



**HAL**  
open science

## On XPath with transitive axes and data tests

Diego Figueira

► **To cite this version:**

Diego Figueira. On XPath with transitive axes and data tests. ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS), Jun 2013, New York, NY, United States. 10.1145/2463664.2463675 . hal-01795148

**HAL Id: hal-01795148**

**<https://hal.science/hal-01795148v1>**

Submitted on 18 May 2018

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

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

# On XPath with Transitive Axes and Data Tests<sup>\*</sup>

Diego Figueira  
University of Edinburgh  
Edinburgh, UK

## ABSTRACT

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

## Categories and Subject Descriptors

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

## General Terms

Algorithms, Languages

## Keywords

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

---

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

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

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

## 1. INTRODUCTION

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

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

## Formalisms for trees with data values

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

### *First-order logic.*

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

### *Automata.*

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

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

### *XPath.*

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

## Contribution

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

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

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

Due to space limitations, all the proofs are included in the Appendix, although the main intuitions, general strategy, and decomposition of the proof in simpler problems are presented in the body of the paper.

## 2. PRELIMINARIES

### Basic notation.

Let  $\mathbb{N}_0 \stackrel{\text{def}}{=} \{0, 1, 2, \dots\}$ ,  $\mathbb{N} \stackrel{\text{def}}{=} \{1, 2, 3, \dots\}$ , and let  $[n] \stackrel{\text{def}}{=} \{1, \dots, n\}$  for any  $n \in \mathbb{N}$ . We fix once and for all  $\mathbb{D}$  to be any infinite domain of data values; for simplicity in our examples we will consider  $\mathbb{D} = \mathbb{N}_0$ . In general we use the symbols  $\mathbb{A}$ ,  $\mathbb{B}$  for finite alphabets, and the symbols  $\mathbb{E}$  and  $\mathbb{F}$  for any kind of alphabet. By  $\mathbb{E}^*$  we denote the set of finite sequences over  $\mathbb{E}$ , by  $\mathbb{E}^+$  the set of finite sequences with at least one element over  $\mathbb{E}$ . We write ‘ $\epsilon$ ’ for the empty sequence and

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

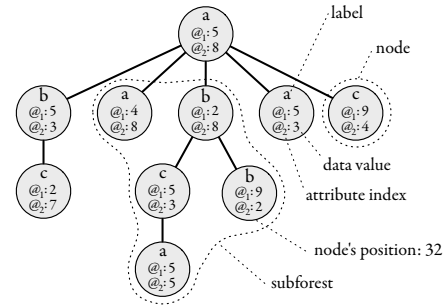


Figure 1: A multi-attribute data tree, where  $\mathbb{A} = \{a, b, c\}$  and  $k = 2$ .

‘ $\cdot$ ’ as the concatenation operator between sequences. By  $|S|$  we denote the length of  $S$  (if  $S$  is a finite sequence), or its cardinality (if  $S$  is a set). We use  $(a_i)_{i \in \{j, \dots, j+n\}}$  as short for  $a_j a_{j+1} \dots a_{j+n}$ .

### Unranked finite trees with data.

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

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

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

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

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

### XPath.

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

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

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

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

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

## 3. PROOF SKETCH

The main contribution of this paper is the following.

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

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

We first introduce *forest profiles* in Section 4, which constitute the algebra domain. A forest profile is an abstraction of a multi-attribute data forest inside a context, where the context consists of the two (possibly empty) forests that are

to the left and to the right. Figure 3 depicts one such possible forest, together with the left and right context forests. A forest profile contains, for each data value  $d$  and each path expression, the information of whether  $d$  can be reached by the path expression, and *where* it can be reached, either

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

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

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

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

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

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

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

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

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

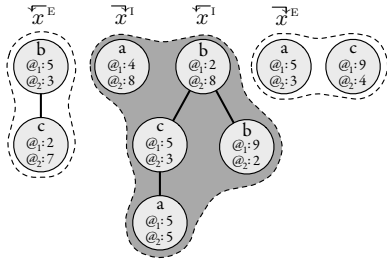


Figure 3: A multi-attribute data forest, with its left and right forests.

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

In Section 6 we show that the satisfiability problem for XPath( $* \leftarrow, \downarrow_*, \rightarrow^*, =$ ) can be reduced to the derivation problem in EXPSPACE. In Section 6.1 we show a normal form, called *direct unnested normal form*, where direct unnested path expressions correspond, precisely, to the pattern expressions used in the forest profiles (basically all path expressions are of the form already described). We then show in Section 6.2 that one can reduce, in EXPSPACE, the satisfiability problem for formulas in this normal form into the derivation problem, obtaining a 3EXPSPACE decidability procedure for SAT-XPath( $* \leftarrow, \downarrow_*, \rightarrow^*, =$ ), obtaining Theorem 3.1.

## 4. FOREST PROFILES

We define abstractions of forests of multi-attribute data trees. These are called *forest profiles*. They are the main construct in our solution. One must think of a forest profile

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

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

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

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

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

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

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

The intended meaning is that the first component operates on the siblings, the second on a downward path, and the third retrieves a data value from an attribute index. We will sometimes use the symbol  $\bar{\alpha}$  to represent elements from  $\Pi$ , or  $(\alpha, \beta, i)$  if we need to explicit the components of the composed pattern.

A **forest profile**  $f$  is a tuple

$$f = (\overleftarrow{\chi}^e, \overrightarrow{\chi}^i, \overleftarrow{\chi}^i, \overrightarrow{\chi}^e, R)$$

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

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

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

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

## 4.1 Rigid and flexible values

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

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

We formalize this by defining an ordering on profiles corresponding to the operation just described, so that the forest profile abstracting  $\bar{t}'$  is bigger than the profile abstracting  $\bar{t}$ . We make explicit (in Lemma 5.1) the aforementioned argument as a monotonicity property of the algebra.

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

## 4.2 Algebra

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

### Preliminaries

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

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

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

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

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

### Fingerprints

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

Let  $\mathcal{A} = \{\overleftarrow{\circ}, \overrightarrow{\circ}, \circ^\downarrow, \circ\}$ . Given a profile  $f \in \mathfrak{F}$  and  $a \in \mathcal{A}$ , we define the set  $f \cdot \chi_a$  as

- $\overleftarrow{\chi}^i \cup \llbracket f \rrbracket \cdot \overleftarrow{\chi}^e$  if  $a = \overleftarrow{\circ}$ ,
- $\overleftarrow{\chi}^i \cup \{f\} \cdot \overleftarrow{\chi}^e$  if  $a = \overrightarrow{\circ}$ ,
- $\{(d, \alpha', \beta, i) \in \mathbb{D} \times \Pi \mid \exists \alpha. (d, \alpha, \beta, i) \in \overleftarrow{\chi}^i \cup \overrightarrow{\chi}^i\}$  if  $a = \circ^\downarrow$ , or
- $\{(d, \alpha, \beta, i) \in \overleftarrow{\chi}^i \cup \overleftarrow{\chi}^i \mid \beta = \epsilon\}$  if  $a = \circ$ .

Note that  $f \cdot \chi_a(\alpha, \beta, i)$  is independent of  $\alpha$  when  $a = \circ^\downarrow$ , but it takes an element of  $\Pi$  as argument for the sake of uniformity of notation. The **fingerprint** of a profile  $f$ , noted  $\xi(f)$ , is an element of

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

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

We fix the set of **consistent fingerprints**, as a set of fingerprints  $\Gamma \subseteq \mathcal{F}$ . The usefulness of this set will become apparent in the reduction from XPath to the derivation problem of forest profiles in Section 6.2, but we can anticipate

that this set will represent all the profiles abstracting multi-attribute data trees that do not contradict the formula we are trying to satisfy. For the moment, however, the reader may simply consider  $\Gamma$  as a given arbitrary set of fingerprints.

### Concatenation

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

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

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

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

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

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

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

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

Notice that

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

### Rooting

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

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

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

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

### 4.3 The derivation problem

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

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

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

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

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

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

## 5. COMPUTING DERIVABLE PROFILES

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

### 5.1 Preliminaries

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

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

### 5.2 Ordering on profiles

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

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

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

We write  $f \preceq f'$  if  $f \preceq f''$  for some  $f'' \sim f'$ . We say that a set of forest profiles  $G \subseteq \mathfrak{F}$  is **upward closed** (resp. **downward**



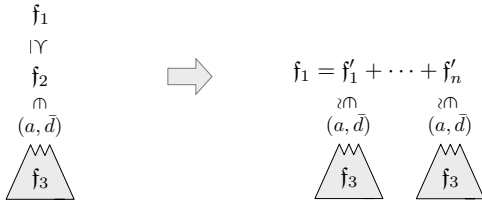


Figure 4: Statement of Lemma 5.1.

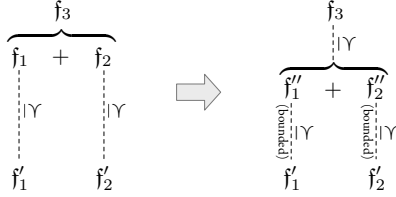


Figure 5: Statement of Lemma 5.3.

**closed**) with respect to  $\succsim$ , if for every  $f \in G$  and  $f' \succsim f$  (resp.  $f \succsim f'$ ), we have  $f' \in G$ . We write

$$\begin{aligned} \uparrow G &\stackrel{\text{def}}{=} \{f \in \mathfrak{F} \mid f \succsim f' \text{ for some } f' \in G\} \\ \downarrow G &\stackrel{\text{def}}{=} \{f \in \mathfrak{F} \mid f' \succsim f \text{ for some } f' \in G\} \end{aligned}$$

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

### 5.3 Monotonicity properties

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

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

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

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

The lemma above implies that the set of derivable profiles is upward closed. (Details in the Appendix.)

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

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

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

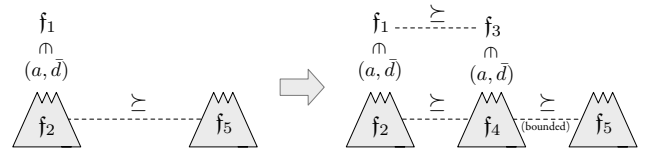


Figure 6: Statement of Lemma 5.4.

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

A similar lemma holds for the rooting operation.

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

### 5.4 Bounding the rigid values

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

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

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

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

By the Lemma just shown, we have the following

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

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

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

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

### 5.5 Restricting to minimal elements

Thanks to the result from the previous section stating that  $\mathcal{D}_b$  is upward closed, we can now show that we can work only with the minimal elements of  $\mathfrak{F}_b$ . The main necessary property concerns all those profiles  $f' \in \mathfrak{F}_b$  that are ' $\succsim$ -related' to a profile  $f'' \in \mathcal{D}_b$ , in the sense that  $f' \succsim f \succsim f'' \in \mathcal{D}_b$  for

some  $f$ . (Note that this set of profiles is precisely  $\uparrow\downarrow\mathfrak{D}_b$ .) The property states that the forest profiles algebra preserves the  $\sim$ -relatedness.

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

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

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

## 5.6 The algorithm

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

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

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

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

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

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

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

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

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

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

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

PROPOSITION 5.12. The derivation problem is decidable in  $2\text{EXPSpace}$ .

## 6. FROM XPATH TO FOREST PROFILES

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

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

### 6.1 Normal forms

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

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

#### Preliminaries

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

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

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

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

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

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

#### Direct normal form

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

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

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

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

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

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

**LEMMA 6.1 (DIRECT NORMAL FORM).** *There exists an exponential time translation that for every node expression  $\varphi \in \text{XPath}(*\leftarrow, \downarrow_*, \rightarrow^*, =)$  returns an equivalent node expression  $\psi$  in direct normal form.*

### Unnested normal form

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

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

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

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

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

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

Then, we obtain the following.

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

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

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

## 6.2 Reduction to the derivation problem

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

Let us fix  $\phi = \phi_1 \wedge \phi_2$  in direct unnested normal form, where  $\mathbb{A}$  as the finite alphabet,  $\mathbf{k}$  as the number of attributes,

$\mathbb{D}$  as any infinite domain, and  $\mathbb{A}_{\text{root}}$  is the set of all  $a \in \mathbb{A}$  that make  $\phi_1$  true.

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

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

Note that  $P_o$  is injective.

Let us define  $\mathcal{P}_\phi$  as the set of patterns consisting of

- the constant  $\top$  and the empty string  $\varepsilon$ ,
- $\psi$ , for every  $\psi \in \mathcal{B}(\mathbb{A})$  that is a subformula of  $\phi$ ,
- every  $\alpha \in (\mathcal{B}(\mathbb{A}))^*$  so that  $P_{\rightarrow^*}(\alpha)$ ,  $P_{*\leftarrow}(\alpha)$ , or  $P_{\downarrow_*}(\alpha)$  is a substring of  $\phi$ .

It follows that  $\mathcal{P}_\phi$  is finite and subword-closed.

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

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

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

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

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

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

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

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

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

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

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

### Abstractions.

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

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

as the forest profile that abstracts the forest  $\bar{\mathfrak{t}}$  in the context of the forests  $\bar{\mathfrak{t}}_l$  to the left and  $\bar{\mathfrak{t}}_r$  to the right. We have already discussed the idea of this abstraction in Section 4. For example, for the forest of Figure 3, assuming  $\mathcal{P} = \{\top, b \cdot c, b, c, \epsilon\}$ , we would obtain an abstraction where

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

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

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

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

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

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

## 7. DISCUSSION

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

### Upward axes

One natural question that stems from the result presented here is whether it can be extended to work with an upward axis as well. However, we claim (without a proof) that already  $SAT\text{-}XPath(\uparrow^*, \rightarrow^*, =)$  has a non-primitive recursive

lower bound. Indeed, this can be proved by reusing the results on lower bounds of [14]. The cited work shows that XPath with one non-reflexive transitive axis is enough to prove non-primitive recursiveness provided that the axis is functional (*i.e.*, the transitive closure of an axis like  $\rightarrow, \leftarrow, \uparrow$  but unlike  $\downarrow$ ). Here, however, we feature *reflexive*-transitive axes instead of only transitive. Therefore, in principle we cannot use this result. However, one can somehow code  $\uparrow^+$  with  $\rightarrow^*[a]\uparrow^*[-a]$  for some label  $a$ . We leave the proof of this claim for the journal version of the present work.

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

### Future work

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

## 8. REFERENCES

- [1] Vince Bárány, Mikołaj Bojańczyk, Diego Figueira, and Paweł Parys. Decidable classes of documents for XPath. In *IARCS Annual Conference on Foundations of Software Technology and Theoretical Computer Science (FSTTCS'12)*, Leibniz International Proceedings in Informatics (LIPIcs), Hyderabad, India, 2012. Leibniz-Zentrum für Informatik.
- [2] Michael Benedikt, Wenfei Fan, and Floris Geerts. XPath satisfiability in the presence of DTDs. *Journal of the ACM*, 55(2):1–79, 2008.
- [3] Mikołaj Bojańczyk and Sławomir Lasota. An extension of data automata that captures XPath. In *Annual IEEE Symposium on Logic in Computer Science (LICS '10)*, 2010.
- [4] Mikołaj Bojańczyk, Anca Muscholl, Thomas Schwentick, and Luc Segoufin. Two-variable logic on data trees and XML reasoning. *Journal of the ACM*, 56(3):1–48, 2009.
- [5] James Clark and Steve DeRose. XML path language (XPath). Website, 1999. W3C Recommendation. <http://www.w3.org/TR/xpath>.
- [6] Claire David, Leonid Libkin, and Tony Tan. Efficient reasoning about data trees via integer linear programming. *ACM Transactions on Database Systems*, 37(3):19, 2012.
- [7] Wenfei Fan, Chee Yong Chan, and Minos N. Garofalakis. Secure XML querying with security views. In *ACM SIGACT-SIGMOD-SIGART*

- International Conference on Management of Data (SIGMOD'04)*, pages 587–598. ACM Press, 2004.
- [8] Diego Figueira. Satisfiability of downward XPath with data equality tests. In *ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS'09)*, pages 197–206. ACM Press, 2009.
- [9] Diego Figueira. Forward-XPath and extended register automata on data-trees. In *International Conference on Database Theory (ICDT'10)*. ACM Press, 2010.
- [10] Diego Figueira. *Reasoning on Words and Trees with Data*. Phd thesis, Laboratoire Spécification et Vérification, ENS Cachan, France, December 2010.
- [11] Diego Figueira. A decidable two-way logic on data words. In *Annual IEEE Symposium on Logic in Computer Science (LICS'11)*, pages 365–374, Toronto, Canada, 2011. IEEE Computer Society Press.
- [12] Diego Figueira. Alternating register automata on finite data words and trees. *Logical Methods in Computer Science*, 8(1), 2012.
- [13] Diego Figueira. Decidability of downward XPath. *ACM Trans. Comput. Log.*, 13(4), 2012.
- [14] Diego Figueira and Luc Segoufin. Future-looking logics on data words and trees. In *Int. Symp. on Mathematical Foundations of Comp. Sci. (MFCS'09)*, volume 5734 of *LNCS*, pages 331–343. Springer, 2009.
- [15] Diego Figueira and Luc Segoufin. Bottom-up automata on data trees and vertical XPath. In *International Symposium on Theoretical Aspects of Computer Science (STACS'11)*, Leibniz International Proceedings in Informatics (LIPIcs). Leibniz-Zentrum für Informatik, 2011.
- [16] Floris Geerts and Wenfei Fan. Satisfiability of XPath queries with sibling axes. In *International Symposium on Database Programming Languages (DBPL'05)*, volume 3774 of *Lecture Notes in Computer Science*, pages 122–137. Springer, 2005.
- [17] Georg Gottlob, Christoph Koch, and Reinhard Pichler. Efficient algorithms for processing XPath queries. *ACM Transactions on Database Systems*, 30(2):444–491, 2005.
- [18] Marcin Jurdziński and Ranko Lazić. Alternating automata on data trees and xpath satisfiability. *ACM Trans. Comput. Log.*, 12(3):19, 2011.
- [19] Michael Kaminski and Nissim Francez. Finite-memory automata. *Theoretical Computer Science*, 134(2):329–363, 1994.
- [20] Michael Kaminski and Tony Tan. Tree automata over infinite alphabets. In *Pillars of Computer Science*, volume 4800 of *Lecture Notes in Computer Science*, pages 386–423. Springer, 2008.
- [21] Wim Martens and Frank Neven. Frontiers of tractability for typechecking simple xml transformations. *J. Comput. Syst. Sci.*, 73(3):362–390, 2007.
- [22] Frank Neven, Thomas Schwentick, and Victor Vianu. Finite state machines for strings over infinite alphabets. *ACM Trans. Comput. Log.*, 5(3):403–435, 2004.
- [23] Tony Tan. An automata model for trees with ordered data values. In *Annual IEEE Symposium on Logic in Computer Science (LICS'12)*, pages 586–595. IEEE Computer Society Press, 2012.

## APPENDIX

We organize this appendix with the missing proofs, grouped by sections on the main body. In each case we repeat the statements for the reader's convenience.

### A. APPENDIX TO SECTION 4

LEMMA A.1 (DISTRIBUTIVITY OF  $(\cdot, \cup)$ ). For every  $P, P' \subseteq \mathcal{P}$  and  $\chi, \chi' \subseteq \mathbb{D} \times \Pi$

1.  $P \cdot (\chi \cup \chi') = P \cdot \chi \cup P \cdot \chi'$ ,
2.  $(P \cup P') \cdot \chi = P \cdot \chi \cup P' \cdot \chi$ .

PROOF. We have that

$$\begin{aligned}
 & P \cdot (\chi \cup \chi') \\
 &= \chi \cup \chi' \cup \{(d, \alpha' \cdot \alpha, \beta, i) \mid (d, \alpha, \beta, i) \in \chi \cup \chi', \alpha' \in P\} && \text{(by definition of '·')} \\
 &= \chi \cup \{(d, \alpha' \cdot \alpha, \beta, i) \mid (d, \alpha, \beta, i) \in \chi, \alpha' \in P\} \cup \\
 &\quad \chi' \cup \{(d, \alpha' \cdot \alpha, \beta, i) \mid (d, \alpha, \beta, i) \in \chi', \alpha' \in P\} \\
 &= P \cdot \chi \cup P \cdot \chi'. && \text{(by definition of '·')}
 \end{aligned}$$

which proves the first statement.

On the other hand,

$$\begin{aligned}
 & (P \cup P') \cdot \chi \\
 &= \chi \cup \{(d, \alpha' \cdot \alpha, \beta, i) \mid (d, \alpha, \beta, i) \in \chi, \alpha' \in P \cup P'\} \\
 &= \chi \cup \{(d, \alpha' \cdot \alpha, \beta, i) \mid (d, \alpha, \beta, i) \in \chi, \alpha' \in P\} \cup \\
 &\quad \chi \cup \{(d, \alpha' \cdot \alpha, \beta, i) \mid (d, \alpha, \beta, i) \in \chi, \alpha' \in P'\} \\
 &= P \cdot \chi \cup P' \cdot \chi
 \end{aligned}$$

which proves the second statement.  $\square$

LEMMA A.2. If  $f_1 + f_2 = f_3$  and  $f_1, f_2 \in \mathfrak{F}$ , then  $f_3 \in \mathfrak{F}$ .

PROOF. If  $\overline{\chi}_3^E(\alpha, \beta, i) = \{d\}$ , then  $\overline{\chi}_1^E(\alpha, \beta, i) = \{d\}$ , and hence  $d \in R_1$  and therefore  $d \in R_3 = R_1$ . A symmetrical argument applies for  $\overline{\chi}_3^E$  and  $\overline{\chi}_2^E$ .  $\square$

LEMMA A.3 (DISTRIBUTIVITY OF  $(\cdot, +)$ ). For every  $f_1, f_2 \in \mathfrak{F}$  and  $\chi \subseteq \mathbb{D} \times \Pi$

- $[f_1 + f_2] \cdot \chi = [f_1] \cdot ([f_2] \cdot \chi)$ ,
- $\langle [f_1 + f_2] \cdot \chi = \langle [f_2] \cdot (\langle [f_1] \cdot \chi)$ .

PROOF. Let  $f_+ = f_1 + f_2$ .

$$\begin{aligned}
 & [f_1 + f_2] \cdot \chi \\
 &= \{\alpha \mid (d, \alpha, \beta, i) \in \overline{\chi}_+^1\} \cdot \chi && \text{(by definition of } [\ ] \cdot \text{)} \\
 &= \{\alpha \mid (d, \alpha, \beta, i) \in \overline{\chi}_1^1 \cup [f_1] \cdot \overline{\chi}_2^1\} \cdot \chi && \text{(by definition of } f_1 + f_2 \text{)} \\
 &= (\{\alpha \mid (d, \alpha, \beta, i) \in \overline{\chi}_1^1\} \cup \{\alpha \mid (d, \alpha, \beta, i) \in [f_1] \cdot \overline{\chi}_2^1\}) \cdot \chi \\
 &= ([f_1] \cdot \{\alpha \mid (d, \alpha, \beta, i) \in [f_1] \cdot \overline{\chi}_2^1\}) \cdot \chi && \text{(by } (\cdot, \cup) \text{ distributivity)} \\
 &= [f_1] \cdot \{\alpha \mid (d, \alpha, \beta, i) \in [f_1] \cdot \overline{\chi}_2^1\} \cdot \chi \\
 &= [f_1] \cdot \{(d, \alpha' \cdot \alpha, \beta, i) \mid (d, \alpha, \beta, i) \in \chi, \\
 &\quad \alpha' \in \{\alpha \mid (d, \alpha, \beta, i) \in [f_1] \cdot \overline{\chi}_2^1\}\} \\
 &= [f_1] \cdot \chi \cup \{(d, \alpha' \cdot \alpha, \beta, i) \mid (d, \alpha, \beta, i) \in \chi, \\
 &\quad \alpha' \in \{\alpha \mid (d, \alpha, \beta, i) \in \{\alpha \mid (d, \alpha, \beta, i) \in \overline{\chi}_1^1 \cdot \overline{\chi}_2^1\}\}\} \\
 &= [f_1] \cdot \chi \cup \{(d, \alpha' \cdot \alpha, \beta, i) \mid (d, \alpha, \beta, i) \in \chi, \\
 &\quad \alpha' \in \{\alpha \cdot \alpha' \mid (d, \alpha, \beta, i) \in \overline{\chi}_1^1 \wedge \\
 &\quad ((d', \alpha', \beta', i') \in \overline{\chi}_2^1 \vee \alpha' = \epsilon)\}\} \\
 &= [f_1] \cdot \chi \cup \{(d, \alpha' \cdot \alpha, \beta, i) \mid (d, \alpha, \beta, i) \in \chi, \\
 &\quad \alpha' \in \{\alpha \cdot \alpha' \mid \alpha \in [f_1] \cdot \alpha', \alpha' \in [f_2] \cdot \{\epsilon\}\}\} \\
 &= [f_1] \cdot \chi \cup \{(d, \alpha' \cdot \alpha'' \cdot \alpha, \beta, i) \mid (d, \alpha, \beta, i) \in \chi, \\
 &\quad \alpha' \in [f_1] \cdot \alpha'', \alpha'' \in [f_2] \cdot \{\epsilon\}\} \\
 &= [f_1] \cdot \chi \cup
 \end{aligned}$$

$$\begin{aligned}
& \{(d, \alpha' \cdot \alpha, \beta, i) \mid (d, \alpha, \beta, i) \in \chi, \alpha' \in [f_1] \} \cup \\
& \{(d, \alpha' \cdot \alpha'' \cdot \alpha, \beta, i) \mid (d, \alpha, \beta, i) \in \chi, \alpha' \in [f_1], \alpha'' \in [f_2] \} \\
= & [f_1] \cdot \chi \cup [f_1] \cdot \chi \cup \\
& \{(d, \alpha' \cdot \alpha'' \cdot \alpha, \beta, i) \mid (d, \alpha, \beta, i) \in \chi, \alpha' \in [f_1], \alpha'' \in [f_2] \} \\
= & [f_1] \cdot \chi \cup \\
& [f_1] \cdot \{(d, \alpha'' \cdot \alpha, \beta, i) \mid (d, \alpha, \beta, i) \in \chi, \alpha'' \in [f_2] \} \\
= & [f_1] \cdot \chi \cup [f_1] \cdot ([f_2] \cdot \chi) \\
= & [f_1] \cdot ([f_2] \cdot \chi)
\end{aligned}$$

An identical reasoning can be used to prove that  $\llbracket f_1 + f_2 \rrbracket \cdot \chi = \llbracket f_2 \rrbracket \cdot (\llbracket f_1 \rrbracket \cdot \chi)$ .  $\square$

LEMMA A.4 (+ ASSOCIATIVITY). *For every  $f_1, f_2, f_3 \in \mathfrak{F}$ ,  $(f_1 + f_2) + f_3 = f_1 + (f_2 + f_3)$ .*

PROOF. Let  $f_{12} = f_1 + f_2$ ,  $f_{23} = f_2 + f_3$ .

We first show that if  $(f_1 + f_2) + f_3$  can be applied, then  $f_1 + (f_2 + f_3)$  can also be applied. Condition (a) is immediate, since  $R_1 = R_2 = R_3$ . Now we show that conditions (b) and (c) hold for  $f_2 + f_3$ . Since  $\overrightarrow{\chi}_2^E = \overrightarrow{\chi}_{12}^E$  and since the conditions between  $f_{12} + f_3$  hold, we have that  $\overrightarrow{\chi}_2^E = \overrightarrow{\chi}_3^1 \cup [f_3] \cdot \overrightarrow{\chi}_3^E$  and hence that condition (b) holds for  $f_2 + f_3$ .

To show that condition (c) holds, note that

$$\begin{aligned}
\overleftarrow{\chi}_3^E &= \overleftarrow{\chi}_{12}^1 \cup \llbracket f_{12} \rrbracket \cdot \overleftarrow{\chi}_{12}^E && \text{(by condition (c) of } f_{12} + f_3) \\
&= \overleftarrow{\chi}_2^1 \cup \llbracket f_2 \rrbracket \cdot \overleftarrow{\chi}_1^1 \cup \llbracket f_{12} \rrbracket \cdot \overleftarrow{\chi}_1^E && \text{(by definition of } f_1 + f_2) \\
&= \overleftarrow{\chi}_2^1 \cup \llbracket f_2 \rrbracket \cdot \overleftarrow{\chi}_1^1 \cup \llbracket f_2 \rrbracket \cdot (\llbracket f_1 \rrbracket \cdot \overleftarrow{\chi}_1^E) && \text{(by } (\cdot, +) \text{ distributivity)} \\
&= \overleftarrow{\chi}_2^1 \cup \llbracket f_2 \rrbracket \cdot (\overleftarrow{\chi}_1^1 \cup \llbracket f_1 \rrbracket \cdot \overleftarrow{\chi}_1^E) && \text{(by } (\cdot, \cup) \text{ distributivity)} \\
&= \overleftarrow{\chi}_2^1 \cup \llbracket f_2 \rrbracket \cdot \overleftarrow{\chi}_2^E && \text{(by condition (c) of } f_1 + f_2)
\end{aligned}$$

Hence, we obtain that  $\overleftarrow{\chi}_3^E = \overleftarrow{\chi}_2^1 \cup \llbracket f_2 \rrbracket \cdot \overleftarrow{\chi}_2^E$ , and thus condition (c) holds for  $f_2 + f_3$ .

On the other hand, assuming that  $f_1 + (f_2 + f_3)$ , we can show that the conditions to apply  $(f_1 + f_2) + f_3$  hold, in the same way as before, since all the definitions are symmetrical, where  $\overleftarrow{\chi}_1^1 = \overleftarrow{\chi}_2^1 \cup [f_2] \cdot \overleftarrow{\chi}_2^E$  and  $\overleftarrow{\chi}_2^E = \overleftarrow{\chi}_1^1 \cup \llbracket f_1 \rrbracket \cdot \overleftarrow{\chi}_1^E$ .

We now show that  $(f_1 + f_2) + f_3 = f_1 + (f_2 + f_3)$ . Let  $f_{1(23)} = f_1 + (f_2 + f_3)$  and  $f_{(12)3} = (f_1 + f_2) + f_3$ .

We have that  $\overrightarrow{\chi}_{1(23)}^E = \overrightarrow{\chi}_{23}^E = \overrightarrow{\chi}_3^E = \overrightarrow{\chi}_{(12)3}^E$  by item condition (+2) of  $+$ . In a similar way, by condition (+3) we have that  $\overleftarrow{\chi}_{1(23)}^E = \overleftarrow{\chi}_{(12)3}^E$ . On the other hand, we have

$$\begin{aligned}
\overrightarrow{\chi}_{1(23)}^1 &= \overrightarrow{\chi}_1^1 \cup [f_1] \cdot \overrightarrow{\chi}_{23}^1 && \text{(by condition (+4))} \\
&= \overrightarrow{\chi}_1^1 \cup [f_1] \cdot (\overrightarrow{\chi}_2^1 \cup [f_2] \cdot \overrightarrow{\chi}_3^1) && \text{(by condition (+4))} \\
&= \overrightarrow{\chi}_1^1 \cup [f_1] \cdot \overrightarrow{\chi}_2^1 \cup [f_1] \cdot ([f_2] \cdot \overrightarrow{\chi}_3^1) && \text{(by } (\cdot, \cup) \text{ distributivity)} \\
&= \overrightarrow{\chi}_1^1 \cup [f_1] \cdot \overrightarrow{\chi}_2^1 \cup [f_1] \cdot \overrightarrow{\chi}_3^1 && \text{(by } (\cdot, +) \text{ distributivity)} \\
&= \overrightarrow{\chi}_{12}^1 \cup [f_1] \cdot \overrightarrow{\chi}_3^1 && \text{(by condition (+4))} \\
&= \overrightarrow{\chi}_{(12)3}^1 && \text{(by condition (+4))}
\end{aligned}$$

and by a similar reasoning using condition (+5) we obtain that  $\overleftarrow{\chi}_{1(23)}^E = \overleftarrow{\chi}_{(12)3}^E$ . Hence,  $(f_1 + f_2) + f_3 = f_1 + (f_2 + f_3)$ .  $\square$

LEMMA A.5. *For every  $f_1, f_2 \in (a, \bar{d})f_3$  and every  $\chi \subseteq \mathbb{D} \times \Pi$ ,*

$$\begin{aligned}
& [f_1] \cdot [f_2] \cdot \chi = [f_1] \cdot \chi = [f_2] \cdot \chi = \\
& \llbracket f_1 \rrbracket \cdot \llbracket f_2 \rrbracket \cdot \chi = \llbracket f_1 \rrbracket \cdot \chi = \llbracket f_2 \rrbracket \cdot \chi.
\end{aligned}$$

PROOF. By condition (c) of rooting and definition of  $[ \cdot ]$ ,  $\llbracket \cdot \rrbracket$ , we have that  $\llbracket f_1 \rrbracket = [f_1] \cdot \llbracket f_2 \rrbracket = [f_2] \cdot \llbracket f_1 \rrbracket = \sigma_a$ . By definition of  $\sigma_a$ , it is immediate that  $\sigma_a \cdot \sigma_a \cdot \chi = \sigma_a \cdot \chi$ . Hence, the statement follows.  $\square$

LEMMA A.6. *For every  $f_1 \in (a, \bar{d})f_2$ ,*

$$\overrightarrow{\chi}_1^1 = [f_1] \cdot \overrightarrow{\chi}_1^1 \quad \overleftarrow{\chi}_1^1 = [f_1] \cdot \overleftarrow{\chi}_1^1.$$

PROOF. By condition (c) of rooting,  $[f_1] \cdot \sigma_a = \sigma_a$  and since  $\overrightarrow{\chi}_1^1 = \{(d, \alpha, \beta \cdot \gamma, i) \mid \exists \alpha'. (d, \alpha', \gamma, i) \in \overrightarrow{\chi}_2^1 \cup \overrightarrow{\chi}_2^1, \alpha, \beta \in \sigma_a\} \cup \bigcup_{i \in [k]} (\{d(i)\} \times (\sigma_a \setminus \{\epsilon\}) \times \sigma_a \times \{i\})$  it follows that  $\sigma_a \cdot \overrightarrow{\chi}_1^1 = \overrightarrow{\chi}_1^1$ , obtaining that  $\overrightarrow{\chi}_1^1 = [f_1] \cdot \overrightarrow{\chi}_1^1$ . Since  $\overleftarrow{\chi}_1^1 = \overleftarrow{\chi}_1^1$  by condition (c), we also obtain that  $\overleftarrow{\chi}_1^1 = [f_1] \cdot \overleftarrow{\chi}_1^1$ .  $\square$

## B. APPENDIX TO SECTION 5

LEMMA B.1.  $\preceq$  and  $\lesssim$  are partial orders over  $\mathfrak{F}$ .

PROOF. If  $f \preceq f' \preceq f''$ , suppose  $f' = f[d_1 \mapsto d_1, d'_1] \cdots f[d_n \mapsto d_n, d'_n]$  and  $f'' = f'[e_1 \mapsto e_1, e'_1] \cdots f'[e_m \mapsto e_m, e'_m]$ . Therefore,  $f'' = f[d_1 \mapsto d_1, d'_1] \cdots f[d_n \mapsto d_n, d'_n][e_1 \mapsto e_1, e'_1] \cdots f[e_m \mapsto e_m, e'_m]$  and hence  $f \preceq f''$ .

On the other hand, it is plain that  $\sim$  preserves transitivity and reflexivity, and that  $f \preceq f$ .  $\square$

LEMMA B.2. *If  $f_1 \preceq f_2$  then  $\llbracket f_1 \rrbracket = \llbracket f_2 \rrbracket$  and  $\llbracket f_1 \rrbracket = \llbracket f_2 \rrbracket$ .*

PROOF. This is immediate from the definition of  $\llbracket \cdot \rrbracket$ ,  $\llbracket \cdot \rrbracket$ , and  $\preceq$ . Note that, since  $f_1 \preceq f_2$  we have  $\{(\alpha, \beta, i) \mid \exists d \in \mathbb{D}.(d, \alpha, \beta, i) \in f_1\} = \{(\alpha, \beta, i) \mid \exists d \in \mathbb{D}.(d, \alpha, \beta, i) \in f_2\}$ . Hence, we obtain that  $\alpha \in \llbracket f_1 \rrbracket$  if and only if there is some  $(d, \alpha, \beta, i) \in \overline{\chi}_1^1$  if and only if there is some  $(d', \alpha, \beta, i) \in \overline{\chi}_2^1$  if and only if  $\alpha \in \llbracket f_2 \rrbracket$ . A similar argument shows that  $\llbracket f_1 \rrbracket = \llbracket f_2 \rrbracket$ .  $\square$

LEMMA B.3. *If  $f \preceq f'$  then  $\xi(f) = \xi(f')$ .*

PROOF. Suppose, without any loss of generality, that  $f' = f[d \mapsto d, d']$ , where  $d \in \text{data}(f) \setminus R$  and  $d' \notin \text{data}(f)$ . Since  $d'$  behaves just as  $d$ , we have that, for any  $\bar{\alpha}, \bar{\alpha}' \in \Pi$ ,  $a, a' \in \{\overleftarrow{\circ}, \overrightarrow{\circ}, \circ^\downarrow, \circ\}$ ,  $d' \in f'.\chi_a(\bar{\alpha}) \cap \chi_{a'}(\bar{\alpha}')$  iff  $d \in f.\chi_a(\bar{\alpha}) \cap f.\chi_{a'}(\bar{\alpha}')$ . Hence,  $f'.\chi_a(\bar{\alpha}) \cap f'.\chi_{a'}(\bar{\alpha}') \neq \emptyset$  iff  $f.\chi_a(\bar{\alpha}) \cap f.\chi_{a'}(\bar{\alpha}') \neq \emptyset$ .

Further, since  $d \notin R$  we have that if  $d \in \chi_a(\bar{\alpha})$  in  $f'$  (or in  $f$ ) then  $|\chi_a(\bar{\alpha})| > 1$ . Hence  $\xi(f)(\bar{\alpha}, a) = \xi(f')(\bar{\alpha}, a) = 2+$  whenever  $d \in f.\chi_a(\bar{\alpha})$ . Further, since  $f'$  preserves all the data values of  $f$  with their descriptions,  $\xi(f)(\bar{\alpha}, a) = \xi(f')(\bar{\alpha}, a)$  whenever  $d \notin f.\chi_a(\bar{\alpha}) \cap f.\chi_{a'}(\bar{\alpha}')$ .  $\square$

LEMMA B.4. *If  $d \in \text{data}(f) \setminus R$  and  $d' \notin \text{data}(f)$ , we have that if  $f \in \mathfrak{F}$  then  $f[d \mapsto d, d'] \in \mathfrak{F}$ .*

PROOF. This is immediate from the definition of validity. Since  $d' \notin \text{data}(f)$ , then  $f' = f[d \mapsto d, d']$  is the result of adding some fresh data value  $d' \notin R$  with the same description as  $d$ . Since  $f' \notin \mathfrak{F}$ , there is some  $\{d''\} = \overline{\chi}^{e'}(\bar{\alpha})$  for  $\bar{\alpha} \in \Pi$  and  $d'' \notin R$ . By definition of  $f'$ , this means that  $d'' \neq d$  and  $d'' \neq d$ , and hence that  $\{d''\} = \overline{\chi}^e(\bar{\alpha})$ , and thus  $f \notin \mathfrak{F}$ , which is a contradiction. A symmetric argument can also be applied to  $\overline{\chi}^e$ .  $\square$

## Appendix to Section 5.3

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

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

PROOF OF LEMMA 5.1. Since  $f_1 \succeq f_2$ , suppose

$$f_1 = f_2[d_1 \mapsto d_1, d'_1] \cdots [d_t \mapsto d_t, d'_t] \quad (\dagger)$$

where  $\{d_1, \dots, d_t\} \subseteq \text{data}(f_2) \setminus R_2$ ;  $\{d'_1, \dots, d'_t\} \cap \text{data}(f_2) = \emptyset$ ; and  $d'_i \neq d'_j$  for all  $i \neq j$ . We define  $n = t + 2$ .

The idea is that we will define  $f'_0, f'_1, \dots, f'_{t+1}$  so that  $f'_i$  is, modulo some renaming of data values, the same as  $f_2$  for all strictly internal values, and with the same descriptions. However,  $f'_i$  has more *external* data values than  $f_2$ —namely  $d'_1, \dots, d'_t$ —that have basically the same descriptions as  $d_1, \dots, d_t$  in  $f_2$ . Although technically we could make use of less profiles, we prefer to define  $t + 1$  different profiles, to preserve some symmetries in the definitions, and thus to simplify proofs. We define  $f'_0, \dots, f'_{t+1}$  in Figure ??.

We shall now show that for every  $f'_i$  there is some  $f''_i$  so that  $f'_i \sim f''_i \in (a, \bar{d})f_3$ .

CLAIM B.4.1. *For every  $i \in [n]$ ,  $f'_i \in (a, \bar{d})f_3$ .*

PROOF. First note that  $f'_0, f'_{t+1} \in (a, \bar{d})f_3$ , since all the internal descriptions of the internal data values are preserved, and  $\xi(f'_0) = \xi(f'_{t+1}) = \xi(f_2)$ . The fact that they all have the same fingerprints is because:

- For every  $a \in \mathcal{A}$ ,  $\bar{\alpha} \in \Pi$ , and every data value  $e \in \text{data}(f_2)$ , we have that  $e \in \chi_a(\bar{\alpha})$  in  $f_2$  iff  $e \in \chi_a(\bar{\alpha})$  in  $f'_0$  iff  $e \in \chi_a(\bar{\alpha})$  in  $f'_{t+1}$ . This is because, since  $\overleftarrow{\chi}_{t+1}'(e) = \overleftarrow{\chi}_0'(e) = \overleftarrow{\chi}_2'(e)$ ,  $\overleftarrow{\chi}_{t+1}'(e) = \overleftarrow{\chi}_0'(e) = \overleftarrow{\chi}_2'(e)$ ,  $\overleftarrow{\chi}_0'(e) = \overleftarrow{\chi}_2'(e)$ ,  $\overleftarrow{\chi}_0'(e) = \overleftarrow{\chi}_2'(e)$ ,  $\overleftarrow{\chi}_{t+1}'(e) = \overleftarrow{\chi}_2^e(e)$ , we have that: for all  $a \in \{\circ, \circ^\downarrow, \overleftarrow{\circ}\}$ ,  $\chi_a(e)$  in  $f'_0$  is equal to  $\chi_a(e)$  in  $f_2$ ; and for all  $a \in \{\circ, \circ^\downarrow, \overleftarrow{\circ}\}$ ,  $\chi_a(e)$  in  $f'_{t+1}$  is equal to  $\chi_a(e)$  in  $f_2$ . And further,

– if  $a = \overrightarrow{\circ}$ ,

$$\begin{aligned} f'_0.\chi_a(e) &= \overleftarrow{\chi}_0'(e) \cup \llbracket f'_0 \rrbracket \cdot \overleftarrow{\chi}_0^{e'}(e) \\ &= \overleftarrow{\chi}_2'(e) \cup \llbracket f_2 \rrbracket \cdot (\overleftarrow{\chi}_2'(e) \cup \llbracket f_2 \rrbracket \cdot \overleftarrow{\chi}_2^e(e)) \\ &= \overleftarrow{\chi}_2'(e) \cup \llbracket f_2 \rrbracket \cdot \overleftarrow{\chi}_2^e(e) && \text{(by Lemmas ?? and ??)} \\ &= f_2.\chi_a(e) \end{aligned}$$

– if  $a = \overleftarrow{\circ}$ ,

$$f'_{t+1}.\chi_a(e) = \overleftarrow{\chi}_0'(e) \cup \llbracket f'_{t+1} \rrbracket \cdot \overleftarrow{\chi}_{t+1}^{e'}(e)$$



$$f'_0 : \begin{cases} R'_0 = R_2, \\ \bar{\chi}'_0(e) = (\bar{\chi}_2^E(e), \bar{\chi}_2^1(e), \bar{\chi}_2^1(e), \bar{\chi}_2^1(e) \cup [f_2] \cdot \bar{\chi}_2^E(e)) \\ \quad \text{for every } e \in \text{data}(f_2), \\ \bar{\chi}'_0(d'_j) = (\bar{\chi}_2^E(d_j), \emptyset, \emptyset, \bar{\chi}_2^1(d_j) \cup [f_2] \cdot \bar{\chi}_2^E(d_j)) \\ \quad \text{for every } j \in [t]. \end{cases}$$

For every  $i \in \{1, \dots, t\}$  we define  $f'_i$  as

$$f'_i : \begin{cases} R'_i = R_2, \\ \bar{\chi}'_i(e) = (\bar{\chi}_2^1(e) \cup \llbracket f_2 \rrbracket \cdot \bar{\chi}_2^E(e), \bar{\chi}_2^1(e), \\ \quad \bar{\chi}_2^1(e), \bar{\chi}_2^1(e) \cup [f_2] \cdot \bar{\chi}_2^E(e)) \\ \quad \text{for every } e \in \text{data}(f_2) \setminus \{d_i\}, \\ \bar{\chi}'_i(d_i) = (\bar{\chi}_2^1(d_i) \cup \llbracket f_2 \rrbracket \cdot \bar{\chi}_2^E(d_i), \emptyset, \\ \quad \emptyset, \bar{\chi}_2^1(d_i) \cup [f_2] \cdot \bar{\chi}_2^E(d_i)), \\ \bar{\chi}'_i(d'_j) = (\llbracket f_2 \rrbracket \cdot \bar{\chi}_2^E(d_j), \bar{\chi}_2^1(d_j), \bar{\chi}_2^1(d_j), [f_2] \cdot \bar{\chi}_2^E(d_j)), \\ \bar{\chi}'_i(d'_j) = (\bar{\chi}_2^1(d_j) \cup \llbracket f_2 \rrbracket \cdot \bar{\chi}_2^E(d_j), \emptyset, \emptyset, [f_2] \cdot \bar{\chi}_2^E(d_j)) \\ \quad \text{for every } j < i, \\ \bar{\chi}'_i(d'_j) = (\llbracket f_2 \rrbracket \cdot \bar{\chi}_2^E(d_j), \emptyset, \emptyset, \bar{\chi}_2^1(d_j) \cup [f_2] \cdot \bar{\chi}_2^E(d_j)) \\ \quad \text{for every } j > i. \end{cases}$$

$$f'_{t+1} : \begin{cases} R'_{t+1} = R_2, \\ \bar{\chi}'_{t+1}(e) = (\bar{\chi}_2^1(e) \cup \llbracket f_2 \rrbracket \cdot \bar{\chi}_2^E(e), \bar{\chi}_2^1(e), \bar{\chi}_2^1(e), \bar{\chi}_2^E(e)) \\ \quad \text{for every } e \in \text{data}(f_2), \\ \bar{\chi}'_{t+1}(d'_j) = (\bar{\chi}_2^1(d_j) \cup \llbracket f_2 \rrbracket \cdot \bar{\chi}_2^E(d_j), \emptyset, \emptyset, \bar{\chi}_2^E(d_j)) \\ \quad \text{for every } j \in [t]. \end{cases}$$

**Figure 7: Definition of  $f'_0$ ,  $f'_{t+1}$  and  $f'_i$  for all  $i \in [t]$ .**

$$\begin{aligned} &= \bar{\chi}_2^1(e) \cup \llbracket f_2 \rrbracket \cdot (\bar{\chi}_2^1(e) \cup \llbracket f_2 \rrbracket \cdot \bar{\chi}_2^E(e)) \\ &= \bar{\chi}_2^1(e) \cup \llbracket f_2 \rrbracket \cdot \bar{\chi}_2^E(e) && \text{(by Lemmas ?? and ??)} \\ &= f_2 \cdot \chi_a(e). \end{aligned}$$

- The data values  $d'_j$  do not change the fingerprint  $\xi(f_2)$ , because  $d'_j$  is added to all those  $\chi_a(\bar{\alpha})$  where  $d_j \in f_2 \cdot \chi_a(\bar{\alpha})$ , and  $f_2 \cdot \chi_a(\bar{\alpha})$  has already two data values (otherwise  $d_j$  would be rigid). Therefore, for any  $a \in \mathcal{A}$ ,  $\bar{\alpha} \in \Pi$ , if some  $d'_j \in f'_0 \cdot \chi_a(\bar{\alpha})$ , then  $|f_2 \cdot \chi_a(\bar{\alpha})| \geq 2$ , and  $f_2 \cdot \chi_a(\bar{\alpha}) \subseteq f'_0 \cdot \chi_a(\bar{\alpha})$  (thus,  $|f'_0 \cdot \chi_a(\bar{\alpha})| \geq 2$ ). Otherwise, if  $d'_j \notin f'_0 \cdot \chi_a(\bar{\alpha})$  for every  $j \in [t]$ , then  $f'_0 \cdot \chi_a(\bar{\alpha}) = f_2 \cdot \chi_a(\bar{\alpha})$ . Moreover, for any  $a, a' \in \mathcal{A}$ ,  $\bar{\alpha}, \bar{\alpha}' \in \Pi$ , if  $d'_j$  is in some  $f'_0 \cdot \chi_a(\bar{\alpha}) \cap f'_0 \cdot \chi_{a'}(\bar{\alpha}')$  it means that  $d_j \in f_2 \cdot \chi_a(\bar{\alpha}) \cap f_2 \cdot \chi_{a'}(\bar{\alpha}')$ . Hence,  $|f'_0 \cdot \chi_a(\bar{\alpha}) \cap f'_0 \cdot \chi_{a'}(\bar{\alpha}')| \geq 1$  iff  $|f_2 \cdot \chi_a(\bar{\alpha}) \cap f_2 \cdot \chi_{a'}(\bar{\alpha}')| \geq 1$ . As a result the fingerprint of  $f_2$  and  $f'_0$  are equal, and hence  $\xi(f'_0) = \xi(f_2) \in \Gamma$ . The same reasoning applies to  $f'_{t+1}$ .

Thus,  $f'_0, f'_{t+1} \in (a, \bar{d})f_3$ .

We now show that for every  $i \in [t]$ , there is  $f''_i$  so that  $f'_i \sim f''_i \in (a, \bar{d})f_3$ . We define  $f''_i$  as  $f''_i[d_i \mapsto d'_i]$  for every  $i \in [t]$ . Observe that  $f''_i \sim f'_i$  for every  $i \in [t]$ . Since every  $d_i$  is not rigid,  $f'_i \in \mathfrak{F}$  and hence  $f''_i \in \mathfrak{F}$ . Further, note that  $f''_i$  and  $f_2$  share the same set of internal data values, and that the internal descriptions of these internal data values are the same in  $f''_i$  and  $f_2$ . Also in this case we have that  $\xi(f''_i) = \xi(f_2)$  for the same reason as before:

- For all  $e \in \text{data}(f_2)$ ,  $a \in \mathcal{A}$  and  $\bar{\alpha} \in \Pi$  we have that  $e \in \chi_a(\bar{\alpha})$  in  $f_2$  iff  $e \in \chi_a(\bar{\alpha})$  in  $f''_i$ .
- The fresh data values  $d'_j$  only add data values to sets  $\chi_a(\bar{\alpha})$  that already have at least 2 data values. And  $d'_j \in f'_i \cdot \chi_a(\bar{\alpha}) \cap f'_i \cdot \chi_{a'}(\bar{\alpha}')$  iff  $d_j \in f_2 \cdot \chi_a(\bar{\alpha}) \cap f_2 \cdot \chi_{a'}(\bar{\alpha}')$ .

This shows that  $\xi(f'_i) = \xi(f_2) \in \Gamma$ .

Thus,  $f''_i \in (a, \bar{d})f_3$  since  $f_2 \in (a, \bar{d})f_3$ .  $\square$

We now check that  $f_1 = f'_0 + \dots + f'_{t+1}$ .

CLAIM B.4.2.  $f_1 = f'_0 + \dots + f'_{t+1}$ .

PROOF. For each data value  $d$ , we are going to show that conditions (b) and (c) of concatenation hold, and that  $\bar{\chi}_1(d)$  equals to the profile of  $d$  in  $f'_1 + \dots + f'_{t+1}$ .

- Let us take any data value  $d \in \text{data}(f_2) \setminus \{d_1, \dots, d_t\}$ . We first check that condition (b) holds for  $d$ , in other words, that for every  $i \in \{0, \dots, t\}$ ,

$$\bar{\chi}_{i+1}^1(d) \cup [f_{i+1}] \cdot \bar{\chi}_{i+1}^E(d) = \bar{\chi}_i^E(d).$$

We have

$$\begin{aligned}
& \overrightarrow{\chi}_{i+1}^1(d) \cup [f_{i+1}] \cdot \overrightarrow{\chi}_{i+1}^E(d) \\
&= \overrightarrow{\chi}_{i+1}^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_{i+1}^E(d) && \text{(by Lemma ??)} \\
&= \overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_{i+1}^E(d) && \text{(by definition)} \\
&= \begin{cases} \overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_2^E(d) & \text{if } i = t, \text{ or} \\ \overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot (\overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_2^E(d)) & \text{otherwise.} \end{cases}
\end{aligned}$$

If  $i = t$ , we further have that  $\overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_2^E(d) = \overrightarrow{\chi}_i^E(d)$ , verifying (b). Otherwise, if  $i \neq t$ ,

$$\begin{aligned}
& \overrightarrow{\chi}_{i+1}^1(d) \cup [f_{i+1}] \cdot \overrightarrow{\chi}_{i+1}^E(d) = \dots = \\
&= \overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot (\overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_2^E(d)) \\
&= \overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot [f_2] \cdot \overrightarrow{\chi}_2^E(d) && \text{(by } (\cdot, \cup) \text{ distributivity)} \\
&= \overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_2^E(d) && \text{(by Lemma ??)} \\
&= [f_2] \cdot \overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_2^E(d) && \text{(since } \overrightarrow{\chi}_2^1(d) \subseteq [f_2] \cdot \overrightarrow{\chi}_2^1(d)) \\
&= \overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_2^E(d) && \text{(by Lemma ??)} \\
&= \overrightarrow{\chi}_i^E(d). && \text{(by definition)}
\end{aligned}$$

In any case, condition (b) holds between  $f'_i$  and  $f'_{i+1}$  for the data value  $d$ . Note that since all definitions are symmetrical, it also follows that condition (c) holds for  $d$ .

We now check that, if we call  $f_+$  to  $f'_1 + \dots + f'_n$ , then  $\overrightarrow{\chi}_+(d) = \overrightarrow{\chi}_1(d)$ . Here we only deal with  $\overrightarrow{\chi}_+^E$  and  $\overrightarrow{\chi}_+^1$  because the cases for  $\overrightarrow{\chi}_+^E$  and  $\overrightarrow{\chi}_+^1$  are symmetrical. By definition of  $+$ , we have that

$$\begin{aligned}
& \overrightarrow{\chi}_+^E(d) = \overrightarrow{\chi}_2^E(d) && \text{(by (+2))} \\
&= \overrightarrow{\chi}_1^E(d), && \text{(by (??), since } d \in \text{data}(f_2)) \\
& \overrightarrow{\chi}_+^1(d) = (\overrightarrow{\chi}_2^1 \cup [f_2] \cdot \overrightarrow{\chi}_2^1 \cup [f_2] \cdot [f_2] \cdot \overrightarrow{\chi}_2^1 \cup \dots \cup \\
&\quad \underbrace{[f_2] \cdot \dots \cdot [f_2] \cdot \overrightarrow{\chi}_2^1}_{n-1 \text{ times}}(d)) && \text{(by (+4))} \\
&= ([f_2] \cdot \overrightarrow{\chi}_2^1)(d) && \text{(by Lemma ??)} \\
&= \overrightarrow{\chi}_2^1(d) && \text{(by Lemma ??)} \\
&= \overrightarrow{\chi}_1^E(d). && \text{(by (??), since } d \in \text{data}(f_2))
\end{aligned}$$

• Suppose now that we have  $d = d_j$  for some  $j \in [t]$ . We check that condition (b) holds for  $d$ , in other words that for every  $i \in \{0, \dots, t\}$ ,

$$\overrightarrow{\chi}_{i+1}^1(d) \cup [f_{i+1}] \cdot \overrightarrow{\chi}_{i+1}^E(d) = \overrightarrow{\chi}_i^E(d).$$

We have

$$\overrightarrow{\chi}_{i+1}^1(d) \cup [f_{i+1}] \cdot \overrightarrow{\chi}_{i+1}^E(d) = \overrightarrow{\chi}_{i+1}^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_{i+1}^E(d) \quad \text{(by Lemma ??)}$$

If  $i = t$ ,

$$\begin{aligned}
& \overrightarrow{\chi}_{i+1}^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_{i+1}^E(d) = \overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_2^E(d) && \text{(by definition)} \\
&= \overrightarrow{\chi}_i^E(d). && \text{(by definition)}
\end{aligned}$$

If  $i \neq t$  and  $j = i + 1$ ,

$$\begin{aligned}
& \overrightarrow{\chi}_{i+1}^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_{i+1}^E(d) = \emptyset \cup [f_2] \cdot (\overrightarrow{\chi}_2^1 \cup [f_2] \cdot \overrightarrow{\chi}_2^E)(d) && \text{(by definition)} \\
&= [f_2] \cdot \overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot [f_2] \cdot \overrightarrow{\chi}_2^E(d) && \text{(by } (\cdot, \cup) \text{ distributivity)} \\
&= [f_2] \cdot \overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_2^E(d) && \text{(by Lemma ??)} \\
&= \overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_2^E(d) && \text{(by Lemma ??)} \\
&= \overrightarrow{\chi}_i^E(d). && \text{(by definition)}
\end{aligned}$$

If  $i \neq t$  and  $j \neq i + 1$ ,

$$\overrightarrow{\chi}_{i+1}^1(d) \cup [f_2] \cdot \overrightarrow{\chi}_{i+1}^E(d) = \overrightarrow{\chi}_2^1(d) \cup [f_2] \cdot (\overrightarrow{\chi}_2^1 \cup [f_2] \cdot \overrightarrow{\chi}_2^E)(d) \quad \text{(by definition)}$$

$$\begin{aligned}
&= \bar{\chi}_2^1(d) \cup [f_2] \cdot \bar{\chi}_2^1(d) \cup [f_2] \cdot [f_2] \cdot \bar{\chi}_2^E(d) && \text{(by } (\cdot, \cup) \text{ distributivity)} \\
&= \bar{\chi}_2^1(d) \cup [f_2] \cdot \bar{\chi}_2^1(d) \cup [f_2] \cdot \bar{\chi}_2^E(d) && \text{(by Lemma ??)} \\
&= \bar{\chi}_2^1(d) \cup \bar{\chi}_2^1(d) \cup [f_2] \cdot \bar{\chi}_2^E(d) && \text{(by Lemma ??)} \\
&= \bar{\chi}_2^1(d) \cup [f_2] \cdot \bar{\chi}_2^E(d) \\
&= \bar{\chi}_i^E(d). && \text{(by definition)}
\end{aligned}$$

Condition (c) follows by symmetry.

We check that  $\bar{\chi}_+^E(d) = \bar{\chi}_1^E(d)$  and  $\bar{\chi}_+^1(d) = \bar{\chi}_1^1(d)$ . By definition of  $+$ , we have that

$$\begin{aligned}
\bar{\chi}_+^E(d) &= \bar{\chi}_2^E(d) && \text{(by (+2))} \\
&= \bar{\chi}_1^E(d), && \text{(by (??), since } d \in \text{data}(f_2)) \\
\bar{\chi}_+^1(d) &= (\bar{\chi}_2^1 \cup [f_2] \cdot \emptyset \cup \dots \cup \underbrace{[f_2] \cdot \dots \cdot [f_2]}_{n-2 \text{ times}} \cdot \emptyset \cup \\
&\quad \underbrace{[f_2] \cdot \dots \cdot [f_2]}_{n-1 \text{ times}} \cdot \bar{\chi}_2^1)(d) && \text{(by (+4))} \\
&= (\bar{\chi}_2^1 \cup [f_2] \cdot \bar{\chi}_2^1)(d) && \text{(by Lemma ??)} \\
&= ([f_2] \cdot \bar{\chi}_2^1)(d) && \text{(since } \bar{\chi}_2^1 \subseteq [f_2] \cdot \bar{\chi}_2^1) \\
&= \bar{\chi}_2^1(d) && \text{(by Lemma ??)} \\
&= \bar{\chi}_1^1(d). && \text{(by (??), since } d \in \text{data}(f_2))
\end{aligned}$$

• Finally, suppose  $d = d'_j$  for some  $j \in [t]$ . We check that condition (b) holds for  $d$ , in other words that for every  $i \in \{0, \dots, t\}$ ,

$$\bar{\chi}_{i+1}^1(d'_j) \cup [f_{i+1}] \cdot \bar{\chi}_{i+1}^E(d'_j) = \bar{\chi}_i^E(d'_j).$$

By definition, we have

$$\bar{\chi}_{i+1}^1(d'_j) \cup [f_{i+1}] \cdot \bar{\chi}_{i+1}^E(d'_j) = \bar{\chi}_{i+1}^1(d'_j) \cup [f_2] \cdot \bar{\chi}_{i+1}^E(d'_j) \quad \text{(by Lemma ??)}$$

If  $i = t$ ,

$$\begin{aligned}
\bar{\chi}_{i+1}^1(d'_j) \cup [f_2] \cdot \bar{\chi}_{i+1}^E(d'_j) &= \emptyset \cup [f_2] \cdot \bar{\chi}_2^E(d'_j) && \text{(by definition)} \\
&= \bar{\chi}_i^E(d'_j). && \text{(by definition)}
\end{aligned}$$

If  $i \neq t$  and  $j = i + 1$ ,

$$\begin{aligned}
\bar{\chi}_{i+1}^1(d'_j) \cup [f_2] \cdot \bar{\chi}_{i+1}^E(d'_j) &= \bar{\chi}_2^1(d'_j) \cup [f_2] \cdot [f_2] \cdot \bar{\chi}_2^E(d'_j) && \text{(by definition)} \\
&= \bar{\chi}_2^1(d'_j) \cup [f_2] \cdot \bar{\chi}_2^E(d'_j) && \text{(by Lemma ??)} \\
&= \bar{\chi}_i^E(d'_j). && \text{(by definition)}
\end{aligned}$$

If  $i \neq t$  and  $j > i + 1$ ,

$$\begin{aligned}
\bar{\chi}_{i+1}^1(d'_j) \cup [f_2] \cdot \bar{\chi}_{i+1}^E(d'_j) &= \emptyset \cup [f_2] \cdot (\bar{\chi}_2^1(d'_j) \cup [f_2] \cdot \bar{\chi}_2^E(d'_j)) && \text{(by definition)} \\
&= [f_2] \cdot \bar{\chi}_2^1(d'_j) \cup [f_2] \cdot [f_2] \cdot \bar{\chi}_2^E(d'_j) && \text{(by } (\cdot, \cup) \text{ distributivity)} \\
&= [f_2] \cdot \bar{\chi}_2^1(d'_j) \cup [f_2] \cdot \bar{\chi}_2^E(d'_j) && \text{(by Lemma ??)} \\
&= \bar{\chi}_2^1(d'_j) \cup [f_2] \cdot \bar{\chi}_2^E(d'_j) && \text{(by Lemma ??)} \\
&= \bar{\chi}_i^E(d'_j). && \text{(by definition)}
\end{aligned}$$

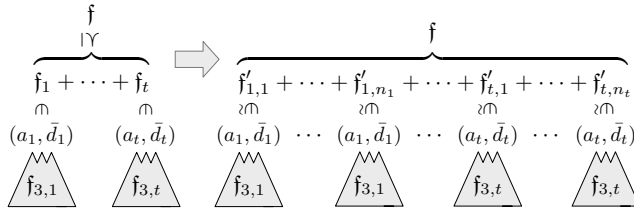
If  $i \neq t$  and  $j < i + 1$ ,

$$\begin{aligned}
\bar{\chi}_{i+1}^1(d'_j) \cup [f_2] \cdot \bar{\chi}_{i+1}^E(d'_j) &= \emptyset \cup [f_2] \cdot [f_2] \cdot \bar{\chi}_2^E(d'_j) && \text{(by definition)} \\
&= [f_2] \cdot \bar{\chi}_2^E(d'_j) && \text{(by Lemma ??)} \\
&= \bar{\chi}_i^E(d'_j). && \text{(by definition)}
\end{aligned}$$

Condition (c) follows by symmetry.

We check that  $\bar{\chi}_+^E(d'_j) = \bar{\chi}_1^E(d'_j)$  and  $\bar{\chi}_+^1(d'_j) = \bar{\chi}_1^1(d'_j)$ . By definition of  $+$ , we have that

$$\begin{aligned}
\bar{\chi}_+^E(d'_j) &= \bar{\chi}_2^E(d'_j) && \text{(by (+2))} \\
&= \bar{\chi}_1^E(d'_j), && \text{(since } \bar{\chi}_1^E(d'_j) = \bar{\chi}_2^E(d'_j) \text{ by (??))}
\end{aligned}$$



**Figure 8: Statement of Lemma ??.**

$$\begin{aligned}
\bar{\chi}_+^1(d'_j) &= (\emptyset \cup [f_2] \cdot \emptyset \cup \dots \cup \underbrace{[f_2] \cdot \dots \cdot [f_2]}_{j \text{ times}} \cdot \emptyset \cup \\
&\quad \underbrace{[f_2] \cdot \dots \cdot [f_2]}_{j+1 \text{ times}} \cdot \bar{\chi}_2^1 \cup \\
&\quad \underbrace{[f_2] \cdot \dots \cdot [f_2]}_{j+2 \text{ times}} \cdot \emptyset \cup \dots \cup \underbrace{[f_2] \cdot \dots \cdot [f_2]}_{n-1 \text{ times}} \cdot \emptyset)(d_j) \quad (\text{by (+4)}) \\
&= ([f_2] \cdot \bar{\chi}_2^1)(d_j) \quad (\text{by Lemma ??}) \\
&= \bar{\chi}_2^1(d_j) \quad (\text{by Lemma ??}) \\
&= \bar{\chi}_1^1(d'_j). \quad (\text{since } \bar{\chi}_1(d'_j) = \bar{\chi}_2(d_j) \text{ by (??)})
\end{aligned}$$

Therefore, we have that condition (a) as well as conditions (b) and (c) hold true for every data value, and that  $f'_0 + \dots + f'_{t+1}$  and  $f_1$  coincide. Thus, the Claim follows.  $\square$

This concludes our proof.

As a corollary of the previous lemma, we obtain the following.

**LEMMA B.5** (FIGURE ??). *For every  $f \succeq f_1 + \dots + f_t$  with  $f_i \in (a_i, \bar{d}_i)f_{3,i}$ , there are  $f'_{i,j} \in (a_i, \bar{d}_i)f_{3,i}$  for every  $i \in [t]$ ,  $j \in [n_i]$  so that*

$$f = f'_{1,1} + \dots + f'_{1,n_1} + \dots + f'_{t,1} + \dots + f'_{t,n_t}.$$

**PROOF OF LEMMA ??.** This is a direct consequence of Lemma 5.1. Since  $f \succeq f_1 + \dots + f_t$ , then

$$f = (f_1 + \dots + f_t)[d_1 \mapsto d_1, d'_1] \dots [d_n \mapsto d_n, d'_n]$$

for some  $d_1, \dots, d_n, d'_1, \dots, d'_n \in \mathbb{D}$  so that  $\{d_1, \dots, d_n\} \cap R = \emptyset$  and  $\{d'_1, \dots, d'_n\} \cap \text{data}(f_1 + \dots + f_t) = \emptyset$ . It follows that

$$\begin{aligned}
f &= (f_1 + \dots + f_t)[d_1 \mapsto d_1, d'_1] \dots [d_n \mapsto d_n, d'_n] \\
&= f_1[d_1 \mapsto d_1, d'_1] \dots [d_n \mapsto d_n, d'_n] + \dots + \\
&\quad f_t[d_1 \mapsto d_1, d'_1] \dots [d_n \mapsto d_n, d'_n].
\end{aligned}$$

Note that  $f_i[d_1 \mapsto d_1, d'_1] \dots [d_n \mapsto d_n, d'_n] \succeq f_i \in (a_i, \bar{d}_i)f_{3,i}$  for every  $i \in [t]$ . We can hence apply Lemma 5.1, and we obtain, for some  $n_i \in \mathbb{N}$ ,  $f'_{i,1}, \dots, f'_{i,n_i} \in (a_i, \bar{d}_i)f_{3,i}$  so that  $f'_{i,1} + \dots + f'_{i,n_i} = f_i[d_1 \mapsto d_1, d'_1] \dots [d_n \mapsto d_n, d'_n]$ . Therefore,

$$\begin{aligned}
f &= f_1[d_1 \mapsto d_1, d'_1] \dots [d_n \mapsto d_n, d'_n] + \dots + \\
&\quad f_t[d_1 \mapsto d_1, d'_1] \dots [d_n \mapsto d_n, d'_n] \\
&= f'_{1,1} + \dots + f'_{1,n_1} + \dots + f'_{t,1} + \dots + f'_{t,n_t},
\end{aligned}$$

which concludes the proof.  $\square$

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

**PROOF OF LEMMA 5.2.** This is a direct consequence of Lemma ???.  $\square$

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

**PROOF OF LEMMA 5.3.** For every description  $\bar{\pi}$  so that  $\bar{\chi}_3(\bar{\pi}) \setminus R_3 \neq \emptyset$ , let  $d_{\bar{\pi}} \in \mathbb{D}$  be a data value so that  $d_{\bar{\pi}} \in \bar{\chi}_3(\bar{\pi}) \setminus R_3$ . Let  $D$  be the set of all these data values. Note that  $|D| \leq |\Pi|^4$ .

We define  $f''_i = (R_3, \bar{\chi}''_i)$ , where

$$\bar{\chi}''_i(d) = \begin{cases} \bar{\chi}'_i(d) & \text{if } d \notin D \\ \bar{\chi}_i(d) & \text{if } d \in D. \end{cases}$$

Note that  $\text{data}(f'_i) \leq \text{data}(f'_i) + |D| \leq \text{data}(f'_i) + |\Pi|^4$ .

CLAIM B.5.1.  $f'_i \preceq f'_i \preceq f_i$  for all  $i \in \{1, 2\}$ .

PROOF. By definition of  $\preceq$ , for every  $d \in D$  there is some  $d' \in \text{data}(f'_i) \setminus R_i$  so that  $\bar{\chi}_i(d) = \bar{\chi}'_i(d')$ . Therefore,  $f'_i = f'_i[d'_1 \mapsto d'_1, d_1] \cdots [d'_n \mapsto d'_n, d_n]$  where  $\{d_1, \dots, d_n\} = D \setminus \text{data}(f'_i)$  and for every  $j \in [n]$ ,  $d'_j \notin R_i$ ,  $d_j \notin \text{data}(f'_i)$ , and  $\bar{\chi}'_i(d'_j) = \bar{\chi}_i(d_j)$ . In other words, we have that  $f'_i \preceq f'_i \preceq f_i$  for all  $i \in \{1, 2\}$ .  $\square$

Therefore,  $f'_i$  is a bounded extension of  $f'_i$  for  $i \in \{1, 2\}$ .

CLAIM B.5.2.  $f'_1 + f'_2 \preceq f_3$ .

PROOF. First, we have that  $f'_1, f'_2 \in \mathfrak{F}$  because otherwise  $f_1$  or  $f_2$  would not be in  $\mathfrak{F}$ .

Since  $f'_i \preceq f_i$ , by Lemma ?? it follows that

$$\langle\langle f'_i \rangle\rangle = \langle\langle f_i \rangle\rangle, \quad \langle\langle f'_i \rangle\rangle = \langle\langle f_i \rangle\rangle. \quad (\dagger)$$

This means that, since  $f'_i$  is the result of adding data values  $d$  with profile  $\bar{\chi}_i(d)$ , if conditions (a), (b), (c) hold for  $f_1 + f_2$  they must also hold for  $f'_1 + f'_2$ . Therefore, all the preconditions to apply  $f'_1 + f'_2$  hold. Let  $f'_3 = f'_1 + f'_2$ .

Note that all the profiles in question share the same set of rigid values,  $R''_3 = R''_1 = R''_2 = R_1 = R_2 = R_3$  by definition of  $\preceq$  and  $+$ . Every data value  $d \in \text{data}(f'_3)$  is so that  $\bar{\chi}'_3(d) = \bar{\chi}_3(d)$ , and with (??) this means that  $\bar{\chi}_3(d) = \bar{\chi}'_3(d)$ , by definition of  $f'_1 + f'_2$ . For every other data value  $d \in \text{data}(f_3) \setminus \text{data}(f'_3)$ , it must be that  $d \notin R_3$  and there must be some  $\hat{d} \in D$  so that  $\bar{\chi}_3(\hat{d}) = \bar{\chi}_3(d)$  by definition of  $D$ . Then, if  $\text{data}(f_3) \setminus \text{data}(f'_3) = \{\hat{d}_1, \dots, \hat{d}_n\}$  we have that

$$f_3 = f'_3[\hat{d}_1 \mapsto \hat{d}_1, d_1] \cdots [\hat{d}_n \mapsto \hat{d}_n, d_n],$$

and hence,  $f'_3 \preceq f_3$ .  $\square$

By Claims ?? and ??, and since  $\text{data}(f'_i) \leq \text{data}(f'_i) + |\Pi|^4$ , the lemma follows.  $\square$

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

PROOF OF LEMMA 5.4. We are given  $f_1, f_2, f_5$  and  $(a, \bar{d})$ . We show that there must be some  $f_3$  and  $f_4$  so that  $f_1 \succeq f_3 \in (a, \bar{d})f_4$  with  $f_4 \succeq f_5$ . Figure 6 contains a graphical representation of the profiles we work with in the proof.

We first define  $f_4$ . For every internal description  $\bar{\pi}$  so that  $\bar{\pi}_1 \cup \bar{\pi}_1 \neq \emptyset$  and  $\bar{\chi}_1(\bar{\pi}) \setminus R_1 \neq \emptyset$ , let  $d_\pi \in D$  be a data value of  $\bar{\chi}_1(\bar{\pi}) \setminus R_1$ . Let  $D$  be the set of all such data values. Note that  $|D| \leq |\Pi|^4$ .

$$f_4 : \begin{cases} R_4 = R_5 \\ \text{for every } d \in D \cup R_1, \bar{\chi}_4(d) = \bar{\chi}_2(d) \\ \text{for every } d \notin D \cup R_1, \bar{\chi}_4(d) = \bar{\chi}_5(d) \end{cases}$$

Note that  $|\text{data}(f_4)| \leq |\text{data}(f_5)| + |\Pi|^4 + |R_1|$ . Since  $f_5 \preceq f_2$ , we have that  $f_5 \preceq f_4 \preceq f_2$ .

CLAIM B.5.3.  $f_5 \preceq f_4 \preceq f_2$ .

PROOF. • We first show that  $f_5 \preceq f_4$ .

First notice that  $R_5 = R_2$  since  $f_5 \preceq f_2$ , and that  $R_4 = R_5$  by definition of  $R_4$ . We then have that for every  $d \in R_5 = R_4$ ,  $\bar{\chi}_4(d) = \bar{\chi}_2(d) = \bar{\chi}_5(d)$ .

We show that for every  $d \in \text{data}(f_5)$ ,  $\bar{\chi}_4(d) = \bar{\chi}_5(d)$ . If  $d \in \text{data}(f_5) \setminus (D \cup R_1)$ , we have that  $\bar{\chi}_4(d) = \bar{\chi}_5(d)$ . If on the other hand  $d \in D \cup R_1$ , we have  $\bar{\chi}_4(d) = \bar{\chi}_2(d)$ . Then  $d \in \text{data}(f_2)$ , and since  $f_5 \preceq f_2$  we have that either  $d \notin \text{data}(f_5)$  or  $\bar{\chi}_2(d) = \bar{\chi}_5(d)$ . Hence, for every  $d \in \text{data}(f_5)$ ,  $\bar{\chi}_4(d) = \bar{\chi}_5(d)$ .

We finish the proof of  $f_4 \preceq f_5$  by showing that for every  $d \in \text{data}(f_4) \setminus R_4$  there is some  $d' \in \text{data}(f_5) \setminus R_4$  so that  $\bar{\chi}_4(d) = \bar{\chi}_5(d')$ . If  $d \notin D \cup R_1$ , then of course we can take  $d' = d$  and  $\bar{\chi}_4(d) = \bar{\chi}_5(d)$ . If  $d \in D \cup R_1$ , we have that  $\bar{\chi}_4(d) = \bar{\chi}_2(d)$  by definition. Since  $f_5 \preceq f_2$ , there must be some  $d' \in \text{data}(f_5) \setminus R_5$  so that  $\bar{\chi}_5(d') = \bar{\chi}_2(d) = \bar{\chi}_5(d)$ , and we are done. This finishes the proof that  $f_5 \preceq f_4$ .

• We now show that  $f_4 \preceq f_2$ .

Again, notice that  $R_4 = R_2 = R_5$ , and since for every  $d \in D$ ,  $\bar{\chi}_4(d)$  is either  $\bar{\chi}_2(d)$  or  $\bar{\chi}_5(d)$ , it follows that for every  $d \in R_4 = R_5$ ,  $\bar{\chi}_4(d) = \bar{\chi}_2(d)$ .

We show that for every  $d \in \text{data}(f_4)$ ,  $\bar{\chi}_4(d) = \bar{\chi}_2(d)$ . If  $d \in D \cup R_1$ , then  $\bar{\chi}_4(d) = \bar{\chi}_2(d)$  by definition. If  $d \in R_4$ , then  $\bar{\chi}_4(d) = \bar{\chi}_2(d)$  as shown before. If  $d \notin D \cup R_1 \cup R_4$ ,  $\bar{\chi}_4(d) = \bar{\chi}_5(d)$ . Notice that  $\bar{\chi}_5(d) \neq (\emptyset, \emptyset, \emptyset, \emptyset)$ , since otherwise  $d \notin \text{data}(f_4)$ . Then,  $d \in \text{data}(f_5)$  and since  $f_5 \preceq f_2$ ,  $\bar{\chi}_5(d) = \bar{\chi}_2(d) = \bar{\chi}_4(d)$ .

We finish the proof of  $f_4 \preceq f_2$  by showing that for every  $d \in \text{data}(f_2) \setminus R_2$  there is some  $d' \in \text{data}(f_4) \setminus R_2$  so that  $\bar{\chi}_2(d) = \bar{\chi}_4(d')$ . If  $d \in D \cup R_1$ , then we can just simply take  $d' = d$  and  $\bar{\chi}_2(d) = \bar{\chi}_4(d)$ . If  $d \notin D \cup R_1$ , then  $\bar{\chi}_4(d) = \bar{\chi}_5(d)$ , notice that  $\bar{\chi}_5(d) \neq (\emptyset, \emptyset, \emptyset, \emptyset)$  because otherwise  $d \notin \text{data}(f_4)$  (since  $d \notin R_4$ ). Hence,  $d \in \text{data}(f_5)$  and since  $f_5 \preceq f_2$ , we have that  $\bar{\chi}_2(d) = \bar{\chi}_5(d)$ , hence  $\bar{\chi}_2(d) = \bar{\chi}_4(d)$ . We just showed that  $f_4 \preceq f_2$ .  $\square$

We define  $f_3$  as follows

$$f_3 : \begin{cases} R_3 = R_1 \\ \text{for every } d \in R_1, \bar{\chi}_3(d) = \bar{\chi}_1(d) \\ \text{for every } d \in \text{data}(f_4), \bar{\chi}_3(d) = \bar{\chi}_1(d) \\ \text{for every } d \in \bar{d}, \bar{\chi}_3(d) = \bar{\chi}_1(d) \\ \text{for every strict external value } d \text{ of } f_1, \bar{\chi}_3(d) = \bar{\chi}_1(d). \end{cases}$$

CLAIM B.5.4.  $f_3 \preceq f_1$ .

PROOF. By definition we have that  $R_1 = R_3$ , and that  $\bar{\chi}_3(d) = \bar{\chi}_1(d)$  for all  $d \in R_1$ . Also by definition, for every  $d \in \text{data}(f_3)$  we have  $\bar{\chi}_3(d) = \bar{\chi}_1(d)$ .

We must show that for every  $d \in \text{data}(f_1) \setminus (R_1 \cup \text{data}(f_3))$  there is some  $d' \in \text{data}(f_3) \setminus R_1$  so that  $\bar{\chi}_1(d) = \bar{\chi}_3(d')$ . Take any such  $d$ . Note that  $d$  must be necessarily an internal data value, since  $f_3$  contains any strict external data value of  $f_1$ . Hence, let  $\bar{\pi} = \bar{\chi}_1(d)$ , where  $\bar{\pi}_1^1 \cup \bar{\pi}_1^1 \neq \emptyset$ . By definition of  $D$ , there must be some  $d_{\bar{\pi}} \in D$  so that  $\bar{\chi}_1(d_{\bar{\pi}}) = \bar{\pi}$  and  $d_{\bar{\pi}} \notin R_1$ . Since  $\bar{\pi}$  is internal and  $f_1 \in (a, \bar{d})f_2$ ,  $d_{\bar{\pi}} \in \bar{d} \cup \text{data}(f_2)$ . If  $d_{\bar{\pi}} \in \bar{d}$ , we have, by definition of  $f_3$ , that  $\bar{\chi}_3(d_{\bar{\pi}}) = \bar{\chi}_1(d_{\bar{\pi}}) = \bar{\pi} = \bar{\chi}_1(d)$ . Hence, there is such  $d' \in \text{data}(f_3) \setminus R_1$  so that  $\bar{\chi}_1(d) = \bar{\chi}_3(d')$ . Otherwise, suppose  $d_{\bar{\pi}} \in \text{data}(f_2)$ , which means, since  $\bar{\pi}_1^1 \cup \bar{\pi}_1^1 \neq \emptyset$ , that  $\bar{\chi}_2(d_{\bar{\pi}}) \neq (\emptyset, \emptyset, \emptyset, \emptyset)$ . By definition of  $f_4$ ,  $\bar{\chi}_4(d_{\bar{\pi}}) = \bar{\chi}_2(d_{\bar{\pi}}) \neq (\emptyset, \emptyset, \emptyset, \emptyset)$ , thus  $d_{\bar{\pi}} \in \text{data}(f_4)$ . Then, by definition of  $f_3$  we have that  $\bar{\chi}_3(d_{\bar{\pi}}) = \bar{\chi}_1(d_{\bar{\pi}}) = \bar{\pi} = \bar{\chi}_1(d)$ . Hence, in this case we also have that there is such  $d' \in \text{data}(f_3) \setminus R_1$  so that  $\bar{\chi}_1(d) = \bar{\chi}_3(d')$ .  $\square$

CLAIM B.5.5.  $f_3 \in (a, \bar{d})f_4$ .

PROOF. Since  $f_4 \preceq f_2$  and  $f_1 \in (a, \bar{d})f_2$ , we have that

- $\xi(f_4) \in \Gamma$  since  $\xi(f_2) \in \Gamma$ , by Lemma ??, and
- $\bar{\chi}_4^E = \bar{\chi}_4^E = \emptyset$  since  $\bar{\chi}_2^E = \bar{\chi}_2^E = \emptyset$ .

Therefore, conditions (a) and (b) hold for  $f_4$ . We must show condition (c), that is,  $\bar{\chi}_3^1 = \bar{\chi}_3^1 = \{(d, \alpha, \beta\gamma, i) \in \mathbb{D} \times \Pi \mid \exists \alpha'. (d, \alpha', \gamma, i) \in \bar{\chi}_4^1 \cup \bar{\chi}_4^1, \alpha, \beta \in \sigma_a\} \cup \bigcup_{i \in [k]} (\{d(i)\} \times (\sigma_a \setminus \{\epsilon\}) \times \sigma_a \times \{i\})$

Take any data value of  $d \in \bar{d} \cup \text{data}(f_4)$ . We have that  $\bar{\chi}_3^1(d) = \bar{\chi}_1^1(d)$  by definition. We also have that  $\bar{\chi}_1^1(d) = \bar{\chi}_1^1(d) = \{(\alpha, \beta\gamma, i) \in \Pi \mid \exists \alpha'. (d, \alpha', \gamma, i) \in \bar{\chi}_2^1 \cup \bar{\chi}_2^1, \alpha, \beta \in \sigma_a\} \cup \bigcup_{d=\bar{d}(i), i \in [k]} ((\sigma_a \setminus \{\epsilon\}) \times \sigma_a \times \{i\})$ . Since  $f_2$  and  $f_4$  coincide in  $d$ , we also have that the set above is equal to  $\{(\alpha, \beta\gamma, i) \in \Pi \mid \exists \alpha'. (d, \alpha', \gamma, i) \in \bar{\chi}_4^1 \cup \bar{\chi}_4^1, \alpha, \beta \in \sigma_a\} \cup \bigcup_{d=\bar{d}(i), i \in [k]} ((\sigma_a \setminus \{\epsilon\}) \times \sigma_a \times \{i\})$ . This, together with the fact that all internal values of  $f_3$  are in  $\text{data}(f_4) \cup \bar{d}$ , implies that  $\bar{\chi}_3^1 = \bar{\chi}_3^1 = \{(d, \alpha, \beta\gamma, i) \in \mathbb{D} \times \Pi \mid \exists \alpha'. (d, \alpha', \gamma, i) \in \bar{\chi}_4^1 \cup \bar{\chi}_4^1, \alpha, \beta \in \sigma_a\} \cup \bigcup_{i \in [k]} (\{d(i)\} \times (\sigma_a \setminus \{\epsilon\}) \times \sigma_a \times \{i\})$  as desired. Hence, condition (c) holds, and we have that  $f_3 \in (a, \bar{d})f_4$ .  $\square$

This concludes the proof.  $\square$

## Appendix to Section 5.4

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

PROOF OF LEMMA 5.5. Suppose we have  $f = f_1 + \dots + f_n$ . Let us first show that for every  $(\alpha, \beta, i) \in \Pi$  there can be at most one data value  $d \in \mathbb{D}$  so that there is some  $j \in [n]$  with  $\bar{\chi}_j^E(\alpha, \beta, i) = \{d\}$ . By means of contradiction, if there were two distinct data values  $d, d'$  then there would be two  $f_j, f_k$  with  $j \neq k$  so that  $\bar{\chi}_j^E(\alpha, \beta, i) = \{d\}$  and  $\bar{\chi}_k^E(\alpha, \beta, i) = \{d'\}$ . Suppose without any loss of generality that  $j < k$ . Then, by definition of  $+$ , we must have that—since  $d \in \bar{\chi}_j^E(\alpha, \beta, i)$  and  $f_j + \dots + f_k$  is defined— $d \in \bar{\chi}_k^E(\alpha, \beta, i)$ . This is in contradiction with the fact that  $\bar{\chi}_k^E(\alpha, \beta, i) = \{d'\}$ . The same happens by symmetry with  $\bar{\chi}^E$ .

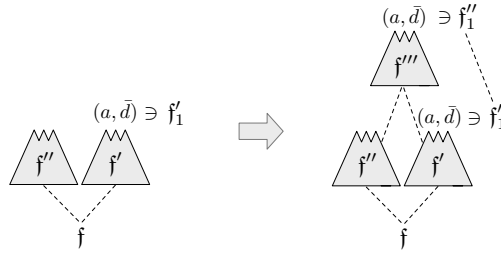
Let us define

$$\begin{aligned} R_{\bar{\chi}_1, \dots, \bar{\chi}_n}^l &\stackrel{\text{def}}{=} \{d \in \mathbb{D} \mid \bar{\chi}_j^E(\alpha, \beta, i) = \{d\} \\ &\quad \text{for some } (\alpha, \beta, i) \in \Pi \text{ and } j \in [n]\}, \\ R_{\bar{\chi}_1, \dots, \bar{\chi}_n}^r &\stackrel{\text{def}}{=} \{d \in \mathbb{D} \mid \bar{\chi}_j^E(\alpha, \beta, i) = \{d\} \\ &\quad \text{for some } (\alpha, \beta, i) \in \Pi \text{ and } j \in [n]\}, \\ R_{\bar{\chi}_1, \dots, \bar{\chi}_n} &\stackrel{\text{def}}{=} R_{\bar{\chi}_1, \dots, \bar{\chi}_n}^l \cup R_{\bar{\chi}_1, \dots, \bar{\chi}_n}^r. \end{aligned}$$

By the discussion before, it follows that  $|R_{\bar{\chi}_1, \dots, \bar{\chi}_n}| \leq 2|\Pi|$ . Consider the profiles  $f', f'_1, \dots, f'_n$  to be as  $f, f_1, \dots, f_n$  but with  $R_{\bar{\chi}_1, \dots, \bar{\chi}_n}$  as the set of rigid values. By construction of  $R_{\bar{\chi}_1, \dots, \bar{\chi}_n}$  we have that  $f' = f'_1 + \dots + f'_n$ . In other words we have the following.

CLAIM B.5.6. *For every  $(R, \bar{\chi}) = (R, \bar{\chi}_1) + \dots + (R, \bar{\chi}_n)$  we have that*

$$(R_{\bar{\chi}_1, \dots, \bar{\chi}_n}, \bar{\chi}) = (R_{\bar{\chi}_1, \dots, \bar{\chi}_n}, \bar{\chi}_1) + \dots + (R_{\bar{\chi}_1, \dots, \bar{\chi}_n}, \bar{\chi}_n).$$



**Figure 9: Statement of Lemma ??.**

Using the above statement we can now prove the lemma. Let  $t$  be a derivation tree for a root profile  $(R, f)$ . Let  $t'$  be the derivation tree that results from replacing in  $t$  the profile labels  $(R, \bar{\chi}_1), \dots, (R, \bar{\chi}_n)$  of any maximal sequence of siblings with the labels

$$(R_{\bar{\chi}_1, \dots, \bar{\chi}_n}, \bar{\chi}_1), \dots, (R_{\bar{\chi}_1, \dots, \bar{\chi}_n}, \bar{\chi}_n).$$

By the claim above,  $t'$  is a derivation tree for  $(R_{\bar{\chi}}, f)$ . (Note that this is true because in the rooting operation there is no restriction on which should be the set of rigid values of the parent profile.) It is immediate that  $(R, f)$  is a root profile if and only if  $(R_{\bar{\chi}}, f)$  is a root profile. Thus, the lemma follows.  $\square$

## Appendix to Section 5.5

Lemma 5.9 is the consequence of the following properties.

**LEMMA B.6** (FIGURE ??). *For every  $f, f', f'', f'_1 \in \mathfrak{F}_b$  and  $(a, \bar{d}) \in \mathbb{A} \times \mathbb{D}^k$  so that  $f \preceq f''$ ,  $f \preceq f'$  and  $f'_1 \in (a, \bar{d})f'$ , there are  $f''', f''_1 \in \mathfrak{F}_b$  so that  $f'' \preceq f'''$ ,  $f' \preceq f'''$ ,  $f'_1 \in (a, \bar{d})f'''$  and  $f'_1 \preceq f''_1$ .*

**PROOF OF LEMMA ??.** Without any loss of generality, let us assume that  $f \preceq f'$  and  $f \preceq f''$  so that  $data(f''') \cap data(f') = data(f'') \cap data(f'_1) = data(f)$ . We define  $f'''$  so that  $R''' = R$  and

$$\bar{\chi}'''(d) = \begin{cases} \bar{\chi}'(d) & \text{if } d \in data(f'), \\ \bar{\chi}''(d) & \text{else, if } d \in data(f''). \end{cases}$$

The following statement follows straight from this definition.

**CLAIM B.6.1.**  $f'' \preceq f'''$ ,  $f' \preceq f'''$ .

**PROOF.** By definition of  $f \preceq f''$  and  $f \preceq f'$  for every data value  $d \in data(f') \cup data(f'')$  there is a data value  $f(d) \in data(f)$  so that  $\bar{\chi}(d) = \bar{\chi}'(d) = \bar{\chi}''(d)$  and  $d \in R$  iff  $f(d) \in R$ . Then,

- $f''' = f'[f(d_1) \mapsto f(d_1), d_1] \cdots [f(d_n) \mapsto f(d_n), d_n]$  for  $\{d_1, \dots, d_n\} = data(f'') \setminus data(f')$ , and
- $f''' = f''[f(d_1) \mapsto f(d_1), d_1] \cdots [f(d_n) \mapsto f(d_n), d_n]$  for  $\{d_1, \dots, d_n\} = data(f') \setminus data(f'')$ .

Thus,  $f'' \preceq f'''$ ,  $f' \preceq f'''$ .  $\square$

Remember that by definition of  $f' \preceq f'''$ , for every data value  $d \in data(f''')$  there must be some data value  $f(d) \in data(f')$  so that  $\bar{\chi}'''(d) = \bar{\chi}'(d)$ . We can further assume that for every  $d \in data(f''')$ ,  $d \in R'_1$  iff  $f(d) \in R'_1$  (i.e., we can simply define  $f(d) = d$  for all  $d \in R'_1 \cap data(f''')$ ). We define  $f''_1$  as  $R''_1 = R'_1$  and

$$\bar{\chi}''_1(d) = \begin{cases} \bar{\chi}'_1(d) & \text{if } d \in data(f'_1) \\ \bar{\chi}'_1(f(d)) & \text{else, if } d \in data(f'''). \end{cases}$$

**CLAIM B.6.2.**  $f''_1 \in (a, \bar{d})f'''$ .

**PROOF.** For every data value of  $d \in data(f'_1)$  we have that  $\bar{\chi}''_1(d)$  and  $\bar{\chi}'''(d)$  are equal to  $\bar{\chi}'_1(d)$  and  $\bar{\chi}'(d)$ , and it therefore verifies the conditions imposed by the rooting operation. Further, any other data value of  $data(f''_1)$  (namely any data value from  $data(f''') \setminus data(f'_1)$ ) behaves as some data value  $f(d)$  in  $f'_1$ . That is  $\bar{\chi}''_1(d)$  and  $\bar{\chi}'''(d)$  are equal to  $\bar{\chi}'_1(f(d))$  and  $\bar{\chi}'(f(d))$ . It then follows that the data value must verify the conditions imposed by the rooting operation. Finally, by Lemma ??,  $\xi(f'_1) \in \Gamma$  and  $f'_1 \preceq f''_1$ , we have that  $\xi(f''_1) \in \Gamma$ . Hence,  $f''_1 \in (a, \bar{d})f'''$ .  $\square$

By definition of  $f''_1$ , it also follows that  $f'_1 \preceq f''_1$ .

**CLAIM B.6.3.**  $f'_1 \preceq f''_1$ .

**PROOF.** It is immediate from the definition that

$$f''_1 = f'_1[f(\hat{d}_1) \mapsto f(\hat{d}_1), \hat{d}_1, \dots, f(\hat{d}_s) \mapsto f(\hat{d}_s), \hat{d}_s],$$

where  $\{\hat{d}_1, \dots, \hat{d}_s\} = data(f''') \setminus data(f'_1)$ . By definition of  $f$ ,  $f(\hat{d}_1), \dots, f(\hat{d}_s) \in data(f') \setminus R'_1$ , and hence  $f'_1 \preceq f''_1$ .  $\square$

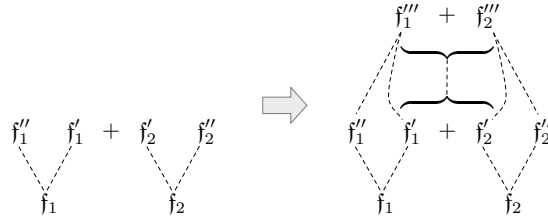


Figure 10: Statement of Lemma ??.

This concludes the proof.  $\square$

LEMMA B.7.  $\mathbf{R}_{up}(\uparrow\downarrow\mathcal{D}_b) \subseteq \downarrow\mathcal{D}_b$ .

PROOF OF LEMMA ???. Let  $f_1' \in \mathbf{R}_{up}(\uparrow\downarrow\mathcal{D}_b)$ . This means that there is some  $f'' \in \mathcal{D}_b$  and  $f$  so that  $f \lesssim f''$ ,  $f \lesssim f'$ , and  $f_1' \in (a, \bar{d})f'$  for some  $(a, \bar{d}) \in \mathbb{A} \times \mathbb{D}^k$ . We can then apply Lemma ??, obtaining that there is some  $f'''$  and  $f_1''$  so that  $f'' \lesssim f'''$ ,  $f' \lesssim f'''$ ,  $f_1' \lesssim f_1''$ , and  $f_1'' \in (a, \bar{d})f'''$ . Since  $\mathcal{D}_b$  is upward-closed by Remark 5.8,  $f''' \in \mathcal{D}_b$ , and therefore  $f_1'' \in \mathcal{D}_b$  as well since  $f_1'' \in (a, \bar{d})f'''$ . Thus, as  $f_1' \lesssim f_1''$ , we obtain  $f_1' \in \downarrow\mathcal{D}_b$ .  $\square$

LEMMA B.8 (FIGURE ??). For every  $i \in \{1, 2\}$  and  $f_i, f_i'' \in \mathfrak{F}_b$  so that  $f_i \lesssim f_i'$ ,  $f_i \lesssim f_i''$  and  $f_1' + f_2'$  is defined, there are  $f_1''', f_2''' \in \mathfrak{F}_b$  so that  $f_1' + f_2' \lesssim f_1''' + f_2'''$  and  $f_i' \lesssim f_i'''$ ,  $f_i'' \lesssim f_i'''$  for every  $i \in \{1, 2\}$ .

PROOF OF LEMMA ???. Without any loss of generality, assume that  $f_i \leq f_i'$ ,  $f_i \leq f_i''$  and  $data(f_i'') \cap data(f_i') = data(f_i)$  for every  $i \in \{1, 2\}$ . Remember that by definition of  $f_i \leq f_i''$ , for every data value  $d \in data(f_i'')$  there must be some data value  $f(d) \in data(f_i)$  so that  $\bar{\chi}_i(d) = \bar{\chi}_i'(d) = \bar{\chi}_i''(d)$ . We then define  $f_i'''$  for every  $i \in \{1, 2\}$  so that  $R_i''' = R_i$  and

$$\bar{\chi}_i'''(d) = \begin{cases} \bar{\chi}_i'(d) & \text{if } d \in data(f_i'), \\ \bar{\chi}_i''(d) & \text{else, if } d \in data(f_i''), \\ \bar{\chi}_i'(f(d)) & \text{else, if } d \in data(f_{3-i}'), \\ (\emptyset, \emptyset, \emptyset, \emptyset) & \text{otherwise.} \end{cases}$$

In the definition above, notice that for every  $d \in data(f_{3-i}')$  we have that  $f(d) \in data(f_{3-i})$  and hence  $f(d) \in data(f_{3-i}')$  which, by definition of  $f_1' + f_2'$ , means that  $f(d) \in data(f_i')$ .

The following claim follows straight from the definition just given.

CLAIM B.8.1.  $f_i' \leq f_i'''$ ,  $f_i'' \leq f_i'''$  for all  $i \in \{1, 2\}$ .

PROOF. Since  $f_i \leq f_i''$  and  $f_i \leq f_i'$ , there is some  $\hat{f}$  so that for every  $d \in data(f_i')$ ,  $\hat{f}(d) \in data(f_i)$  so that  $\bar{\chi}_i''(d) = \bar{\chi}_i(\hat{f}(d)) = \bar{\chi}_i''(\hat{f}(d))$  and  $d \in R_i$  iff  $\hat{f}(d) \in R_i$ .

For every  $i \in \{1, 2\}$ , it is immediate from the definition of  $f_i'''$  that

$$f_i''' = f_i''[f(\hat{d}_1) \mapsto f(\hat{d}_1), \hat{d}_1] \cdots [f(\hat{d}_n) \mapsto f(\hat{d}_n), \hat{d}_n] [f(d_1) \mapsto f(d_1), d_1] \cdots [f(d_m) \mapsto f(d_m), d_m]$$

where  $\{d_1, \dots, d_m\} = data(f_{3-i}')$  and  $\{\hat{d}_1, \dots, \hat{d}_n\} = data(f_i'') \setminus data(f_i')$ .

Since  $f_i \leq f_i''$  and  $f_i \leq f_i'$ , there is some  $\hat{f}$  so that for every  $d \in data(f_i'')$ ,  $\hat{f}(d) \in data(f_i)$  so that  $\bar{\chi}_i''(d) = \bar{\chi}_i(\hat{f}(d)) = \bar{\chi}_i''(\hat{f}(d)) = \bar{\chi}_i''(\hat{f}(d))$  and  $d \in R_i$  iff  $\hat{f}(d) \in R_i$ . For every  $i \in \{1, 2\}$ , it is immediate from the definition of  $f_i'''$  that

$$f_i''' = f_i''[\hat{f}(\hat{d}_1) \mapsto \hat{f}(\hat{d}_1), d_1] \cdots [\hat{f}(\hat{d}_n) \mapsto \hat{f}(\hat{d}_n), d_n] [f(d_1) \mapsto f(d_1), d_1] \cdots [f(d_m) \mapsto f(d_m), d_m]$$

where each  $d_i$  and  $\hat{d}_i$  is so that  $\{d_1, \dots, d_m\} = data(f_{3-i}')$  and  $\{\hat{d}_1, \dots, \hat{d}_n\} = data(f_i'')$ .  $\square$

We are then left with the following easy claim.

CLAIM B.8.2.  $f_1' + f_2' \leq f_1''' + f_2'''$ .

PROOF. For every data value  $d \in data(f_1') \cup data(f_2')$  it is easy to see that the conditions of  $+$  apply for  $f_1'''$  and  $f_2'''$  since they have the same description for  $d$ . For any other data value  $d \in data(f_1'') \cup data(f_2'')$  we have that  $f_1'''$  and  $f_2'''$  behave just as  $f_1'$  and  $f_2'$  for the data value  $f(d)$ . Therefore, the conditions of  $+$  hold, and  $f_1''' + f_2'''$  is well defined. Moreover,  $f_1' + f_2' \leq f_1''' + f_2'''$  since, by definition of  $f_1'''$ ,  $f_2'''$  we have

$$f_1''' + f_2''' = (f_1' + f_2')[f(d_1) \mapsto f(d_1), d_1] \cdots [f(d_n) \mapsto f(d_n), d_n]$$

for  $\{d_1, \dots, d_n\} = data(f_1''' + f_2''') \setminus data(f_1' + f_2')$ .  $\square$

This concludes the proof of the lemma.  $\square$

LEMMA B.9.  $\mathbf{R}_+(\uparrow\downarrow\mathcal{D}_b) \subseteq \downarrow\mathcal{D}_b$ .



PROOF OF LEMMA ???. Let  $f'_3 \in \mathbf{R}_+(\uparrow\downarrow\mathfrak{D}_b)$ . Then, there must be some  $f'_1, f'_2 \in \mathfrak{F}_b$  so that  $f'_1 + f'_2 = f'_3$  so that there exist  $f_1, f_2, f''_1, f''_2$  where  $f''_i \in \mathfrak{D}_b$ ,  $f_i \lesssim f'_i$ ,  $f_i \lesssim f''_i$  for all  $i \in \{1, 2\}$ . Therefore, by Lemma ??, there must be some  $f''_1, f''_2 \in \mathfrak{F}_b$  so that  $f'_3 = f'_1 + f'_2 \lesssim f''_1 + f''_2$  and  $f''_i \lesssim f'_i$ ,  $f''_i \lesssim f'_i$  for all  $i \in \{1, 2\}$ . Since  $f''_1, f''_2 \in \mathfrak{D}_b$ , and since  $\mathfrak{D}_b$  is upward-closed by Remark 5.8, we have that  $f''_1, f''_2 \in \mathfrak{D}_b$  and hence  $f''_1 + f''_2 \in \mathfrak{D}_b$ . Hence, since  $f'_3 \lesssim f''_1 + f''_2$  it follows that  $f'_2 \in \downarrow\mathfrak{D}_b$ .  $\square$

PROOF OF LEMMA 5.9. Immediate from Lemmas ?? and ??.  $\square$

## Appendix to Section 5.6

Thanks to the property  $\mathbf{R}(\uparrow\downarrow\mathfrak{D}_b) \subseteq \downarrow\mathfrak{D}_b$  one can show that  $\mathbf{C}_{k_0}^\sim = \text{MIN}(\downarrow\mathfrak{D}_b)$  and that  $\mathbf{C}_{i+1}$  can be computed from  $\mathbf{C}_i$  in 2EXPSpace in  $|\mathbf{C}_i|$  and  $|\mathcal{P}|$ . It is possible to test from the set  $\text{MIN}(\downarrow\mathfrak{D}_b)$  whether there is a root derivable profile in  $\mathfrak{D}_b$ , hence obtaining that the derivation problem is decidable.

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

PROOF OF LEMMA 5.11. By Lemma 5.9 we know that  $\mathbf{R}(\uparrow\downarrow\mathfrak{D}_b) \subseteq \downarrow\mathfrak{D}_b$ . Therefore, since  $\mathbf{C}_0 \subseteq \text{MIN}(\downarrow\mathfrak{D}_b)$ , we have that for every  $i$ ,  $\mathbf{C}_i \subseteq \text{MIN}(\downarrow\mathfrak{D}_b)$ . We then need to show that  $\text{MIN}(\downarrow\mathfrak{D}_b) \subseteq \mathbf{C}_{k_0}^\sim$ . Suppose, by means of contradiction, that there is some  $f \in \mathfrak{F}_b$  so that  $f \in \text{MIN}(\downarrow\mathfrak{D}_b)$  but  $f \notin \mathbf{C}_{k_0}^\sim$ . Then, since  $\{f_\emptyset\} = \mathbf{C}_0 \subseteq \mathbf{C}_{k_0}$  there must be one such  $f$  so that  $f \lesssim f'$  for some  $f' \in \mathfrak{D}_b$  where  $f' = f_1 + f_2$  for  $f_1, f_2 \in \mathbf{C}_{k_0}^\sim$ , or  $f' \in (a, \bar{d})f''$  for  $f'' \in \mathbf{C}_{k_0}^\sim$ ,  $(a, \bar{d}) \in \mathbb{A} \times \mathbb{D}^k$ .

- If  $f' = f_1 + f_2$  where  $f_1, f_2 \in \mathbf{C}_{k_0}^\sim$ , then  $f' \in \mathbf{R}_+(\uparrow\downarrow\mathbf{C}_{k_0})$ , and hence  $f \in \text{MIN}(\mathbf{R}_+(\uparrow\downarrow\mathbf{C}_{k_0})) \subseteq \mathbf{C}_{k_0+1}^\sim = \mathbf{C}_{k_0}^\sim$ , which is an absurd.
- If  $f' \in (a, \bar{d})f''$  where  $f'' \in \mathbf{C}_{k_0}^\sim$ , then  $f' \in \mathbf{R}_{up}^{(a, \bar{d})}(\uparrow\downarrow\mathbf{C}_{k_0})$ , and hence  $f \in \text{MIN}(\mathbf{R}_{up}^{(a, \bar{d})}) \subseteq \mathbf{C}_{k_0+1}^\sim = \mathbf{C}_{k_0}^\sim$ , which is also an absurd.

Therefore,  $\text{MIN}(\downarrow\mathfrak{D}_b) \subseteq \mathbf{C}_{k_0}^\sim$  and thus  $\mathbf{C}_{k_0}^\sim = \text{MIN}(\downarrow\mathfrak{D}_b)$ .  $\square$

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

LEMMA B.10.  $k_0$  is bounded by a doubly exponential function on  $|\mathcal{P}|$ .

PROOF OF LEMMA ??. Remember that  $\mathfrak{F}_b$  is the set of profiles that have less than  $2|\Pi|$  rigid values. Then,

$$(2^\Pi \times 2^\Pi \times 2^\Pi \times 2^\Pi)^{2|\Pi|}$$

represent all the possible profiles of the rigid values. We must also remember which profiles have either 0, or 1 or more flexible values, and this information can be represented with an element of

$$2^{2^\Pi \times 2^\Pi \times 2^\Pi \times 2^\Pi}.$$

Therefore, there are doubly exponentially many elements in  $\text{MIN}(\mathfrak{F}_b)/\sim$ , and hence  $k_0$  is bounded by a doubly exponential function in  $|\Pi|$ . Since  $|\Pi|$  is polynomial in  $|\mathcal{P}|$ , the statement follows.  $\square$

LEMMA B.11. Given  $G \subseteq \mathfrak{F}_b$  and  $f \in \mathfrak{F}_b$ , testing  $f \in \text{MIN}(\downarrow\mathbf{R}(\uparrow\downarrow G))$  is computable in 2EXPSpace in  $|G|$  and  $|\mathcal{P}|$ .

PROOF OF LEMMA ??. First we check that  $f \in \text{MIN}(\downarrow\mathfrak{F}_b)$ .

To check  $f \in \text{MIN}(\downarrow\mathbf{R}_+(\uparrow\downarrow G))$ , we verify if  $f \lesssim f_1 + f_2$  in 2EXPSpace, where  $f_1, f_2$  is a pair of the (doubly exponentially many) bounded extensions of elements of  $\text{MIN}(\downarrow G)$ . By Lemma 5.3, there are such  $f_1, f_2$  if, and only if,  $f \in \text{MIN}(\downarrow\mathbf{R}_+(\uparrow\downarrow G))$ .

Finally, to check  $f \in \text{MIN}(\downarrow\mathbf{R}_{up}(\uparrow\downarrow G))$ , we verify that  $f \lesssim f_3 \in \text{MIN}((a, \bar{d})f_4)$  in 2EXPSpace, where  $f_4$  is a bounded extension of a profile of  $\text{MIN}(\downarrow G)$ . If the condition holds of course  $f \in \text{MIN}(\downarrow\mathbf{R}_{up}(\uparrow\downarrow G))$ . On the other hand, by Lemma 5.4 cum Remark 5.6, if  $f \in (a, \bar{d})f_2$  with  $f_2 \lesssim f_5$ , then  $f \lesssim f_3 \in \text{MIN}((a, \bar{d})f_4)$  for some bounded extension  $f_4$  of  $f_5$ , and hence the condition holds.  $\square$

By the lemma above, we immediately obtain the following.

LEMMA B.12.  $\mathbf{C}_{i+1}$  can be computed from  $\mathbf{C}_i$  in 2EXPSpace in  $|\mathbf{C}_i|$  and  $|\mathcal{P}|$ .

PROOF OF LEMMA ??. For each  $f \in \text{MIN}(\mathfrak{F}_b)/\sim$  we can check, in 2EXPSpace, whether  $f \in \text{MIN}(\downarrow\mathbf{R}(\uparrow\downarrow\mathbf{C}_i))$  thanks to Lemma ??. We can therefore compute  $\text{MIN}(\downarrow\mathbf{R}(\uparrow\downarrow\mathbf{C}_i))/\sim$  in 2EXPSpace, and thus we can also compute  $\mathbf{C}_{i+1}$  in 2EXPSpace.  $\square$

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

PROOF OF PROPOSITION 5.12. We can compute all

$$\mathbf{C}_0, \dots, \mathbf{C}_{k_0}$$

in 2EXPSpace by Lemmas ?? and ??. Since  $\mathbf{C}_{k_0}^\sim = \text{MIN}(\downarrow\mathfrak{D}_b)$  by Lemma 5.11, it follows that  $\text{MIN}(\downarrow\mathfrak{D}_b)$  is hence computable in 2EXPSpace.

In order to test if there is a derivable profile in  $\mathfrak{D}_b$  we choose some  $f \in \text{MIN}(\{f \in \mathfrak{F}_b \mid \xi(f) \in \Gamma\})$  with no external values (i.e., so that  $\overleftarrow{\chi}^E = \overrightarrow{\chi}^E = \emptyset$ ),  $f' \in \text{MIN}(\downarrow\mathfrak{D}_b)$  and  $(a, \bar{d}) \in \mathbb{A}_{root} \times \mathbb{D}^k$ , and test, in 2EXPSpace (Lemma ??), whether  $f \in \text{MIN}(\downarrow\mathbf{R}_{up}^{(a, \bar{d})}(\uparrow\downarrow\{f'\}))$ .

CLAIM B.12.1. The following statements are equivalent.

1. There is some  $f \in \text{MIN}(\{f \in \mathfrak{F} \mid \xi(f) \in \Gamma\})$  with no external values so that

- $f \in \text{MIN}(\downarrow \mathbf{R}_{up}^{(a,\bar{d})}(\uparrow \downarrow \{f'\})),$
- $(a, \bar{d}) \in \mathbb{A}_{root} \times \mathbb{D}^k$  and
- $f' \in \text{MIN}(\downarrow \mathfrak{D}_b).$

2. There is a derivable root profile in  $\mathfrak{D}_b$ .

PROOF. [?? $\Rightarrow$ ??] Assume first that condition ?? holds. We change the names for the profiles to make the explanation clear. Suppose that there is some  $f'_1 \in \text{MIN}(\downarrow \mathbf{R}_{up}^{(a,\bar{d})}(\uparrow \downarrow \{f\}))$  with no external data values, for some  $(a, \bar{d}) \in \mathbb{A}_{root} \times \mathbb{D}^k$ ,  $f'_1 \in \text{MIN}(\{f \in \mathfrak{F}_b \mid \xi(f) \in \Gamma\})$ , and  $f \in \text{MIN}(\downarrow \mathfrak{D}_b)$ . This means that there must be some  $f'' \in \mathfrak{D}_b$  so that  $f \lesssim f''$  and some  $f'$  so that  $f \lesssim f'$  and  $f'_1 \in (a, \bar{d})f'$ . We can then apply Lemma ?? obtaining that there is some  $f'''$  and  $f''_1$  so that  $f'' \lesssim f'''$ ,  $f' \lesssim f'''$ ,  $f'_1 \lesssim f''_1$ , and  $f''_1 \in (a, \bar{d})f'''$ . Since  $\mathfrak{D}_b$  is upward closed (Remark 5.8),  $f''' \in \mathfrak{D}_b$  and hence  $f''_1 \in \mathfrak{D}_b$ . By definition of  $\lesssim$ , it follows that, since  $f'_1$  has no external data values,  $f''_1$  has no external data values either. Thus, there is a derivable root profile, namely  $f''_1$ .

[?? $\Leftarrow$ ??] Suppose, on the other hand, that condition ?? holds. If there is a derivable root profile in  $\mathfrak{D}_b$ , then there are  $f, f' \in \mathfrak{D}_b$  and  $(a, \bar{d}) \in \mathbb{A}_{root} \times \mathbb{D}^k$  so that  $\overleftarrow{\chi}^e = \overrightarrow{\chi}^e = \emptyset$  and  $f \in (a, \bar{d})f'$ . Let  $f'_1 \in \text{MIN}(\downarrow \{f'\})/\sim$ , that is,  $f'_1$  is a minimal element corresponding to  $f'$ . Also, let  $f_1 \in \text{MIN}(\downarrow \{f\})/\sim$ , that is,  $f_1$  is a minimal element corresponding to  $f$ . Of course, it follows that  $f_1 \in \text{MIN}(\downarrow \mathbf{R}_{up}^{(a,\bar{d})}(\uparrow \downarrow \{f_1\}))$ . Note that  $f_1 \in \text{MIN}(\{f \in \mathfrak{F}_b \mid \xi(f) \in \Gamma\})$ , and that  $f_1$  has no external values since  $f$  has no external values. Hence condition ?? holds.  $\square$

Since there is a derivable root profile in  $\mathfrak{D}_b$  if and only if there is a derivable root profile in  $\mathfrak{D}$  by Lemma 5.7, the proposition follows.  $\square$

## C. APPENDIX TO SECTION 6

### Appendix to Section 6.1

LEMMA 6.1 (DIRECT NORMAL FORM). *There exists an exponential time translation that for every node expression  $\varphi \in \text{XPath}(*\leftarrow, \downarrow_*, \rightarrow^*, =)$  returns an equivalent node expression  $\psi$  in direct normal form.*

PROOF OF LEMMA 6.1. The idea is that every data test expression  $\langle \alpha @_i = \beta @_j \rangle$  is translated into a big disjunction of expressions in direct normal form, where loops in the tree navigation of  $\alpha, \beta$  are factored as node expressions, as done in (‡).

For any finite alphabet  $\mathbb{B}$ , we define **tree order morphisms** between forests over a powerset alphabet  $2^{\mathbb{B}}$ . Given two forests  $\bar{t}, \bar{t}'$  over  $2^{\mathbb{B}}$ , a tree order morphism from  $\bar{t}$  to  $\bar{t}'$  is a function  $f$  from the nodes of  $\bar{t}$  to the nodes of  $\bar{t}'$  so that

- for every node  $x$  of  $\bar{t}$ , the label of  $x$  in  $\bar{t}$  is a subset of the label of  $f(x)$  in  $\bar{t}'$ , and
- for every two nodes  $x, y$  of  $\bar{t}$ , if  $(x, y)$  is in the reflexive-transitive closure of the next-sibling relation (resp. of the child relation) in  $\bar{t}$ , then  $(f(x), f(y))$  is also in the reflexive-transitive closure of the next-sibling relation (resp. of the child relation) in  $\bar{t}'$ .

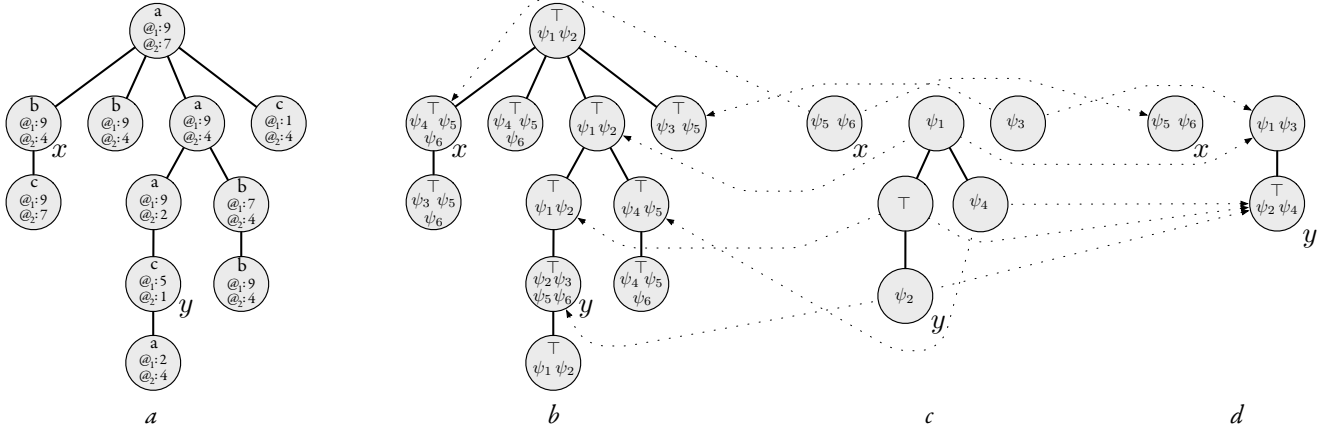
We say that  $\bar{t}', (x', y')$  is a **contraction** of  $\bar{t}, (x, y)$  if there is a tree order morphism  $f$  from  $\bar{t}$  to  $\bar{t}'$  so that  $f(x) = x'$  and  $f(y) = y'$ . We also say that  $\bar{t}, (x, y)$  is an **expansion** of  $\bar{t}', (x', y')$ . Notice that the function need not be surjective, and hence the fact that one forest is a contraction of another does not have any implication on the sizes of the forests: a forest may have less or more nodes than its contraction. For example, both the tree of Figure ??-b and the forest of Figure ??-d are contractions of the forest of Figure ??-c.

Let  $\alpha$  be a path expression, and let  $ne(\alpha)$  be the set of all node expressions of  $\alpha$ . Consider  $\hat{\mathbf{t}}^\alpha$  the tree over the alphabet  $2^{ne(\alpha)}$ , where every node is labeled by those  $\psi \in ne(\alpha)$  that are true at the node. For example, for the multi-attribute data tree  $\mathbf{t}$  defined in Figure ??-a and  $\alpha$  as defined in Figure ??-e, the tree  $\hat{\mathbf{t}}^\alpha$  is the one of Figure ??-b.

Let  $\bar{t}$  be any forest over  $2^{ne(\alpha)}$ , and let  $x, y$  be two nodes of  $\bar{t}$ , so that  $x$  is a root. For such a forest and nodes, one can build a direct path expression  $\mu_{\bar{t}, x, y}^\alpha$  that tests if there is a path between  $x$  and  $y$  satisfying the relative order of nodes satisfying the labels as in the tree. For example, for the forest  $\bar{t}$  of Figure ??-c and the nodes  $x, y$  depicted, the corresponding expression  $\hat{\mu}_{\bar{t}, x, y}^\alpha$  would be the one appearing in Figure ??-e.

CLAIM C.0.1. *A direct path expression  $\mu_{\bar{t}, x, y}^\alpha$  can be built from  $\bar{t}$  in polynomial time, such that for every multi-attribute data tree  $\mathbf{t}$  we have that  $\mathbf{t}, (x', y') \models \mu_{\bar{t}, x, y}^\alpha$  if and only if  $\hat{\mathbf{t}}^\alpha, (x', y')$  is a contraction of  $\bar{t}, (x, y)$ .*

PROOF. Given a forest  $\bar{t}$  and nodes  $x, y$  so that  $x$  is a root, one can build  $\mu_{\bar{t}, x, y}^\alpha$  iteratively. We first build a path expression that starts in  $x$  and ends in  $y$  and checks all the labels between  $x$  and  $y$  in the unique path between  $x$  and  $y$  that corresponds to a direct navigation. For example if  $\bar{t}, x, y$  are as in Figure ??-a, we build a path as in Figure ??-b that tests all the node expressions of the nodes in



$$e \left\{ \begin{array}{l} \alpha = \leftarrow[\psi_6] \rightarrow^* [\psi_3] \leftarrow [\psi_7] \rightarrow^* [\psi_1] \downarrow_* [\psi_4] \leftarrow [\top] \downarrow_* [\psi_2] \\ \mu_{\bar{t},x,y}^\alpha = [\psi_6 \wedge \psi_7] \rightarrow^* [\psi_1 \wedge \langle \rightarrow^* [\psi_3] \rangle] \downarrow_* [\top \wedge \langle \rightarrow^* [\psi_4] \rangle] \downarrow_* [\psi_2] \\ \text{where } \psi_1 = a \quad \psi_3 = c \quad \psi_5 = \neg a \\ \psi_2 = \langle \downarrow_* [a] \rangle \quad \psi_4 = b \quad \psi_6 = \neg(\leftarrow[a]) \end{array} \right.$$

Figure 11: Example for the proof of Lemma 6.1.

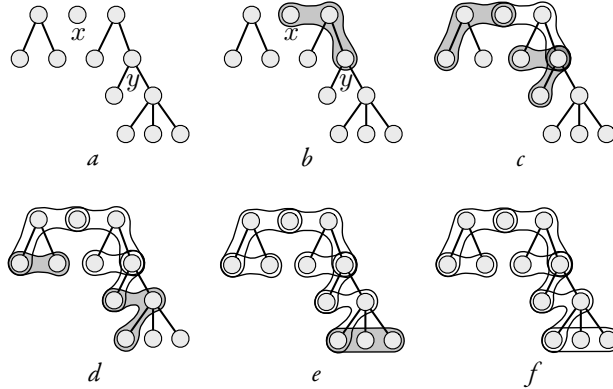


Figure 12: Idea of construction of  $\hat{\mu}_{\bar{t},x,y}^\alpha$ .

the path. We then build a more complex path expression nesting expressions that test the existence of paths as depicted in Figure ??-c. We iterate until we have covered the whole tree  $\bar{t}$ , at each iteration, we add tests for paths that are at a deeper nesting degree in the resulting direct path expression. For example, the paths of Figure ??-c are at depth 1, the paths of Figure ??-d at depth 2, and the path of Figure ??-e at depth 3. Note that if we apply this construction to the forest of Figure ??-c and the path  $\alpha$  of Figure ??-e, we obtain the expression  $\mu_{\bar{t},x,y}^\alpha$  of Figure ??-e.  $\square$

Given a multi-attribute data tree  $\mathbf{t}$  so that  $\mathbf{t}, (x, y) \models \alpha$ , we call a *witness forest* of  $\alpha$  for  $\mathbf{t}, (x, y)$  to a forest  $\bar{t}$  over the alphabet  $2^{ne(\alpha)}$  together with some nodes  $x, y$ , so that  $\bar{t}$  contains only the nodes involved in the satisfaction of  $\alpha$ , and they are labeled with the node expressions that they must verify. For example, given the multi-attribute data tree  $\mathbf{t}$  of Figure ??-a, and  $\alpha$  as defined in Figure ??-e, a possible witness forest of  $\alpha$  for  $\mathbf{t}, (x, y)$  is depicted in Figure ??-c. Of course, there may be several witness forests of  $\alpha$  for a given  $\mathbf{t}, (x, y)$ . Let  $\|\alpha\|$  be the size of  $\alpha$ , computed as the number of axes in  $\alpha$ , irrespective of the size of node expressions. For example the path expression  $\alpha$  defined in Figure ??-e is so that  $\|\alpha\| = 7$ . Notice that the number of nodes of any witness forest of  $\alpha$  for  $\mathbf{t}, (x, y)$  is bounded by  $\mathbf{p}(\|\alpha\|)$  for some polynomial  $\mathbf{p}(\cdot)$ .

Given a path expression  $\alpha$ , let  $\hat{\alpha}$  be the path expression over the alphabet  $2^{ne(\alpha)}$ , where every node expressions  $\psi$  is replaced by  $\bigvee_{S \subseteq ne(\alpha), \psi \in S} S$ , that is, it is treated as a disjunction of tests for labels from the alphabet. Notice that any witness forest  $\bar{t}, (x', y')$  of  $\alpha$  for  $\mathbf{t}, (x, y)$  is an expansion of  $\hat{\mathbf{t}}^\alpha, (x, y)$ , and we have  $\bar{t}, (x', y') \models \hat{\alpha}$ . We then have the following.

CLAIM C.0.2.  $\mathbf{t}, (x, y) \models \alpha$  if, and only if,  $\bar{t}, (x', y') \models \hat{\alpha}$ , for some expansion  $\bar{t}, (x', y')$  of  $\hat{\mathbf{t}}^\alpha, (x, y)$ .  
 For every contraction  $\bar{t}, (x, y)$  of  $\bar{t}', (x', y')$ , if  $\bar{t}', (x', y') \models \hat{\alpha}$  then  $\bar{t}, (x, y) \models \hat{\alpha}$ . This observation, together with the previous claim and the fact that every witness forest  $\bar{t}, (x', y')$  of  $\alpha$  for  $\mathbf{t}, (x, y)$  is bounded by  $\mathbf{p}(\|\alpha\|)$  yields the following.

CLAIM C.0.3. *There is an expansion  $\bar{t}, (x', y')$  of  $\hat{\mathbf{t}}^\alpha, (x, y)$  of size at most  $\mathbf{p}(\|\alpha\|)$  so that  $\bar{t}, (x', y') \models \alpha$  if, and only if,  $\mathbf{t}, (x, y) \models \alpha$ .*

For every forest  $\bar{t}$  over  $2^{ne(\alpha)}$  of size  $\leq \mathbf{p}(\|\alpha\|)$  and nodes  $x, y$  so that  $\bar{t}, (x, y) \models \hat{\alpha}$ , we say that  $\mu_{\bar{t}, x, y}^\alpha$  is a *direct normal form linearization* of  $\alpha$ . It follows that  $\mathbf{t}, (x, y) \models \mu_{\bar{t}, x, y}^\alpha$  if, and only if, a contraction of  $\bar{t}, (x', y')$  is a witness forest of  $\alpha$  for  $\mathbf{t}, (x, y)$ . Thus, the following claim follows.

CLAIM C.0.4.  *$\mathbf{t}, (x, y) \models \alpha$  if, and only if,  $\mathbf{t}, (x, y) \models \mu_{\bar{t}, x, y}^\alpha$  for some direct normal form linearization  $\mu_{\bar{t}, x, y}^\alpha$  of  $\alpha$ .*

PROOF. If  $\mathbf{t}, (x, y) \models \alpha$ , then by Claim ?? there is a forest  $\bar{t}, (x', y')$  over  $2^{ne(\alpha)}$  (the witness forest) of size  $\leq \mathbf{p}(\|\alpha\|)$  so that  $\bar{t}, (x', y') \models \hat{\alpha}$ . Hence,  $\mu_{\bar{t}, x, y}^\alpha$  is a direct normal form linearization of  $\alpha$ . Since  $\bar{t}, (x', y')$  is an expansion of  $\hat{\mathbf{t}}^\alpha, (x, y)$  and  $\bar{t}, (x', y') \models \hat{\alpha}$ , then  $\mathbf{t}, (x, y) \models \mu_{\bar{t}, x, y}^\alpha$ .

Suppose now that  $\mathbf{t}, (x, y) \models \mu_{\bar{t}, x, y}^\alpha$  for some  $\bar{t}, (x', y')$  over  $2^{ne(\alpha)}$  of size  $\leq \mathbf{p}(\|\alpha\|)$  so that  $\bar{t}, (x', y') \models \hat{\alpha}$ . Then,  $\bar{t}, (x', y')$  must be an expansion of  $\hat{\mathbf{t}}^\alpha, (x, y)$  by Claim ??. Since  $\bar{t}, (x', y') \models \hat{\alpha}$ , we then have that  $\mathbf{t}, (x, y) \models \alpha$  by Claim ??.  $\square$

Note that there are exponentially many trees in  $\|\alpha\|$  over  $ne(\alpha)$  of size  $\leq \mathbf{p}(\|\alpha\|)$ . Hence, all the direct normal form linearizations of a given path expression  $\alpha$  are computable in exponential time  $\|\alpha\|$ . Hence, given  $\varphi \in \text{XPath}(*\leftarrow, \downarrow_*, \rightarrow^*, =)$ , consider  $\varphi'$  as the result of replacing every appearance of  $\langle \alpha @_i = \beta @_j \rangle$  (resp.  $\neq$ ) by a disjunction of  $\langle \mu_{\bar{t}, x, y}^\alpha @_i = \mu_{\bar{t}', x', y'}^\beta @_j \rangle$  (resp.  $\neq$ ) for every direct normal form linearization  $\mu_{\bar{t}, x, y}^\alpha$  of  $\alpha$  and  $\mu_{\bar{t}', x', y'}^\beta$  of  $\beta$ . Since these are all the possible linearizations, the translated formula is satisfied in the same multi-attribute data trees and nodes as the original formula. Further,  $\varphi'$  is in direct normal form. Since every replacement is exponential in  $\|\alpha\|$ , and since  $\|\alpha\|$  does not depend on the size of sub-node expressions, it follows that  $|\varphi'|$  in time exponential in  $|\varphi|$ . Thus, the lemma follows.  $\square$

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

PROOF OF LEMMA 6.2. Given a formula  $\eta$  we define the alphabet  $\mathbb{A}_\varphi$  of the translation  $\varphi$  as all the locally consistent sets of subformulas of  $\eta$ . That is, the sets  $S$  such that for every subformula  $\psi$  of  $\eta$ : (1) if  $\psi = \neg\psi'$  then  $\{\psi', \neg\psi'\} \not\subseteq S$ ; (2) if  $\psi = \psi' \wedge \psi''$  then  $\psi \in S$  iff  $\{\psi', \psi''\} \subseteq S$ ; and (3) if  $\psi = \psi' \vee \psi''$  then  $\psi \in S$  iff  $\psi' \in S$  or  $\psi'' \in S$ .

Given a formula  $\psi$ ,  $tr(\psi)$  denotes the result of replacing

- every instance of a path expression  $o_1[\psi_1]o_2 \cdots o_n[\psi_n]$  in  $\psi$  (where  $o_1, \dots, o_n \in \{*\leftarrow, \downarrow_*, \rightarrow^*\}$ ) which does not appear nested inside another path expression, with the expression

$$o_1[\zeta_{\psi_1}]o_2 \cdots o_n[\zeta_{\psi_n}],$$

and

- every test for label  $a$  which does not appear inside a path expression by  $\zeta_a$ ,

where  $\zeta_\psi \stackrel{\text{def}}{=} \bigvee_{S \in \mathbb{A}_\varphi, \psi \in S} S$ .

To build the formula  $\varphi = \varphi_1 \wedge \varphi_2$  in normal form, we define  $\varphi_1 = \zeta_\eta$ , and we build  $\varphi_2$  as a conjunction of formulas

$$\neg(\downarrow_*[\zeta_\psi \wedge \neg tr(\psi)])$$

for all subformulas  $\psi$  of  $\eta$ . It is easy to see that this translation preserves satisfiability.  $\square$

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

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

PROOF OF COROLLARY 6.3. The blowup in the exponential translation comes only from the formulas  $\zeta_\psi$ . In fact,  $\varphi$  can be symbolically written in polynomial space just as we did, using a symbol  $\zeta_\psi$  instead of a big exponential disjunction. Remark that testing whether a label  $S \in \mathbb{A}_\varphi$  satisfies  $\zeta_\psi$  reduces to testing  $\psi \in S$ .  $\square$

## Appendix to Section 6.2

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

PROOF. If  $\psi \in \mathcal{B}(\mathbb{A})$ , then it is immediate that  $\mathfrak{f} \vdash \psi$  iff  $\xi(\mathfrak{f})((\psi, \epsilon, i), \circ) \neq 0$  for some  $i \in [\mathbf{k}]$ .

Suppose now that  $\psi = \langle \alpha \cdot \beta @_i \neq \gamma \cdot \delta @_j \rangle$  where  $\alpha$  is leftward,  $\epsilon$  or empty and  $\gamma$  is rightward,  $\epsilon$  or empty, and  $\beta, \delta$  are downward or empty. We show that whether  $\mathfrak{f} \vdash \psi$  depends only on  $\xi(\mathfrak{f})$ . We show

one of the possible cases in the definition of  $\vdash$ , while the remaining ones are analogous or simpler. Suppose, for example, that  $\alpha \neq \epsilon$ ,  $\alpha \neq \varepsilon$ ,  $\gamma = \epsilon$ . We then have that  $f \vdash \psi$  iff

- there is at least one data value reachable with  $\alpha\beta@_i$ , that is,  $\xi(f)(\bar{\alpha}, \bar{\circ}) \neq 0$ , with  $\bar{\alpha} = (P_{* \leftarrow}^{-1}(\alpha), P_{\downarrow * }^{-1}(\beta), i)$ , and
- there is at least one data value reachable with  $\gamma\delta@_j$ , that is,  $\xi(f)(\bar{\beta}, \circ^\perp) \neq 0$ , with  $\bar{\beta} = (\top, P_{\downarrow * }^{-1}(\delta), j)$  and either
  - there are two data values reachable with  $\alpha\beta@_i$ , that is,  $\xi(f)(\bar{\alpha}, \bar{\circ}) = 2+$ , or
  - there are two data values reachable with  $\gamma\delta@_j$ , that is,  $\xi(f)(\bar{\beta}, \circ^\perp) = 2+$ , or
  - there is one data value reachable with  $\alpha\beta@_i$  and only one data value reachable with  $\gamma\delta@_j$ , but they are different, that is,  $\xi(f)(\bar{\alpha}, \bar{\circ}) = 1$ ,  $\xi(f)(\bar{\beta}, \circ^\perp) = 1$ , and  $\xi(f)(\bar{\alpha}, \bar{\circ}, \bar{\beta}, \circ^\perp) = 0$ .

These are indeed the necessary and sufficient conditions for the existence of two data values  $d \neq d'$  so that  $d \in \chi_a(\bar{\alpha})$  and  $d' \in \chi_{a'}(\bar{\beta})$ , where  $a = \bar{\circ}$ ,  $a' = \circ^\perp$ . Hence, these are the necessary and sufficient conditions to verify  $f \vdash \psi$ .

The remaining cases are analogous or easier.  $\square$

LEMMA 6.5.  $f \vdash \bigwedge_{a \in \mathbb{A}} (a \Rightarrow \gamma_\varphi(a))$  if, and only if,  $\xi(f) \in \Gamma_\phi$ .

PROOF OF LEMMA 6.5. This is a direct consequence of Lemma 6.4 and the definition of  $\Gamma_\phi$ .  $\square$

### Abstractions.

We now define  $abs(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}, \bar{\mathbf{t}}_r)$  for any forests  $\bar{\mathbf{t}}_l, \bar{\mathbf{t}}, \bar{\mathbf{t}}_r$ . Let  $\mathbf{t}_l = \mathbf{a}_l \otimes \mathbf{d}_l$ ,  $\mathbf{t} = \mathbf{a} \otimes \mathbf{d}$ ,  $\mathbf{t}_r = \mathbf{a}_r \otimes \mathbf{d}_r$  be the multi-attributes data trees  $(a, \bar{d})\bar{\mathbf{t}}_l$ ,  $(a, \bar{d})\bar{\mathbf{t}}$ , and  $(a, \bar{d})\bar{\mathbf{t}}_r$  respectively, for some fixed arbitrary  $(a, \bar{d})$ . We then define  $abs(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}, \bar{\mathbf{t}}_r) = \bar{\chi}$  where

$$\begin{aligned} \bar{\chi}^1 &= \{(\mathbf{d}(y)(i), \alpha, \beta, i) \mid (1, y) \in \llbracket P_{\rightarrow * }(\alpha) \cdot P_{\downarrow * }(\beta) \rrbracket_{\mathbf{t}}, \\ &\quad \text{for } (\alpha, \beta, i) \in \Pi \text{ and } y \in pos(\mathbf{t})\} \\ \bar{\chi}^1 &= \{(\mathbf{d}(y)(i), \alpha, \beta, i) \mid (n, y) \in \llbracket P_{* \leftarrow }(\alpha) \cdot P_{\downarrow * }(\beta) \rrbracket_{\mathbf{t}}, \\ &\quad \text{for } n = |\bar{\mathbf{t}}|, (\alpha, \beta, i) \in \Pi \text{ and } y \in pos(\mathbf{t})\} \\ \bar{\chi}^E &= \{(\mathbf{d}_r(y)(i), \alpha, \beta, i) \mid (1, y) \in \llbracket P_{\rightarrow * }(\alpha) \cdot P_{\downarrow * }(\beta) \rrbracket_{\mathbf{t}_r}, \\ &\quad \text{for } (\alpha, \beta, i) \in \Pi \text{ and } y \in pos(\mathbf{t}_r)\} \\ \bar{\chi}^E &= \{(\mathbf{d}_l(y)(i), \alpha, \beta, i) \mid (n, y) \in \llbracket P_{* \leftarrow }(\alpha) \cdot P_{\downarrow * }(\beta) \rrbracket_{\mathbf{t}_l}, \\ &\quad \text{for } n = |\bar{\mathbf{t}}_l|, (\alpha, \beta, i) \in \Pi \text{ and } y \in pos(\mathbf{t}_l)\}. \end{aligned}$$

For any  $R \subseteq \mathbb{D}$  we define  $abs_R(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}, \bar{\mathbf{t}}_r) \stackrel{def}{=} (abs(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}, \bar{\mathbf{t}}_r), R)$ .

REMARK C.1. Note that  $abs_\emptyset(\epsilon, \epsilon, \epsilon) = f_\emptyset$ .

We show that  $abs$  is basically an algebra morphism between multi-attribute data forests with rooting and concatenation and forest profiles with profile rooting and profile concatenation. This is necessary to show our reduction from SAT-XPath( $*\leftarrow, \downarrow, \rightarrow, =$ ) into the derivation problem for forest profiles.

LEMMA C.2. Given  $R \subseteq \mathbb{D}$ , and multi-attribute data forests  $\bar{\mathbf{t}}, \bar{\mathbf{t}}_l, \bar{\mathbf{t}}_r, \bar{\mathbf{t}}_1, \bar{\mathbf{t}}_2$ ,

1. provided that  $abs_R(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}_1, \bar{\mathbf{t}}_2 \cdot \bar{\mathbf{t}}_r)$ ,  $abs_R(\bar{\mathbf{t}}_l \cdot \bar{\mathbf{t}}_1, \bar{\mathbf{t}}_2, \bar{\mathbf{t}}_r)$ ,  $abs_R(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}_1 \cdot \bar{\mathbf{t}}_2, \bar{\mathbf{t}}_r) \in \mathfrak{F}$ ,

$$abs_R(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}_1, \bar{\mathbf{t}}_2 \cdot \bar{\mathbf{t}}_r) + abs_R(\bar{\mathbf{t}}_l \cdot \bar{\mathbf{t}}_1, \bar{\mathbf{t}}_2, \bar{\mathbf{t}}_r) = abs_R(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}_1 \cdot \bar{\mathbf{t}}_2, \bar{\mathbf{t}}_r).$$

2. Given  $(a, \bar{d}) \in \mathbb{A} \times \mathbb{D}^k$ ,  $R' \subseteq \mathbb{D}$  so that  $abs_{R'}(\epsilon, \bar{\mathbf{t}}, \epsilon) \in \mathfrak{F}$ ,  $abs_R(\bar{\mathbf{t}}_l, (a, \bar{d})\bar{\mathbf{t}}, \bar{\mathbf{t}}_r) \in \mathfrak{F}$ , and  $\xi(abs_R(\bar{\mathbf{t}}_l, (a, \bar{d})\bar{\mathbf{t}}, \bar{\mathbf{t}}_r)) \in \Gamma_\phi$ , we have

$$abs_R(\bar{\mathbf{t}}_l, (a, \bar{d})\bar{\mathbf{t}}, \bar{\mathbf{t}}_r) \in (a, \bar{d}) abs_{R'}(\epsilon, \bar{\mathbf{t}}, \epsilon).$$

PROOF OF LEMMA ??.

???. Condition (a) is obvious because we have the same set  $R$  of rigid values in all profiles. Conditions (b) and (c) follow straight from the semantics of XPath and the definition of  $abs$ . By definition of  $abs$ , the set  $\bar{\chi}_1^E$  where  $f_1 = abs_R(\bar{\mathbf{t}}_l \cdot \bar{\mathbf{t}}_1, \bar{\mathbf{t}}_2, \bar{\mathbf{t}}_r)$  is composed by all the paths that can reach data values in  $\bar{\mathbf{t}}_2 \cdot \bar{\mathbf{t}}_r$ , which is precisely  $\bar{\chi}_2^1 \cup [f_1] \cdot \bar{\chi}_2^E$  where  $f_2 = abs_R(\bar{\mathbf{t}}_l \cdot \bar{\mathbf{t}}_1, \bar{\mathbf{t}}_2, \bar{\mathbf{t}}_r)$ . In the same way,  $\bar{\chi}_2^E$  is composed by all the paths that can reach data values in  $\bar{\mathbf{t}}_l \cdot \bar{\mathbf{t}}_1$ , which is precisely  $\bar{\chi}_1^1 \cup \{f_1\} \cdot \bar{\chi}_1^E$ . Finally, conditions (+2)–(+5) follow from the definition of  $abs$ .

???. Since  $\xi(abs_R(\bar{\mathbf{t}}_l, (a, \bar{d})\bar{\mathbf{t}}, \bar{\mathbf{t}}_r)) \in \Gamma_\phi$ , we only need to show that  $abs_R(\bar{\mathbf{t}}_l, (a, \bar{d})\bar{\mathbf{t}}, \bar{\mathbf{t}}_r)$  is so that the rooting conditions (b) and (c) hold. Condition (b) is immediate since by definition of  $abs$  we have that  $abs_{R'}(\epsilon, \bar{\mathbf{t}}, \epsilon)$  has empty external descriptions. Condition (c) holds by definition of  $abs$  and the semantics of XPath.  $\square$

For any derivation forest  $\bar{t}$  (as defined in Section 4.3) we associate a multi-attribute data forest  $\bar{\mathbf{t}}_t$  which is the result of removing all leaf nodes from  $t$  and projecting the tree onto  $\mathbb{A} \times \mathbb{D}^k$ .

LEMMA C.3. For every boolean combination  $\psi$  of non-recursive subformulas of  $\phi$ , we have that  $abs_R(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}, \bar{\mathbf{t}}_r) \vdash \psi$  if, and only if,  $\bar{\mathbf{t}}_l \cdot \bar{\mathbf{t}}_r, i_0 \models \psi$ , for  $i_0 = |\bar{\mathbf{t}}_l| + 1$  and  $R = data(\bar{\mathbf{t}}_l \cdot \bar{\mathbf{t}}_r)$ .

PROOF OF LEMMA ???. Let  $f = \text{abs}_R(\bar{\mathbf{t}}_l, \mathbf{t}, \bar{\mathbf{t}}_r)$ .

• Suppose first that  $\psi = b$  with  $b \in \mathbb{A}$ . Since  $b \in \mathcal{P}_\phi$ , we have that  $f \vdash b$  iff  $(b, \epsilon, i) \in \bar{\chi}^1(d)$  for some  $d, i$  iff  $(1, y) \in \llbracket \rightarrow^* [b] @_i \rrbracket_{(a, \bar{d}) \mathbf{t}}$  for some  $y$  by definition of  $\text{abs}$  iff  $(1, 1) \in \llbracket \rightarrow^* [b] @_i \rrbracket_{(a, \bar{d}) \mathbf{t}}$  iff  $\mathbf{t} \models b$  iff  $\bar{\mathbf{t}}_l \cdot \mathbf{t} \cdot \bar{\mathbf{t}}_r, i_0 \models b$ , since  $b \in \mathbb{A}$ .

• If  $\psi = \psi_1 \wedge \psi_2$  or  $\psi = \neg \psi'$ , we proceed by induction.

• Now suppose  $\psi = \langle \alpha \cdot \beta @_i = \gamma \cdot \delta @_j \rangle$ , where  $\alpha$  is leftward or empty,  $\gamma$  is rightward or empty, and  $\beta, \delta$  are downward or empty. Let  $\mathbf{t}' = \mathbf{a}' \otimes \mathbf{d}' = (a, \bar{d})(\bar{\mathbf{t}}_l \cdot \mathbf{t} \cdot \bar{\mathbf{t}}_r)$ .

[ $\Rightarrow$ ] If  $f \vdash \psi$  there must be some data value  $d \in \mathbb{D}$  that witnesses this fact, as required by the definition of  $\vdash$ . Suppose first that  $\alpha \neq \epsilon, \alpha \neq \varepsilon$ . Hence, there must be some  $d \in \mathbb{D}$  so that  $(P_{\leftarrow}^{-1}(\alpha), P_{\downarrow}^{-1}(\beta), i) \in (\llbracket f \rrbracket \cdot \bar{\chi}^E \cup \bar{\chi}^1)(d)$ .

1. If  $(P_{\leftarrow}^{-1}(\alpha), P_{\downarrow}^{-1}(\beta), i) \in \bar{\chi}^1(d)$ , by definition of  $\text{abs}$  we have that  $\bar{\chi}^1$  is the set of all  $(\mathbf{d}(y)(k), \alpha_1, \alpha_2, k)$  so that  $(1, y) \in \llbracket P_{\leftarrow}^{-1}(\alpha_1) P_{\downarrow}^{-1}(\alpha_2) @_k \rrbracket_{(a, \bar{d}) \mathbf{t}}$ . It then follows that there must be some  $y'$  so that  $(i_0, y') \in \llbracket \alpha \cdot \beta @_i \rrbracket_{(a, \bar{d})(\bar{\mathbf{t}}_l \cdot \mathbf{t} \cdot \bar{\mathbf{t}}_r)}$  with  $\mathbf{d}'(y')(i) = d$ .
2. If on the other hand  $(P_{\leftarrow}^{-1}(\alpha), P_{\downarrow}^{-1}(\beta), i) \in \llbracket f \rrbracket \cdot \bar{\chi}^E(d)$ , either  $(P_{\leftarrow}^{-1}(\alpha), P_{\downarrow}^{-1}(\beta), i) \in \bar{\chi}^E(d)$ , or there must be some  $\alpha'_1, \alpha'_2$  so that  $\alpha'_1 \in \llbracket f \rrbracket$  and  $(\alpha'_2, P_{\downarrow}^{-1}(\beta), i) \in \bar{\chi}^E(d)$  and  $\alpha'_1 \cdot \alpha'_2 = P_{\leftarrow}^{-1}(\alpha)$ . Suppose that the latter occurs (the former being only easier). Since by definition of  $\text{abs}$  we have that  $\bar{\chi}^E$  is the set of all  $(\mathbf{d}_l(y)(k), \alpha_1, \beta_1, k)$  so that  $(n, y) \in \llbracket P_{\leftarrow}^{-1}(\alpha_1) P_{\downarrow}^{-1}(\beta_1) @_k \rrbracket_{(a, \bar{d}) \bar{\mathbf{t}}_l}$  for  $n = |\bar{\mathbf{t}}_l| = i_0 - 1$  there is some  $y'$  so that  $(i_0 - 1, y') \in \llbracket \alpha_2 \cdot \beta \rrbracket_{(a, \bar{d})(\bar{\mathbf{t}}_l \cdot \mathbf{t} \cdot \bar{\mathbf{t}}_r)}$  with  $\mathbf{d}'(y')(i) = d$  and  $\alpha_2 = P_{\leftarrow}^{-1}(\alpha'_2)$ . Also, since  $\alpha'_1 \in \llbracket f \rrbracket$ , we have  $(i_0, i_0 - 1) \in \llbracket \alpha_1 \rrbracket_{(a, \bar{d})(\bar{\mathbf{t}}_l \cdot \mathbf{t} \cdot \bar{\mathbf{t}}_r)}$  where  $\alpha_1 = P_{\leftarrow}^{-1}(\alpha'_1)$ . Then, there is some  $y'$  so that  $(i_0, y') \in \llbracket \alpha \cdot \beta \rrbracket_{(a, \bar{d})(\bar{\mathbf{t}}_l \cdot \mathbf{t} \cdot \bar{\mathbf{t}}_r)}$  with  $\mathbf{d}'(y')(i) = d$ .

The case where  $\alpha = \epsilon$ , or  $\alpha = \varepsilon, \gamma = \epsilon$  is only easier (it is basically as in the item 1). In any case we obtain that there is some  $y'$  so that  $(i_0, y') \in \llbracket \alpha \cdot \beta \rrbracket_{(a, \bar{d})(\bar{\mathbf{t}}_l \cdot \mathbf{t} \cdot \bar{\mathbf{t}}_r)}$  with  $\mathbf{d}'(y')(i) = d$ .

By similar arguments, we have that there must be some  $z' \in \text{pos}((a, \bar{d})(\bar{\mathbf{t}}_l \cdot \mathbf{t} \cdot \bar{\mathbf{t}}_r))$ ,  $z' \neq \epsilon$  with  $(i_0, z') \in \llbracket \alpha \cdot \beta \rrbracket_{(a, \bar{d})(\bar{\mathbf{t}}_l \cdot \mathbf{t} \cdot \bar{\mathbf{t}}_r)}$  and  $\mathbf{d}'(z')(j) = d$ . Hence,  $\bar{\mathbf{t}}_l \cdot \mathbf{t} \cdot \bar{\mathbf{t}}_r, i_0 \models \langle \alpha \cdot \beta @_i = \gamma \cdot \delta @_j \rangle$ .

[ $\Leftarrow$ ] Suppose that  $\bar{\mathbf{t}}_l \cdot \mathbf{t} \cdot \bar{\mathbf{t}}_r, i_0 \models \langle \alpha \cdot \beta @_i = \gamma \cdot \delta @_j \rangle$ . In other words, there are  $y, z \in \text{pos}(\mathbf{t}')$ ,  $y \neq \epsilon, z \neq \epsilon$ , so that  $(i_0, y) \in \llbracket \alpha \cdot \beta \rrbracket_{\mathbf{t}'}$ ,  $(i_0, z) \in \llbracket \gamma \cdot \delta \rrbracket_{\mathbf{t}'}$ , and  $d = \mathbf{d}'(z)(i) = \mathbf{d}'(y)(j)$  for some  $d \in \mathbb{D}$ . Let  $\alpha' = \top$  if  $\alpha = \varepsilon$  or  $\alpha = \epsilon$  and  $\alpha' = P_{\leftarrow}^{-1}(\alpha)$  otherwise.

1. If  $y$  is inside the subtree  $\mathbf{t}$  of  $\mathbf{t}'$ , then  $(\alpha', P_{\downarrow}^{-1}(\beta), i) \in \bar{\chi}^1(d)$ , by definition of  $\text{abs}$ .
2. If  $y$  is inside the subforest  $\bar{\mathbf{t}}_l$  of  $\mathbf{t}'$ , then we have that either  $(\alpha', P_{\downarrow}^{-1}(\beta), j) \in \bar{\chi}^E(d)$ , or there must be some  $\alpha_1, \alpha_2$  so that  $\alpha_1 \cdot \alpha_2 = \alpha$ ,  $\alpha'_1 \in \llbracket f \rrbracket$  and  $(\alpha'_2, P_{\downarrow}^{-1}(\beta), j) \in \bar{\chi}^E(d)$ , where  $\alpha'_1 = P_{\leftarrow}^{-1}(\alpha_1), \alpha'_2 = P_{\leftarrow}^{-1}(\alpha_2)$ . If the latter case holds (the former being only easier), we have that  $(\alpha'_1 \cdot \alpha'_2, P_{\downarrow}^{-1}(\beta), j) \in (\llbracket f \rrbracket \cdot \bar{\chi}^E)(d)$ , which implies that  $(\alpha', P_{\downarrow}^{-1}(\beta), j) \in (\llbracket f \rrbracket \cdot \bar{\chi}^E)(d)$ .

By similar arguments, we obtain some  $(\gamma', P_{\downarrow}^{-1}(\delta), j) \in (\bar{\chi}^1 \cup \llbracket f \rrbracket \cdot \bar{\chi}^E)(d)$ , where  $\gamma' = \top$  if  $\gamma = \varepsilon$  or  $\gamma = \epsilon$  and  $\gamma' = P_{\leftarrow}^{-1}(\gamma)$  otherwise. Hence,  $f \vdash \langle \alpha \cdot \beta @_i = \gamma \cdot \delta @_j \rangle$ .

• The cases of  $\langle \alpha \cdot \beta @_i = \gamma \cdot \delta @_j \rangle$  where both  $\alpha$  and  $\gamma$  are rightwards or leftwards is similar. The same applies to the inequality formulas  $\langle \alpha \cdot \beta @_i \neq \gamma \cdot \delta @_j \rangle$ .  $\square$

The next lemma basically shows that every derivable root profile is the abstraction of some tree satisfying  $\phi$ .

LEMMA C.4. Let  $f_{\text{root}} \in \mathfrak{D}$  be a root profile, and let  $t$  be its derivation tree. Let  $\bar{t}$  be a maximal subforest of  $t$ , where  $\bar{\mathbf{t}}$  is the multi-attribute data forest associated with  $\bar{t}$ , and  $f_1, \dots, f_n$  are the profile labels of the roots of  $\bar{t}$ . Let  $R$  be the rigid set for  $f_1, \dots, f_n$ . Then,

- $\text{abs}_R(\epsilon, \bar{\mathbf{t}}, \epsilon) = f_1 + \dots + f_n$ , and
- $\bar{\mathbf{t}}, i \models \phi_2$  for every  $1 \leq i \leq |\bar{\mathbf{t}}|$ .

PROOF OF LEMMA ???. We will make use of the following easy fact.

CLAIM C.4.1. Given  $f_1, \dots, f_n \in \mathfrak{F}$ , so that  $f_+ = f_1 + \dots + f_n$  where  $\bar{\chi}_+^E = \bar{\chi}_+^E = \emptyset$ , then the sets  $\{\bar{\chi}_i^1, \bar{\chi}_i^1 \mid i \in [n]\}$  determine the sets  $\{\bar{\chi}_i^E, \bar{\chi}_i^E \mid i \in [n]\}$ .

PROOF. This is by conditions (b) and (c) of profile concatenation:

- $\bar{\chi}_i^1$  and  $\bar{\chi}_i^E$  determine  $\bar{\chi}_{i+1}^E$ . Hence,  $\{\bar{\chi}_i^E \mid i \in [n]\}$  are determined by  $\bar{\chi}_1^E$  and  $\{\bar{\chi}_i^1 \mid i \in [n]\}$ .
- $\bar{\chi}_{i+1}^1$  and  $\bar{\chi}_{i+1}^E$  determine  $\bar{\chi}_i^E$ . Hence,  $\{\bar{\chi}_i^E \mid i \in [n]\}$  are determined by  $\bar{\chi}_n^E$  and  $\{\bar{\chi}_i^1 \mid i \in [n]\}$ .

Once we fix  $\bar{\chi}_n^E = \bar{\chi}_1^E = \emptyset$ , then  $\{\bar{\chi}_i^1, \bar{\chi}_i^1 \mid i \in [n]\}$  determine  $\{\bar{\chi}_i^E, \bar{\chi}_i^E \mid i \in [n]\}$ .  $\square$

We now prove the main statement by generalized induction on the height of  $\bar{t}$ . Suppose that  $\bar{t}$  has height  $m \geq 1$ . Then there must be forests  $\bar{\mathbf{t}}_1, \dots, \bar{\mathbf{t}}_n$  so that

$$\bar{\mathbf{t}} = (a_1 \bar{d}_1) \bar{\mathbf{t}}_1 \cdots (a_n \bar{d}_n) \bar{\mathbf{t}}_n$$

and each  $\bar{\mathbf{t}}_i$  has height less than  $m$ . Let

$$\hat{f}_i = \text{abs}_R((a_j, \bar{d}_j) \bar{\mathbf{t}}_j)_{j < i}, (a_i, d_i) \bar{\mathbf{t}}_i, ((a_j, \bar{d}_j) \bar{\mathbf{t}}_j)_{j > i}$$

for every  $i$ . We then have the following.

CLAIM C.4.2. For every  $i \in [n]$ ,  $f_i = \hat{f}_i$ .

PROOF. For every  $\bar{t}_i \neq \epsilon$ , we have that  $f'_1 + \dots + f'_l = \text{abs}_{R_i}(\epsilon, \bar{t}_i, \epsilon)$  by inductive hypothesis, where  $f'_1, \dots, f'_l$  are the profile labels of the children of the root of the  $i$ th tree of  $\bar{t}$ . Then, since  $f_i \in (a_i, \bar{d}_i)(f'_1 + \dots + f'_l)$  by definition of derivation tree, we have that  $f_i \in (a_i, \bar{d}_i)\text{abs}_{R_i}(\epsilon, \bar{t}_i, \epsilon)$ . If  $\bar{t}_i = \epsilon$ , then, since  $\text{abs}_\emptyset(\epsilon, \epsilon, \epsilon) = f_\emptyset$  and  $f_i \in (a_i, \bar{d}_i)f_\emptyset$ , we also have that  $f_i \in (a_i, \bar{d}_i)\text{abs}_{R_i}(\epsilon, \bar{t}_i, \epsilon)$ . Hence, for all  $i \in [n]$ ,  $f_i \in (a_i, \bar{d}_i)\text{abs}_{R_i}(\epsilon, \bar{t}_i, \epsilon)$ .

Now, since

- all the internal descriptions of any forest profile inside  $(a_i, \bar{d}_i)\text{abs}_\emptyset(\epsilon, \bar{t}_i, \epsilon)$  is completely determined by  $(a_i, \bar{d}_i)$  and  $\text{abs}_\emptyset(\epsilon, \bar{t}_i, \epsilon)$ , and hence  $f_i$  and  $\hat{f}_i$  have the same internal descriptions,
- the external descriptions of  $f_1 + \dots + f_n$  are  $\emptyset$ ,
- $\text{abs}_R(\epsilon, (a_1, \bar{d}_1)\bar{t}_1 \dots (a_n, \bar{d}_n)\bar{t}_n, \epsilon) = \hat{f}_1 + \dots + \hat{f}_n$  by Lemma ??-??, and the external descriptions of  $\hat{f}_1 + \dots + \hat{f}_n$  are  $\emptyset$

we can apply Claim ??, concluding that  $f_i = \hat{f}_i$  for every  $i$ .  $\square$

Therefore, since  $\text{abs}_R(\epsilon, (a_1, \bar{d}_1)\bar{t}_1 \dots (a_n, \bar{d}_n)\bar{t}_n, \epsilon) = \hat{f}_1 + \dots + \hat{f}_n$  by Lemma ??-??, we have that

$$\text{abs}_R(\epsilon, (a_1, \bar{d}_1)\bar{t}_1 \dots (a_n, \bar{d}_n)\bar{t}_n, \epsilon) = f_1 + \dots + f_n.$$

On the other hand, since  $\xi(f_1), \dots, \xi(f_n) \in \Gamma_\phi$  by condition (a) of rooting, we have  $\hat{f}_i = f_i \vdash \bigwedge_{a \in \mathbb{A}_\phi} (a \Rightarrow \gamma_\phi(a))$  for every  $i$ . Hence, by Lemma ??,  $(a_1, \bar{d}_1)\bar{t}_1 \dots (a_n, \bar{d}_n)\bar{t}_n, i \models \bigwedge_{a \in \mathbb{A}} (a \Rightarrow \gamma_\phi(a))$  for every  $i$ . Also, by inductive hypothesis, for all  $\bar{t}_i \neq \epsilon$  we have that  $\bar{t}_i, j \models \phi_2$  for all  $1 \leq j \leq |\bar{t}_i|$ . Then, by definition of  $\phi_2$ , it follows that

$$(a_1, \bar{d}_1)\bar{t}_1 \dots (a_n, \bar{d}_n)\bar{t}_n, i \models \phi_2$$

for every  $i$ .  $\square$

We now show that if  $\phi$  is satisfiable, then there must be a derivable root profile.

LEMMA C.5. Let  $\mathbf{t}_\phi$  be a tree so that  $\mathbf{t}_\phi \models \phi$ . For every maximal subforest  $\bar{\mathbf{t}}_l \cdot \bar{\mathbf{t}}_r$  of  $\mathbf{t}_\phi$  so that  $\bar{\mathbf{t}} \neq \epsilon$ , there is a derivable forest profile  $f \in \mathfrak{D}$  so that  $f = \text{abs}_R(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}, \bar{\mathbf{t}}_r)$ , where  $R = \text{data}(\bar{\mathbf{t}}_l \cdot \bar{\mathbf{t}}_r)$ . Further, if  $\bar{\mathbf{t}}$  is a tree,  $f$  is so that  $f \in (a, \bar{d})f'$ , with  $f' \in \mathfrak{D}$ .

PROOF OF LEMMA ??. We proceed by induction on the size of  $\bar{\mathbf{t}}$ , defined as its number of nodes. Suppose first that  $\bar{\mathbf{t}}$  consists in one tree with just one node  $(a, \bar{d})$ . We show that  $\text{abs}_R(\bar{\mathbf{t}}_l, (a, \bar{d}), \bar{\mathbf{t}}_r) \in (a, \bar{d})f_\emptyset$ .

Since  $\bar{\mathbf{t}}_l \cdot (a, \bar{d}) \cdot \bar{\mathbf{t}}_r, j \models \phi_2$ , for  $j = |\bar{\mathbf{t}}_l| + 1$ , we have that  $\text{abs}_R(\bar{\mathbf{t}}_l, \mathbf{t}, \bar{\mathbf{t}}_r) \vdash \gamma_\phi(a)$  by Lemma ??. In other words, we have  $\xi(\text{abs}_R(\bar{\mathbf{t}}_l, \mathbf{t}, \bar{\mathbf{t}}_r)) \in \Gamma_\phi$ . Then, by Lemma ??-?? we have  $\text{abs}_R(\bar{\mathbf{t}}_l, (a, \bar{d}), \bar{\mathbf{t}}_r) \in (a, \bar{d})\text{abs}_\emptyset(\epsilon, \epsilon, \epsilon)$ . Since  $f_\emptyset = \text{abs}_\emptyset(\epsilon, \epsilon, \epsilon)$  (Remark ??), we have that  $(a, \bar{d})\text{abs}_\emptyset(\epsilon, \epsilon, \epsilon) = (a, \bar{d})f_\emptyset$ . Therefore,  $\text{abs}_R(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}, \bar{\mathbf{t}}_r) \in \mathfrak{D}$ .

Suppose now that  $\bar{\mathbf{t}}$  consists in one tree  $(a, \bar{d})\bar{\mathbf{t}}'$  where  $\bar{\mathbf{t}}' \neq \epsilon$ . By inductive hypothesis there is some  $f' \in \mathfrak{D}$  so that  $f' = \text{abs}_{R'}(\epsilon, \bar{\mathbf{t}}', \epsilon)$  for some  $R'$ . By Lemma ?? since  $\bar{\mathbf{t}}_l \cdot (a, \bar{d})\bar{\mathbf{t}}' \cdot \bar{\mathbf{t}}_r, j \models \phi_2$  for  $j = |\bar{\mathbf{t}}_l| + 1$ , we have that  $\xi(\text{abs}_R(\bar{\mathbf{t}}_l, (a, \bar{d})\bar{\mathbf{t}}', \bar{\mathbf{t}}_r)) \in \Gamma_\phi$ . We then have that

$$\text{abs}_R(\bar{\mathbf{t}}_l, (a, \bar{d})\bar{\mathbf{t}}', \bar{\mathbf{t}}_r) \in (a, \bar{d})f' = (a, \bar{d})\text{abs}_{R'}(\epsilon, \bar{\mathbf{t}}', \epsilon)$$

again by Lemma ??-??. Therefore,  $\text{abs}_R(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}, \bar{\mathbf{t}}_r) \in \mathfrak{D}$ .

Suppose now that  $\bar{\mathbf{t}}$  consists in more than one tree, let  $\bar{\mathbf{t}}_1$  and  $\bar{\mathbf{t}}_2$  be non-empty forests so that  $\bar{\mathbf{t}} = \bar{\mathbf{t}}_1 \cdot \bar{\mathbf{t}}_2$ . By applying twice the inductive hypothesis there must be  $f_1, f_2 \in \mathfrak{D}$  so that  $f_1 = \text{abs}_R(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}_1, \bar{\mathbf{t}}_2 \cdot \bar{\mathbf{t}}_r)$  and  $f_2 = \text{abs}_R(\bar{\mathbf{t}}_l \cdot \bar{\mathbf{t}}_1, \bar{\mathbf{t}}_2, \bar{\mathbf{t}}_r)$ . As  $\text{abs}_R(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}_1 \cdot \bar{\mathbf{t}}_2, \bar{\mathbf{t}}_r) = \text{abs}_R(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}_1, \bar{\mathbf{t}}_2 \cdot \bar{\mathbf{t}}_r) + \text{abs}_R(\bar{\mathbf{t}}_l \cdot \bar{\mathbf{t}}_1, \bar{\mathbf{t}}_2, \bar{\mathbf{t}}_r)$  by Lemma ??-??. It follows that  $\text{abs}_R(\bar{\mathbf{t}}_l, \bar{\mathbf{t}}_1 \cdot \bar{\mathbf{t}}_2, \bar{\mathbf{t}}_r) = f_1 + f_2 \in \mathfrak{D}$  since  $f_1, f_2 \in \mathfrak{D}$ .  $\square$

COROLLARY 6.6. There is a derivable root forest profile if, and only if,  $\phi$  is satisfiable.

PROOF OF COROLLARY 6.6. [ $\Rightarrow$ ] Let  $f$  be a derivable root forest profile. Then, there must be a derivation tree  $t$  for  $f = (R, \bar{\chi})$ . Let  $\mathbf{t}$  be the multi-attribute data tree associated to  $t$ . By Lemma ??, we have that  $\text{abs}_R(\epsilon, \mathbf{t}, \epsilon) = f$  and that  $\mathbf{t} \models \phi_2$ . Since  $f$  is a root profile, the root of  $\mathbf{t}$  must have a label  $(a, \bar{d})$  with  $a \in \mathbb{A}_{\text{root}}$ , and hence  $\mathbf{t} \models \phi_1$  as well. Therefore,  $\phi$  is satisfiable.

[ $\Leftarrow$ ] Let  $\mathbf{t}$  be so that  $\mathbf{t} \models \phi$ . Then by Lemma ?? we have that  $f = \text{abs}_{\text{data}(\mathbf{t})}(\epsilon, \mathbf{t}, \epsilon) \in \mathfrak{D}$ . Further, since  $\mathbf{t}$  is a tree, we have that in fact  $f = (a, \bar{d})f'$  for some  $f' \in \mathfrak{D}$ . Further, since  $\mathbf{t} \models \phi_1$ ,  $a \in \mathbb{A}_{\text{root}}$ . Hence,  $f$  is a derivable root profile.  $\square$