$  (Proposition 31). Then, we prove that any  $\leftrightarrow_{\ell}^{\downarrow}$ -invariant  $\varphi(x) \in \text{FO}(\sigma)$  that is  $\ell$ -local is  $\leftrightarrow_{\ell}^{\downarrow}$ -invariant (see Proposition 32 below). Finally, we show that any  $\text{FO}(\sigma)$ -definable property which is  $\leftrightarrow_{\ell}^{\downarrow}$ -invariant is definable in  $\ell\text{-XPath}_{=}(↓)$  (see Proposition 33 below).  $\square$

## 4.2 Vertical XPath

The analog of Theorem 30 fails for  $\text{XPath}_{=}(↑↓)$ . This basically because the property  $p(x) = \text{“the tree where } x \text{ belongs contains a label } a\text{”}$  is  $\leftrightarrow_{\ell}^{\uparrow\downarrow}$ -invariant while it is not expressible in  $\text{XPath}_{=}(↑↓)$ . Notice that, however,  $p$  is not  $\leftrightarrow_{\ell}^{\downarrow}$ -invariant. This is as expected, since otherwise  $p$  would be expressed by a  $\text{XPath}_{=}(↓)$  formula. Indeed, consider two data trees with only two nodes, the root and a leaf, labeled respectively  $a, b$  in one tree, and  $b, b$  in the other tree. Note that the leaves of these trees are  $\leftrightarrow_{\ell}^{\downarrow}$ -bisimilar although they don't coincide in the property  $p$ .

**Lemma 34.** *The  $\text{FO}(\sigma)$ -formula  $(\exists x) P_a(x)$  is  $\leftrightarrow_{\ell}^{\uparrow\downarrow}$ -invariant though not logically equivalent over [finite] data-trees to any node expression of  $\text{XPath}_{=}(↑↓)$ .*

*Proof.* Let  $\varphi(x)$  be the  $\text{FO}(\sigma)$ -formula for *there is a node labeled } a \text{ in the tree}*, i.e.,  $\varphi(x) = (\exists y) P_a(y)$ . We prove that  $\varphi$  is  $\leftrightarrow_{\ell}^{\uparrow\downarrow}$ -invariant over [finite] data-trees, though it is not logically equivalent over [finite] data-trees to any node expression of  $\text{XPath}_{=}(↑↓)$ .

To see that  $\varphi$  is  $\leftrightarrow_{\ell}^{\uparrow\downarrow}$ -invariant over [finite] data-trees, take  $\mathcal{T}, u$  and  $\mathcal{T}', u'$  such that  $\mathcal{T}, u \leftrightarrow_{\ell}^{\uparrow\downarrow} \mathcal{T}', u'$  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 20,  $\mathcal{T}', u' \models \uparrow^n \downarrow^m a$  and so  $\mathcal{T}' \models \varphi(u')$ .

Assume by contradiction that there is  $\psi \in \text{XPath}_{=}(↑↓)$  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 Observation 19 we have  $\mathcal{T}, u \leftrightarrow_{r, s, k}^{\uparrow\downarrow} (\mathcal{T}|_r^s u), u$ . Since  $\mathcal{T}, u \models \psi$ , by Theorem 20, 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  $a$ .  $\square$

Hence  $\text{XPath}_=(\uparrow\downarrow)$  is not the fragment of  $\text{FO}(\sigma)$  which is  $\Leftrightarrow^{\downarrow\uparrow}$ -invariant over [finite] data-trees. However, the following analog of Proposition 33 (needed for the proof of Theorem 30) still holds for the case of  $\text{XPath}_=(\uparrow\downarrow)$ : For any  $r, s, k$ , every first-order formula  $\Leftrightarrow_{r,s,k}^{\downarrow\uparrow}$ -invariant is equivalent to a  $\text{XPath}_=(\uparrow\downarrow)$  formula.

**Proposition 35.** *Let  $k' = k \cdot (r + s + 2)$ . If  $\varphi(x) \in \text{FO}(\sigma)$  is  $\Leftrightarrow_{r,s,k'}^{\downarrow\uparrow}$ -invariant over [finite] data-trees, then there is  $\psi \in (r, s, k)\text{-XPath}_=(\uparrow\downarrow)$  such that  $\text{Tr}_x(\psi)$  is logically equivalent to  $\varphi$  over [finite] data-trees.*

*Proof.* By Corollary 17, for every data tree  $\mathcal{T}$  and  $u \in T$  there is a node expression  $\chi_{r,s,k,\mathcal{T},u}$  of  $(r, s, k)\text{-XPath}_=(\uparrow\downarrow)$  such that  $\mathcal{T}, u \equiv_{r,s,k}^{\downarrow\uparrow} \mathcal{T}', u'$  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  $\chi_{r,s,k,\mathcal{T},u} \in (r, s, k)\text{-XPath}_=(\uparrow\downarrow)$  and, by Proposition 16,  $\equiv_{r,s,k}^{\downarrow\uparrow}$  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 33, as we show next. Let us see that  $\varphi \models \text{Tr}_x(\psi)$ . Suppose  $\mathcal{T} \models \varphi(u)$ . Since  $\mathcal{T}, u \models \chi_{r,s,k,\mathcal{T},u}$ , 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  $\mathcal{T}, u \equiv_{r,s,k}^{\downarrow\uparrow} \mathcal{T}', u'$  and since  $\varphi$  is  $\Leftrightarrow_{r,s,k \cdot (r+s+2)}^{\downarrow\uparrow}$ -invariant (and hence  $\equiv_{r,s,k}^{\downarrow\uparrow}$ -invariant by Theorem 20-2) we conclude  $\mathcal{T} \models \varphi(u)$ .  $\square$

Notice that the counterexample in Lemma 34 is an unrestricted, existential formula. One may wonder if it might be possible to extend the expressive power of  $\text{XPath}_=(\uparrow\downarrow)$  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 30 fails. Formally, consider the logic  $\text{XPath}_=(\uparrow\downarrow E)$ , which results from adding the operator  $E$  to  $\text{XPath}_=(\uparrow\downarrow)$  with the following semantics:  $\llbracket E\varphi \rrbracket^{\mathcal{T}} = T$  if  $\llbracket \varphi \rrbracket^{\mathcal{T}} \neq \emptyset$ , and  $\llbracket E\varphi \rrbracket^{\mathcal{T}} = \emptyset$  otherwise.

The following lemma shows a counterexample to the analog of Theorem 30, showing that  $\text{XPath}_=(\uparrow\downarrow E)$  is not the fragment of  $\text{FO}(\sigma)$   $\Leftrightarrow^{\downarrow\uparrow}$ -invariant over [finite] data-trees.

**Lemma 36.** *The  $\text{FO}(\sigma)$ -formula  $(\exists y, z) [y \approx z \wedge P_a(y) \wedge P_b(z)]$  is  $\Leftrightarrow^{\downarrow\uparrow}$ -invariant though not logically equivalent over [finite] data-trees to any node expression of  $\text{XPath}_=(\uparrow\downarrow E)$ .*

*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  $\text{XPath}_=(\downarrow, \uparrow, E)$ . Suppose, by means of contradiction, that there is a node expression  $\psi \in \text{XPath}_=(\downarrow, \uparrow, E)$  expressing  $\varphi$ , with  $\text{vd}(\psi) = (r, s)$  ( $\text{vd}(\cdot)$  for  $\text{XPath}_=(\downarrow, \uparrow, E)$  is defined as before together with 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  $\text{XPath}_{=}(\downarrow, \uparrow, E)$ .  $\square$

## 5. Applications

We devote this section to exemplify how the model theoretic tools we developed can be used to show expressiveness results for  $\text{XPath}_{=}$ . 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.

### 5.1 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. Observation 4, 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 3 gives a schematic representation of this operation.

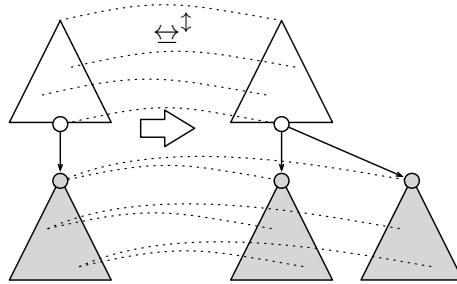


Figure 3: Closure under subtree replication.

**Proposition 37.**  $\text{XPath}_{=}(\downarrow, \uparrow, \downarrow_*, \uparrow^*)$  is closed under subtree replication, i.e. if  $\mathcal{T}'$  is a subtree replication of  $\mathcal{T}$ , and  $u \in T$  then  $\mathcal{T}', u \equiv^{\downarrow, \uparrow, \downarrow_*, \uparrow^*} \mathcal{T}, u$ .

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

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

$Z$  is depicted as dotted lines in Figure 3.  $\square$

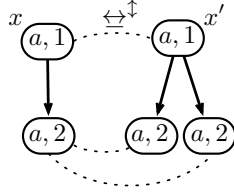
## 5.2 Non-expressivity Results

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  $(\forall x)[P_a(x) \Rightarrow (\exists y)[P_b(y) \wedge x \approx y]]$ .

**Proposition 38.**

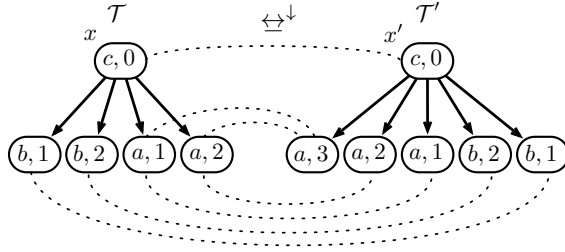
1.  $key(a)$  is not expressible in  $XPath_{=}(\downarrow\uparrow\downarrow_*\uparrow^*)$ .
2.  $fk(a, b)$  is expressible in  $XPath_{=}(\downarrow\uparrow\downarrow_*\uparrow^*)$  but it is not expressible in  $XPath_{=}(\downarrow\downarrow_*)$  or  $XPath_{=}(\downarrow\uparrow\downarrow_*\uparrow^*)^+$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

*Proof.* The first item follows from Proposition 37. Since the logic is closed under subtree replication, the trees below 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  $\neg\langle\uparrow^*\downarrow_*[a \wedge \neg\langle\varepsilon = \uparrow^*\downarrow_*[b]\rangle]\rangle$ . However, this property cannot be expressed in  $XPath_{=}(\downarrow\downarrow_*)$  because the models  $\mathcal{T}$  and  $\mathcal{T}'$  below 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 28 it follows that  $fk(a, b)$  is not expressible in  $XPath_{=}(\downarrow\downarrow_*)$ .

Finally, suppose there exists  $\psi \in XPath_{=}(\downarrow\uparrow\downarrow_*\uparrow^*)^+$  expressing  $fk(a, b)$ . Since  $\mathcal{T}$  is a substructure of  $\mathcal{T}'$  we have  $\mathcal{T}, x \xrightarrow{\downarrow\uparrow} \mathcal{T}', x$  by Lemma 25. By Theorem 28(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.

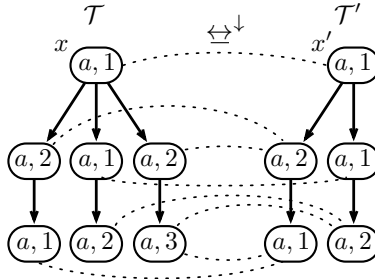
**Proposition 39.**

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

*Proof.* For 1, one can check that  $\mathcal{T}, x \models \varphi$  iff  $\mathcal{T}, x$  satisfies  $dist_3$ , for

$$\varphi = \langle \varepsilon \neq \downarrow\downarrow[\langle \varepsilon \neq \uparrow[\langle \varepsilon \neq \uparrow \rangle]] \rangle.$$

Let us see 2. Consider the data trees  $\mathcal{T}, x$  and  $\mathcal{T}', x'$  depicted below. 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  $\mathcal{T}, x \xleftrightarrow{\downarrow} \mathcal{T}', x'$  via  $Z \subseteq T \times T'$  defined by

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

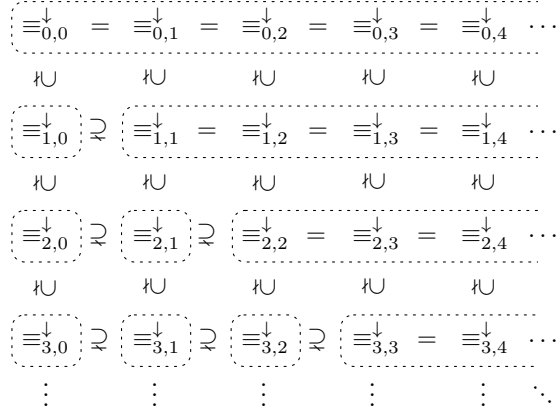
where  $h(y)$  is the height of  $y$ , i.e., the distance from  $y$  to the root of the corresponding tree ( $Z$  is depicted as dotted lines in picture above). Since  $\mathcal{T}, x$  satisfies  $dist_3$  but  $\mathcal{T}', x'$  does not, from Theorem 28 it follows that  $dist_3$  is not expressible in  $XPath_{=}(\downarrow\downarrow_*)$ .

For 3, one can verify that  $\mathcal{T}, x \xrightarrow{\uparrow\downarrow} \mathcal{T}', x'$  via  $Z$  as defined above. If  $dist_3$  were definable in  $XPath_{=}(\downarrow\uparrow\downarrow_*\uparrow^*)^+$  via  $\psi$  and the fact that  $\mathcal{T}, x \models \psi$ , by Theorem 28(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 38-2,  $\overline{dist_3}$  is not expressible in  $XPath_{=}(\downarrow\uparrow\downarrow_*\uparrow^*)^+$ .  $\square$

### 5.3 Expressiveness Hierarchies

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


 Figure 4: Hierarchy of  $\text{XPath}_{=}(↓)$ .

- **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.

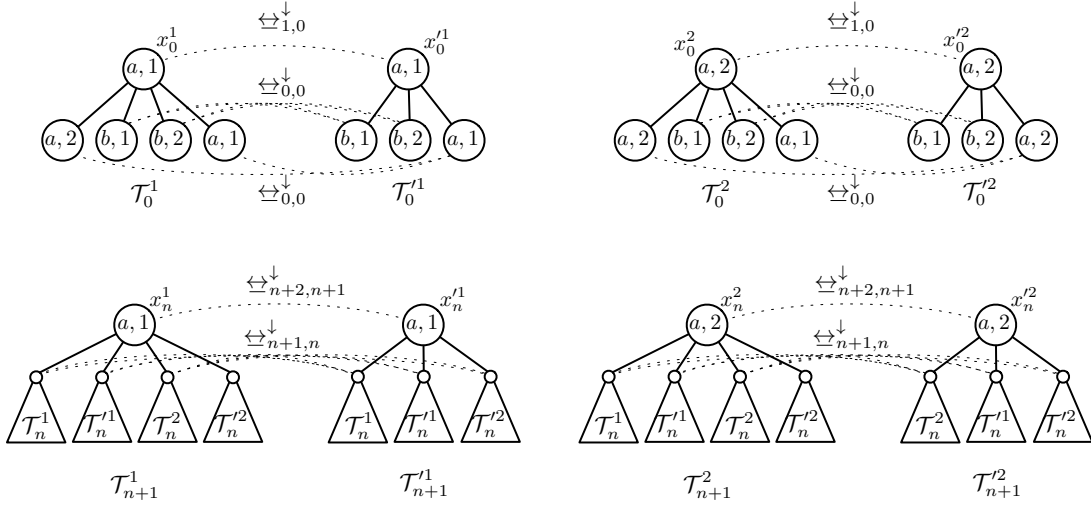
Following the same ideas used in Propositions 9 and 11, it is easy to show that  $(\ell, k)$ -bisimulations characterize  $(\ell, k)$ -equivalence.

**Proposition 40.**  $\mathcal{T}, u \Leftrightarrow_{\ell, k}^{\downarrow} \mathcal{T}', u'$  iff  $\mathcal{T}, u \equiv_{\ell, k}^{\downarrow} \mathcal{T}', u'$ .

The following theorem characterizes when an increase in nesting depth results in an increase in expressive power (see Figure 4). We speculate that a similar hierarchy holds in the absence of data values, but this is not a direct consequence of our result.

**Theorem 41.** For all  $\ell, k \geq 0$ ,  $i \geq 1$ ,  $\equiv_{\ell, 0}^{\downarrow} \supsetneq \equiv_{\ell, 1}^{\downarrow} \supsetneq \cdots \supsetneq \equiv_{\ell, \ell}^{\downarrow} = \equiv_{\ell, \ell+i}^{\downarrow}$ , and  $\equiv_{\ell, k}^{\downarrow} \supsetneq \equiv_{\ell+i, k}^{\downarrow}$ .

*Proof.* Consider the data trees  $\mathcal{T}_n^i$ ,  $\mathcal{T}_n^i$  ( $n \geq 0, i \in \{1, 2\}$ ) defined 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  $\mathcal{T}_k^1, x_k^1 \equiv_{k+1, k}^\downarrow \mathcal{T}_k^{\prime 1}, x_k^{\prime 1}$  but  $\mathcal{T}_k^1, x_k^1 \not\equiv_{k+1, k+1}^\downarrow \mathcal{T}_k^{\prime 1}, x_k^{\prime 1}$ .

That  $\mathcal{T}_k^1, x_k^1 \not\equiv_{k+1, k+1}^\downarrow \mathcal{T}_k^{\prime 1}, x_k^{\prime 1}$  results 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  $k+1$ ,

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

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

To show  $\mathcal{T}_k^1, x_k^1 \equiv_{k+1, k}^\downarrow \mathcal{T}_k^{\prime 1}, x_k^{\prime 1}$  we use Proposition 40 and show  $\mathcal{T}_k^1, x_k^1 \stackrel{\downarrow}{\simeq}_{k+1, k} \mathcal{T}_k^{\prime 1}, x_k^{\prime 1}$ . Note that  $\mathcal{T}_k^1$  and  $\mathcal{T}_k^2$  (resp.  $\mathcal{T}_k^{\prime 1}$  and  $\mathcal{T}_k^{\prime 2}$ ) 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 42.** *Note that the set of immediate subtrees of the roots of  $\mathcal{T}_k^1, \mathcal{T}_k^{\prime 1}, \mathcal{T}_k^2, \mathcal{T}_k^{\prime 2}$  are the same as those of  $\mathcal{T}_k^{\prime 1}, \mathcal{T}_k^2, \mathcal{T}_k^{\prime 2}$  (and of  $\mathcal{T}_k^1, \mathcal{T}_k^{\prime 1}, \mathcal{T}_k^2$ ) by construction.*

We show  $\mathcal{T}_k^1, x_k^1 \stackrel{\downarrow}{\simeq}_{k+1, k} \mathcal{T}_k^{\prime 1}, x_k^{\prime 1}$ . For every  $j \leq k+1, t \leq k$ , let  $Z_{j,t}$  be the set of all pairs  $(x, y) \in T_k^1 \times T_k^{\prime 1}$  so that  $x$  and  $y$  are some  $x_{k'}^i$  or  $x_{k'}^{\prime i}$  for  $i \in \{1, 2\}$  and  $k' \geq t$  (the notation  $x_{k'}^i$  and  $x_{k'}^{\prime i}$  does not necessarily identify a unique node but, possibly, many; the intended meaning is that they can refer to *any* of them). Observe that

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

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  $x_{t'}^1$  and  $u'$  is  $x_{t'}^{\prime 1}$  for some  $t' \leq t$ , the other cases being similar or simpler.

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

- If  $v$  is inside the subtree  $\mathcal{T}'_{t-1}$  of  $\mathcal{T}_t^1$ , but it is not  $x_{t-1}^2$ , then we choose  $v'$  as the corresponding node inside the subtree  $\mathcal{T}'_{t-1}$  of  $\mathcal{T}_t^1$ . Remember that  $\mathcal{T}'_{t-1}$  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.  $data(v) = data(v')$  by Observation 42. Furthermore, since every node of  $\mathcal{T}'_{t-1}$  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  $\mathcal{T}'_{t-1}$ ,  $x_{t-1}^2$ . Again, we have that  $data(v') = data(v)$  and by construction that  $vZ_{j,t-1}v'$ . Thus, by (4),  $vZ_{j-1,t-1}v'$ .
- Finally, if  $v$  falls outside  $\mathcal{T}'_{t-1}$ , we choose  $v'$  as the same node in  $\mathcal{T}_t^1$ , where of course we will have that  $data(v) = data(v')$  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  $\mathcal{T}_k^1, x \leftrightarrow_{k+1,k}^\downarrow \mathcal{T}_k^1, x'$  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  $\ell$ .

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

For the inductive case, let  $Z_{j,t} = \leftrightarrow_{j,t}^\downarrow$  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)  $data(v) = data(w')$  iff  $data(v') = data(w')$ ,
- (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,  $\mathcal{T}, (\xrightarrow{i}v) \leftrightarrow_{\ell-n+i,k}^\downarrow \mathcal{T}', (\xrightarrow{i}v')$  and thus  $(\xrightarrow{i}v)Z_{\ell-n+i,k}(\xrightarrow{i}v')$ . By an identical reasoning,  $\mathcal{T}, (\xrightarrow{i}w) \leftrightarrow_{\ell-n+i,k}^\downarrow \mathcal{T}', (\xrightarrow{i}w')$  and thus  $(\xrightarrow{i}w)Z_{\ell-n+i,k}(\xrightarrow{i}w')$ . Thus, the Zig condition for  $\leftrightarrow_{\ell,k+1}^\downarrow$  is verified. The Zag condition holds by symmetry.  $\square$



With respect to vertical XPath, note that since  $\equiv_{r,s,k}^{\downarrow\uparrow} \subseteq \equiv_{r',s',k'}^{\downarrow\uparrow}$  for all  $(r, s, k) \leq (r', s', k')$ , as a consequence of Proposition 15 we obtain that for every  $r, s, k$  with  $r + s \geq 2$  there is some  $k' > k$  so that  $\equiv_{r,s,k}^{\downarrow\uparrow} \supsetneq \equiv_{r,s,k'}^{\downarrow\uparrow}$ . In fact, we conjecture that  $\equiv_{r,s,k}^{\downarrow\uparrow} \supsetneq \equiv_{r,s,k+1}^{\downarrow\uparrow}$  for every  $k$ . We argue that this can be proven through the models  $(\mathcal{T}_n)_n$  in the proof of Proposition 15, by showing that  $\mathcal{T}_k, x_{r',s'} \equiv_{r,s,k}^{\downarrow\uparrow} \mathcal{T}_{k+1}, x_{r',s'}$  but  $\mathcal{T}_k, x_{r',s'} \not\equiv_{r,s,k+1}^{\downarrow\uparrow} \mathcal{T}_{k+1}, x_{r',s'}$  for every  $(r, s) \geq (r', s')$ . The fact that  $\equiv_{r,s,k}^{\downarrow\uparrow} \supsetneq \equiv_{r+1,s,k}^{\downarrow\uparrow}$  and  $\equiv_{r,s,k}^{\downarrow\uparrow} \supsetneq \equiv_{r,s+1,k}^{\downarrow\uparrow}$  are straightforward. We would then obtain the following.

**Conjecture 43.**  $\equiv_{r,s,k}^{\downarrow\uparrow} \supsetneq \equiv_{r',s',k'}^{\downarrow\uparrow}$  for all  $(r, s, k) < (r', s', k')$ ,  $r + s \geq 2$ .

## 6. 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` ( $\rightarrow$ ) and `previous-sibling` ( $\leftarrow$ ). 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. Next we explain why this is so. Consider the following notation for denoting positions of nodes in a tree. Positions are elements of  $(\mathbb{N} \times \mathbb{N})^*$ , where the root's position is the empty string  $\epsilon$ , and the position of any other node in the tree is the concatenation of the position of its parent and the pair  $(l, r)$ , where  $l$  is the number of siblings to the left of the node and  $r$  the number of siblings to its right. For every position  $p \in (\mathbb{N} \times \mathbb{N})^*$  of a tree, there is a path expression  $\alpha_p$  of XPath<sub>=</sub>( $\downarrow \rightarrow$ ) that can access the node in this position (and only this position) from the root. The  $\alpha_p$ 's are defined as  $\alpha_\epsilon = \varepsilon$  for the root;  $\alpha_{(l,r)} = \downarrow \rightarrow^{l+1} [\langle \rightarrow^k \rangle \wedge \neg \langle \rightarrow^{k+1} \rangle]$  for  $(l, r) \in \mathbb{N} \times \mathbb{N}$ ; and  $\alpha_{p \cdot (l,r)} = \alpha_p \alpha_{(l,r)}$  for  $p \in (\mathbb{N} \times \mathbb{N})^*$ ,  $(l, r) \in \mathbb{N} \times \mathbb{N}$ . It is easy to check that testing  $\bigwedge_p \langle \alpha_p [\neg \langle \downarrow \rangle] \rangle$  where  $p$  ranges over all leaf positions of a given tree  $\mathcal{T}$  tests the property that the subtree hanging from the node where the formula is evaluated is structurally isomorphic to  $\mathcal{T}$ . Further, by also adding the tests  $\langle \alpha_p [a] = \alpha_{p'} [a'] \rangle$  for every pair of positions  $p, p'$  with label  $a, a'$  of  $\mathcal{T}$  having the same data value, as well as  $\langle \alpha_p [a] \neq \alpha_{p'} [a'] \rangle$  for those having different data value, yields the property of being *equal* to  $\mathcal{T}$ , up to isomorphism of data values.

In Section 5 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 maximum auto-bisimulation (observe that the identity is always an auto-bisimulation, and that the union of two bisimulations is a bisimulation; therefore there is a ‘maximum’, often called ‘coarsest’, bisimulation).

The idea is to find the coarsest auto-bisimulation  $Z$  over a given data-tree  $\mathcal{T}$ . One cannot simply make the quotient of  $\mathcal{T}$  over  $Z$ , as the result is not necessarily a tree and it is not clear how to assign data values to each class in the quotient. However, one can make a quotient of  $\mathcal{T}$  over the equivalence relation ‘ $Z$  and same data value’. If we do so, we obtain a smaller structure, and then evaluate all queries here.

Determining the maximum auto-(bi)simulation, either downward or vertical, of a finite data  $\mathcal{T}$  tree can be done in polynomial time. A naive algorithm starts by defining  $Z$  as the set of all nodes which satisfy **Harmony**. Each time **Zig** or **Zag** is not satisfied, it removes from  $Z$  the pair responsible for **Zig** or **Zag** not being true. It repeats this until a fixed point of  $Z$  is found; this is the maximum auto-bisimulation in  $\mathcal{T}$ . If one is interested in deciding if two nodes of different finite data trees are bisimilar, one can use the same idea: the answer is ‘yes’ if and only if the fixed point is not the empty set. Since checking the validity of **Zig** or **Zag** is polynomial time computable (because there are linearly many paths in a tree), and at each step the algorithm decreases the size of  $Z$ , the whole process has polynomial time. Better implementations, based on more sophisticated ideas, such as the works of Henzinger, Henzinger, and Kopke (1995) or Dovier, Piazza, and Policriti (2004), can lead to more efficient polynomial time algorithms. We plan to design and implement algorithms for data tree minimization using bisimulation and investigate their computational complexity.

## Acknowledgements

This work was partially supported by grant ANPCyT-PICT-2013-2011, 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”.

## References

- Abiteboul, S., Bourhis, P., Muscholl, A., & Wu, Z. (2013). Recursive queries on trees and data trees. In *International Conference on Database Theory*, pp. 93–104.
- Abriola, S., Descotte, M. E., & Figueira, S. (2014). Definability for downward and vertical XPath on data trees. In *Workshop on Logic, Language, Information and Computation*, Vol. 6642 of *Lecture Notes in Computer Science*, pp. 20–34.
- Abriola, S., Descotte, M. E., & Figueira, S. (2015). Model theory of XPath on data trees. Part II: Binary bisimulation and definability. To appear in *Information and Computation*. <http://www.glyc.dc.uba.ar/santiago/papers/xpath-part2.pdf>.

- Benedikt, M., Fan, W., & Geerts, F. (2008). XPath satisfiability in the presence of DTDs. *Journal of the ACM*, 55(2), 1–79.
- Benedikt, M., & Koch, C. (2008). XPath leashed. *ACM Computing Surveys*, 41(1).
- Blackburn, P., de Rijke, M., & Venema, Y. (2001). *Modal Logic*, Vol. 53 of *Cambridge Tracts in Theoretical Computer Science*. Cambridge University Press.
- Bojańczyk, M., Muscholl, A., Schwentick, T., & Segoufin, L. (2009). Two-variable logic on data trees and XML reasoning. *Journal of the ACM*, 56(3), 1–48.
- Bojańczyk, M., & Parys, P. (2011). XPath evaluation in linear time. *Journal of the ACM*, 58(4), 17.
- Clark, J., & DeRose, S. (1999). XML path language (XPath). Website. W3C Recommendation. <http://www.w3.org/TR/xpath>.
- David, C. (2008). Complexity of data tree patterns over XML documents. In *Mathematical Foundations of Computer Science*, Vol. 5162 of *Lecture Notes in Computer Science*, pp. 278–289. Springer.
- Dawar, A., & Otto, M. (2009). Modal characterisation theorems over special classes of frames. *Annals of Pure and Applied Logic*, 161(1), 1–42.
- Dovier, A., Piazza, C., & Policriti, A. (2004). An efficient algorithm for computing bisimulation equivalence. *Theor. Comput. Sci*, 311, 221–256.
- Figueira, D. (2010). *Reasoning on Words and Trees with Data*. PhD thesis, Laboratoire Spécification et Vérification, ENS Cachan, France.
- Figueira, D. (2012). Decidability of downward XPath. *ACM Transactions on Computational Logic*, 13(4).
- Figueira, D., & Segoufin, L. (2011). Bottom-up automata on data trees and vertical XPath. In *International Symposium on Theoretical Aspects of Computer Science*, Vol. 9 of *LIPICs*, pp. 93–104. Leibniz-Zentrum für Informatik.
- Figueira, D., Figueira, S., & Areces, C. (2014). Basic model theory of XPath on data trees. In *International Conference on Database Theory*, pp. 50–60.
- Figueira, D., & Libkin, L. (2014). Pattern logics and auxiliary relations. In *Logic in Computer Science*, pp. 40:1–40:10.
- Figueira, S., & Gorín, D. (2010). On the size of shortest modal descriptions.. In *Advances in Modal Logic*, Vol. 8, pp. 114–132.
- Forti, M., & Honsell, F. (1983). Set theory with free construction principles. *Annali Scuola Normale Superiore, Pisa*, X(3), 493–522.
- Goranko, V., & Otto, M. (2007). Model theory of modal logic. In P. Blackburn, J. V. B., & Wolter, F. (Eds.), *Handbook of Modal Logic*, Vol. 3 of *Studies in Logic and Practical Reasoning*, chap. 5, pp. 249–329. Elsevier.
- Gottlob, G., Koch, C., & Pichler, R. (2005). Efficient algorithms for processing XPath queries. *ACM Transactions on Database Systems*, 30(2), 444–491.

- Gyssens, M., Paredaens, J., Gucht, D. V., & Fletcher, G. (2006). Structural characterizations of the semantics of XPath as navigation tool on a document. In *Principles of Database Systems*, pp. 318–327. ACM.
- Harel, D. (1984). Dynamic logic. In Gabbay, D., & Guenther, F. (Eds.), *Handbook of Philosophical Logic. Vol. II*, Vol. 165 of *Synthese Library*, pp. 497–604. D. Reidel Publishing Co., Dordrecht. Extensions of classical logic.
- Henzinger, M. R., Henzinger, T. A., & Kopke, P. W. (1995). Computing simulations on finite and infinite graphs. In *Proc. of 36th Annual Symposium on Foundations of Computer Science*, pp. 453–462. IEEE Computer Society Press.
- Hopcroft, J. (1971). *An  $n \log(n)$  algorithm for minimizing states in a finite automaton*. In Z. Kohave, editor, *Theory of Machines and Computations*. Academic Press.
- Jurziński, M., & Lazić, R. (2011). Alternating automata on data trees and xpath satisfiability. *ACM Transactions on Computational Logic*, 12(3), 19.
- Kanellakis, P., & Smolka, S. (1990). CCS expressions, finite state processes, and three problems of equivalence. *Inf. Comput.*, 86(1), 43–68.
- Kurtonina, N., & de Rijke, M. (1997). Simulating without negation. *Journal of Logic and Computation*, 7, 503–524.
- Lutz, C., Piro, R., & Wolter, F. (2011). Description logic tboxes: Model-theoretic characterizations and rewritability. In *International Joint Conference on Artificial Intelligence, Barcelona, Catalonia, Spain, July 16-22, 2011*, pp. 983–988.
- Marx, M. (2004). XPath with conditional axis relations. In *International Conference on Extending Database Technology*, Vol. 2992 of *Lecture Notes in Computer Science*, pp. 477–494. Springer.
- Marx, M., & de Rijke, M. (2005). Semantic characterizations of navigational XPath. *SIGMOD Record*, 34(2), 41–46.
- Milner, R. (1980). *A Calculus of Communicating Systems*, Vol. 92 of *Lecture Notes in Computer Science*. Springer.
- Neven, F., Schwentick, T., & Vianu, V. (2004). Finite state machines for strings over infinite alphabets. *ACM Transactions on Computational Logic*, 5(3), 403–435.
- Olteanu, D. (2007). Forward node-selecting queries over trees. *ACM Transactions on Database Systems*, 32(1), 3.
- Olteanu, D., Meuss, H., Furche, T., & Bry, F. (2002). XPath: Looking forward. In *International Conference on Extending Database Technology*, pp. 109–127.
- Otto, M. (2004a). Elementary proof of the van Benthem-Rosen characterisation theorem. Tech. rep. 2342, Fachbereich Mathematik, Technische Universität Darmstadt.
- Otto, M. (2004b). Modal and guarded characterisation theorems over finite transition systems. *Annals of Pure and Applied Logic*, 130(1-3), 173–205.
- Otto, M. (2006). Bisimulation invariance and finite models. In *Logic Colloquium'02*, Vol. 27 of *Lecture Notes in Logic*, pp. 276–298.

- Paige, R., & Tarjan, R. (1987). Three partition refinement algorithms. *SIAM J. Comput.*, *16*(6), 973–989.
- Park, D. (1981). Concurrency and automata on infinite sequences. In *Theoretical Computer Science*, Vol. 104 of *Lecture Notes in Computer Science*, pp. 167–183. Springer.
- Rosen, E. (1997). Modal logic over finite structures. *Journal of Logic, Language and Information*, *6*(4), 427–439.
- Sangiorgi, D. (2009). On the origins of bisimulation and coinduction. *ACM Transactions on Programming Languages and Systems*, *31*(4).
- ten Cate, B., Fontaine, G., & Litak, T. (2010). Some modal aspects of XPath. *Journal of Applied Non-Classical Logics*, *20*(3), 139–171.
- van Benthem, J. (1976). *Modal Correspondence Theory*. PhD thesis, Universiteit van Amsterdam.