



A Decidable Two-Way Logic on Data Words

Diego Figueira

► To cite this version:

Diego Figueira. A Decidable Two-Way Logic on Data Words. Annual IEEE Symposium on Logic in Computer Science (LICS), Jun 2011, Toronto, United States. 10.1109/LICS.2011.18 . hal-01803452

HAL Id: hal-01803452

<https://hal.science/hal-01803452>

Submitted on 30 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.

A decidable two-way logic on data words

Diego Figueira
University of Warsaw

Abstract—We study the satisfiability problem for a logic on data words. A data word is a finite word where every position carries a label from a finite alphabet and a data value from an infinite domain. The logic we consider is two-way, contains “future” and “past” modalities, where these are considered as reflexive and transitive relations, and data equality and inequality tests. This logic corresponds to the fragment of XPath with the ‘following-sibling-or-self’ and ‘preceding-sibling-or-self’ axes over a data word. We show that this problem is decidable, EXPSPACE-complete. This is surprising considering that with the strict navigation relations (non-reflexive) the satisfiability problem is undecidable. To show this we first reduce the problem to a derivation problem for an infinite transition system, and then we show how to abstract this problem into a reachability problem of a finite transition system.

I. INTRODUCTION

We work on data words. A data word is a finite string where each position carries a label from a finite alphabet and a data value from some infinite domain. This kind of structure has been considered in the realm of semistructured data, timed automata, program verification, and generally in systems manipulating data values. Finding automata and logics with decidable satisfiability problems over data words is an important quest when studying systems manipulating data values.

Suppose that a data word models a database of ‘actions’ that were performed throughout the time, where the linear structure of the data word corresponds to the flow of time. Over these structures one may want to model the behavior of a system, and may want to check that the specified behavior is satisfiable. In other words, that it does not contain any contradiction. The satisfiability problem is crucial to solve this and other static analysis problems.

In this paper we consider a formalism to express properties of data words, where the data values can only be tested for equality or inequality. Logics or automata that can test equality of data values and that can express meaningful properties are important to find, although they become very easily undecidable.

We investigate a natural two-way logic on data words and show that it has a decidable EXPSPACE satisfiability problem. We call this logic *Two-way Path Logic* (or PathLogic for short). PathLogic can be seen as an extended temporal logic [DL09] or as a fragment of XPath [CD99] on data words. It has *future* and *past* modalities to navigate the word and data equality test expressions. For example, it can test the following property of a node: “there is a node x to the left and two nodes $y < z$ to the right with labels a , b , and c respectively, such that x and z have the same data value”. This property is very easily

expressed in PathLogic with the formula

$$\langle \overleftarrow{[a]} = \overrightarrow{[b] \cdot [c]} \rangle. \quad (*)$$

Although we can test the relative appearances of labels, the left and right relations are interpreted as transitive-reflexive relations. Thus, $\langle \overleftarrow{[a]} = \overrightarrow{[b] \cdot [c]} \rangle$ is equivalent to $\langle \overleftarrow{[a]} = \overrightarrow{[b]} \rangle$.

Our main result is the following.

Main Theorem. *The satisfiability problem for PathLogic is EXPSPACE-complete.*

This logic has links with other logics. As mentioned before, PathLogic is inspired in the logic for XML documents XPath with data equality tests and the reflexive-transitive axes ‘following-sibling-or-self’ and ‘preceding-sibling-or-self’¹. We note this axes \rightarrow^* and $*\leftarrow$ respectively for economy of space and XPath($\rightarrow^*, *\leftarrow, =$) to the logic. PathLogic is in fact expressive-equivalent (through an EXPTIME translation) to XPath($\rightarrow^*, *\leftarrow, =$) interpreted over data words. Indeed, PathLogic is a simplification of the syntax of XPath($*\leftarrow, \rightarrow^*, =$) that preserves only the essential. This has also the advantage of simplifying the decidability proofs. In terms of XPath our result is translated as follows.

Theorem 1. *The satisfiability problem for XPath($*\leftarrow, \rightarrow^*, =$) is in 2EXPSPACE.*

This may be surprising given the known results in this area. In fact, [FS09] shows that XPath($\rightarrow^+, =$) is decidable with non-primitive recursive complexity. This is the fragment that contains the non-reflexive ‘following-sibling’ axis. Even more, XPath($^+\leftarrow, \rightarrow^+, =$) or even XPath($*\leftarrow, \rightarrow^+, =$) are undecidable. Simply reducing the navigation by having reflexive-transitive relations instead of only transitive, we change from an undecidable problem to an elementarily decidable one. Thus, the \rightarrow^* and $*\leftarrow$ axes behave surprisingly much better. This result opens some promising conjectures on the decidability of several fragments of XPath on XML documents.

Decidable two-way logics on data words are scarce. Very often two-way logics on data words have an undecidable satisfiability problem. One such prominent example is the case of $LTL_1^\downarrow(F, F^{-1})$, a temporal logic with with future and past modalities and one register to store and compare data values. Whether we interpret the F and F^{-1} modalities as either reflexive or strict, its satisfiability problem remains undecidable. Further, if only one of these modalities is allowed, it becomes

¹Strictly speaking these axes do not exist in XPath [CD99], they must be interpreted as the reflexive version of the ‘following-sibling’ and ‘preceding-sibling’ axes respectively.

decidable, but with non-primitive recursive complexity [FS09]. PathLogic can be seen as a fragment of $LTL_1^\downarrow(F, F^{-1})$ (this is the fragment identified in [FS09] as “simple- $LTL_1^\downarrow(F, F^{-1})$ ”).

Related work

This paper addresses a topic related to the work of [FS09], about the complexity and decidability of logics for data words and trees with transitive relations. [FS09] showed lower bounds for several fragments of LTL_1^\downarrow and XPath with transitive relations for navigation. In particular, it was shown that all fragments of XPath with data equality tests containing transitive and non-reflexive axes \rightarrow^+ , $+\leftarrow$ or \uparrow^+ have a satisfiability problem of non-primitive recursive complexity. To obtain such results it was necessary to study the behavior of XPath on non-branching structures, *i.e.*, on *data words*. The present paper also centers on data words, showing a positive result with the spirit of obtaining in the future positive results for fragments of XPath over XML documents.

Bojańczyk et al. [BMS⁺06] study the satisfiability problem for first order logics with two variables on data words. In particular they study the fragment $FO^2(<, \sim)$. This is first-order logic restricted to two variables, a binary relation $<$ which corresponds to the order in the word, and a binary relation \sim to test for data equality or inequality of nodes. They show that the satisfiability problem is decidable, NEXPTIME-complete. In fact, they show that the decidability is also preserved if we allow any predicate ‘ $+k$ ’, relating any two nodes at distance k in the data word.² This logic is in some sense two-way, since it can navigate both in the future and past directions. Since it also includes equality and inequality of nodes, the transitive $<$ relation and the reflexive-transitive relation \leq are one definable in terms of the other. Although close, $FO^2(<, \sim)$ is incomparable with PathLogic in terms of expressive power. This decidability result was extended to $FO^2(<, \preceq)$, where $x \preceq y$ tests that x carries a data value smaller or equal to that of y . The satisfiability of this logic is EXPSpace [SZ10].

Demri and Lazić [DL09] study the logic LTL_1^\downarrow , which is a two-way temporal logics on data words. It is based on the linear temporal logic LTL with special constructs to store and test for equality of data values. This logic with Future, Next and Until modalities is shown to be decidable [DL09], but any fragment with Future and Past modalities becomes undecidable [FS09]. Nevertheless, [DL09] identifies a decidable fragment that is equivalent to $FO^2(<, \sim, +1, \dots, +k)$, and hence decidable. [KSZ10] also investigates a decidable two-way fragment but this time over *multi-attributes data words*, where every position carries a set of data values. These logics are intimately related with FO^2 and they are hence also incomparable to PathLogic in expressive power.

Organization: In Section II we introduce our logic PathLogic and we briefly comment on the relation with XPath, LTL_1^\downarrow , and FO^2 . Sections III, IV and V are dedicated to prove our main result, that the satisfiability problem of PathLogic

is in EXPSpace. In Section VI we discuss about the lower bound. Section VII contains some corollaries of our result in terms of XPath fragments. Finally, in Section VIII we propose some conjectures and future lines of work. Due to space limitations some proofs are in the Appendix.

II. A TWO-WAY LOGIC

Notation: We consider a finite word over \mathbb{E} as a function $\mathbf{w} : [n] \rightarrow \mathbb{E}$ for some $n > 0$, where $[n] = \{1, \dots, n\}$. We define the set of words as $Words(\mathbb{E}) := \{[n] \rightarrow \mathbb{E} \mid n > 0\}$. We write $\text{pos}(\mathbf{w}) = \{1, \dots, n\}$ to denote the set of **positions** (that is, the domain of \mathbf{w}). Given $\mathbf{w} \in Words(\mathbb{E})$ and $\mathbf{w}' \in Words(\mathbb{F})$ with $\text{pos}(\mathbf{w}) = \text{pos}(\mathbf{w}') = P$, we write $\mathbf{w} \otimes \mathbf{w}' \in Words(\mathbb{E} \times \mathbb{F})$ for the word such that $\text{pos}(\mathbf{w} \otimes \mathbf{w}') = P$ and $(\mathbf{w} \otimes \mathbf{w}')(x) = (\mathbf{w}(x), \mathbf{w}'(x))$. A **data word** is a word $\mathbf{w} = \mathbf{a} \otimes \mathbf{d} \in Words(\mathbb{A} \times \mathbb{D})$, where \mathbb{A} is a finite alphabet of letters and \mathbb{D} is an infinite domain. We use \mathbb{A} for a finite alphabet, \mathbb{D} as an infinite domain, and we refer to its elements with the letters $a, b, c \in \mathbb{A}$ and $d, e \in \mathbb{D}$. We use \mathbf{w} , \mathbf{a} and \mathbf{d} for words of $Words(\mathbb{A} \times \mathbb{D})$, $Words(\mathbb{A})$ and $Words(\mathbb{D})$ respectively. We use the loosely term *node* to denote the pair of label and data value of a position. We write $\wp(S)$ and $\wp_{<\infty}(S)$ to denote the powerset and finite powerset of a set S respectively.

We work with a logic with a simple syntax and same expressive power as XPath($*\leftarrow, \rightarrow^*, =$). We call this logic PathLogic. This is an equivalent formalism which is more convenient to work with for the proofs presented here.

PathLogic is a two-sorted formulas with *path expressions* (that we note α, β) and *node expressions* (φ, ψ). A path expression is a succession of node expressions, noted $\alpha = [\psi_1] \cdot \dots \cdot [\psi_n]$ and we say that a position i of the data word reaches another position j through α if there are witness positions $i \leq i_1 \leq \dots \leq i_n = j$ where φ_k holds true at position i_k , for every k . On the other hand, a node expression is a boolean combination of tests for labels or formulas $\langle \overleftarrow{\alpha} = \overrightarrow{\beta} \rangle$ that demand the existence of one position to the left and one to the right reachable by $\overleftarrow{\alpha}$ and $\overrightarrow{\beta}$ respectively, that carry the same data value. Analogously, $\langle \overleftarrow{\alpha} \neq \overrightarrow{\beta} \rangle$ expresses the same property except that the positions must carry *different* data values. Expressions are formally defined as follows.

$$\begin{aligned} \varphi &::= \langle \overleftarrow{\alpha} = \overrightarrow{\beta} \rangle \mid \langle \overleftarrow{\alpha} \neq \overrightarrow{\beta} \rangle \mid a \mid \neg \varphi \mid \varphi \vee \varphi \\ \alpha, \beta &::= \varepsilon \mid [\varphi] \cdot \alpha \end{aligned}$$

And semantics are defined according to the intuition we just gave. See Figure 1 for a more precise definition.

In this context, we define $\mathbf{w}, i \models \varphi$ iff $i \in \llbracket \varphi \rrbracket^{\mathbf{w}}$, and $\mathbf{w} \models \varphi$ iff $\mathbf{w}, 1 \models \varphi$. We also define the test $\langle \overrightarrow{\alpha} \rangle$, which checks if there is a node reachable by α , as $\langle \overrightarrow{\alpha} = \overrightarrow{\alpha} \rangle$, and likewise for $\langle \overleftarrow{\alpha} \rangle$. The **satisfiability problem** is the problem of, given $\varphi \in \text{PathLogic}$, whether there is a data word \mathbf{w} such that $\mathbf{w} \models \varphi$. Our main result is that this problem is decidable.

Comparison with other logics

a) XPath: PathLogic corresponds precisely to XPath($*\leftarrow, \rightarrow^*, =$) in terms of expressiveness. The

²For this logic no primitive-recursive upper bound is known.

$$\begin{aligned}
 \llbracket a \rrbracket^{\mathbf{w}} &\stackrel{\text{def}}{=} \{i \in \text{pos}(\mathbf{w}) \mid \mathbf{a}(i) = a\} \\
 \llbracket \varphi \vee \psi \rrbracket^{\mathbf{w}} &\stackrel{\text{def}}{=} \llbracket \varphi \rrbracket^{\mathbf{w}} \cup \llbracket \psi \rrbracket^{\mathbf{w}} \\
 \llbracket \neg \varphi \rrbracket^{\mathbf{w}} &\stackrel{\text{def}}{=} \text{pos}(\mathbf{w}) \setminus \llbracket \varphi \rrbracket^{\mathbf{w}} \\
 \llbracket \langle \overleftarrow{\alpha} = \overrightarrow{\beta} \rangle \rrbracket^{\mathbf{w}} &\stackrel{\text{def}}{=} \{i \in \text{pos}(\mathbf{w}) \mid \exists j, k \in \text{pos}(\mathbf{w}), \\
 &\quad (i, j) \in \llbracket \alpha \rrbracket_{\leftarrow}^{\mathbf{w}}, (i, k) \in \llbracket \beta \rrbracket_{\rightarrow}^{\mathbf{w}}, \mathbf{d}(j) = \mathbf{d}(k)\} \\
 \llbracket \langle \overleftarrow{\alpha} \neq \overrightarrow{\beta} \rangle \rrbracket^{\mathbf{w}} &\stackrel{\text{def}}{=} \{i \in \text{pos}(\mathbf{w}) \mid \exists j, k \in \text{pos}(\mathbf{w}), \\
 &\quad (i, j) \in \llbracket \alpha \rrbracket_{\leftarrow}^{\mathbf{w}}, (i, k) \in \llbracket \beta \rrbracket_{\rightarrow}^{\mathbf{w}}, \mathbf{d}(j) \neq \mathbf{d}(k)\} \\
 \llbracket \varepsilon \rrbracket_{\rightarrow}^{\mathbf{w}} &\stackrel{\text{def}}{=} \llbracket \varepsilon \rrbracket_{\leftarrow}^{\mathbf{w}} = \{(i, i) \mid i \in \text{pos}(\mathbf{w})\} \\
 \llbracket [\varphi] \cdot \alpha \rrbracket_{\rightarrow}^{\mathbf{w}} &\stackrel{\text{def}}{=} \{(i, j) \in (\text{pos}(\mathbf{w}))^2 \mid \exists k \geq i \text{ s.t.} \\
 &\quad k \in \llbracket \varphi \rrbracket^{\mathbf{w}}, (k, j) \in \llbracket \alpha \rrbracket_{\rightarrow}^{\mathbf{w}}\} \\
 \llbracket [\varphi] \cdot \alpha \rrbracket_{\leftarrow}^{\mathbf{w}} &\stackrel{\text{def}}{=} \{(i, j) \in (\text{pos}(\mathbf{w}))^2 \mid \exists k \leq i \text{ s.t.} \\
 &\quad k \in \llbracket \varphi \rrbracket^{\mathbf{w}}, (k, j) \in \llbracket \alpha \rrbracket_{\leftarrow}^{\mathbf{w}}\}
 \end{aligned}$$

Fig. 1. Syntax and semantics of PathLogic for a data word $\mathbf{w} = \mathbf{a} \otimes \mathbf{d}$.

only difference is that the syntax is simplified. While $\text{XPath}(*\leftarrow, \rightarrow^*, =)$ may change of direction inside a path expression, PathLogic cannot. But in fact this does not change the expressive power: we can always factorize path expressions with mixed axes into formulas whose path expressions use only one axis. Thus, there is an EXPTIME translation from $\text{XPath}(*\leftarrow, \rightarrow^*, =)$ into PathLogic.

b) LTL_1^\downarrow : PathLogic also corresponds to the fragment of *simple-LTL* $^\downarrow_1$ as defined in [FS09], denoted by $\text{sLTL}_1^\downarrow(\mathbf{F}, \mathbf{F}^{-1})$.³ Here, the \mathbf{F} and \mathbf{F}^{-1} are interpreted as reflexive-transitive relations. A formula of LTL_1^\downarrow in negated normal form is said to be *simple* if (i) there is at most one occurrence of \uparrow within the scope of each occurrence of \downarrow and, (ii) there is no negation between an occurrence of \uparrow and its matching \downarrow , except maybe immediately before \uparrow . We denote by sLTL_1^\downarrow the fragment of LTL_1^\downarrow containing only simple formulas. The logic $\text{sLTL}_1^\downarrow(\mathbf{F}, \mathbf{F}^{-1})$ is equivalent to $\text{XPath}(*\leftarrow, \rightarrow^*, =)$ (and hence also to PathLogic) in terms of expressive power. As shown in [FS09], there is a PTIME translation in both senses.

c) FO^2 : As already mentioned, $\text{FO}^2(<, \sim)$ is incomparable with PathLogic. Indeed, $\text{FO}^2(<, \sim)$ can express that the word contains at least two nodes, while we cannot express this property using only *reflexive*-transitive modalities. In turn, PathLogic can easily express a property like “all nodes labeled a verify the property (\star) ” (as discussed in the Introduction), that is unlikely to be expressible in $\text{FO}^2(<, \sim)$. We argue that over a fixed alphabet, PathLogic captures $\text{FO}^2(<, \sim)$ in terms of expressive power over the class of *duplicate-free data words*, defined as those words that do not contain two successive identical nodes (*i.e.*, that do not contain $(a, d)(a, d)$ as a substring for any a and d).

Lemma 1. *Given a fixed alphabet, PathLogic is more expressive than $\text{FO}^2(<, \sim)$ on duplicate-free data words.*

³This fragment has no connection with the *simple* fragment of LTL_1^\downarrow defined in [DL09].

We solve the satisfiability problem through a series of two reductions. First, we reduce the problem into a derivation problem for an infinite transition system \succrightarrow . And second, we abstract this transition system, reducing this problem into a similar problem for a *finite* transition system \rightarrow_K . In truth, \rightarrow_K is not a finite transition system, but we will see that it can be easily restricted to a finite one.

Organization: In this section we present the general outline of the decidability proof, and we show our first reduction to a derivation problem for \succrightarrow . Section IV contains the most difficult result, which is the reduction to the derivation problem for \rightarrow_K . Finally, in Section V we solve the derivation problem for \rightarrow_K with an EXPSpace procedure.

A. Outline of the proof

We first reduce the problem of testing whether some formula φ is satisfiable into a derivation problem for an infinite state transition system \succrightarrow that depends on φ . This problem consists in testing whether there exists a finite sequence $\mu_1 \succrightarrow \dots \succrightarrow \mu_n$ such that μ_1 and μ_n satisfy certain local properties. The domain of this transition system is the set of *mosaics*: abstractions of a position of a data word, containing all the data values that can be found to the right and to the left, and further with which path expressions of the logic they can be reached. The transition system \succrightarrow corresponds to the property of when two of these mosaics can be *matched*: that is, that could abstract two positions that are one next to another. We show that from a derivation $\mu_1 \succrightarrow \dots \succrightarrow \mu_n$ with some good properties, one can build a data word of length n that satisfies the input formula φ . In turn, if there is no such derivation, φ is unsatisfiable. To obtain this reduction, we work with formulas of PathLogic in a certain normal form. This normal form enables to have a simple definition of the transition \succrightarrow specially convenient for our proofs.

Next, in Section IV, we solve the derivation problem for \succrightarrow by a reduction to a similar derivation problem for another transition system \rightarrow_K . This transition system is parametrized by a set K of data values. The definition of \rightarrow_K introduces an important concept of *rigid* data values, which are data values that play a determined function in a sequence $\mu_1 \succrightarrow \dots \succrightarrow \mu_n$. We denote by K such data values of important interest. Based on this concept we define a quasi-order \leq_K over the set of mosaics, which is monotone with respect to \succrightarrow . This monotonicity serves to finally show the reduction from derivation problem for \succrightarrow to the derivation problem for \rightarrow_K .

The derivation problem for \rightarrow_K may seem in appearance more difficult than for \succrightarrow , since we are dealing now with a *family* of transition systems \rightarrow_K , one for every K . Nevertheless, in Section V we show that the parameter K is harmless and we can bound and fix a value for it. Also, we will see that we only need to work with the minimal classes of equivalences of \leq_K which are also finite and bounded. These observations give us indeed a finite transition system. Hence, we can use a standard reachability algorithm that solves

the derivation problem for \rightarrow_K and hence the satisfiability problem for PathLogic.

Organization of this section: We give a more detailed informal description of the normal form for PathLogic and the transition system \rightarrow , together with some fundamental definitions and notation that we use throughout. We then state the target problem for the reduction.

B. A normal form

We will assume a certain normal form of the formula φ to test for satisfiability. This will simplify the necessary machinery to solve our problem. This normal form consists simply of having formulas without nesting of data tests. That is, we avoid treating formulas like

$$\langle \overleftarrow{[a] = [b]} \rangle = \overrightarrow{[c] \cdot [d]} \rangle.$$

Let us call $BC(labels)$ to all the formulas which are a boolean combination of tests for labels of \mathbb{A} . If a formula is such that all its path expressions α contain only tests for labels (i.e., $\alpha \in (BC(labels))^*$ abusing notation) we call it a **non-recursive** formula.

In the normal form we suppose that $\varphi = \varphi_1 \wedge \varphi_2$ where φ_1 is a non-recursive formula and φ_2 is a conjunction of tests of the form “every node with label a satisfies ψ ” for some non-recursive formula ψ . Formally, φ_2 contains a conjunction of tests of the form

$$\neg \langle \overrightarrow{[a \wedge \neg \psi]} \rangle$$

for $a \in \mathbb{A}$ a label and ψ a non-recursive formula. Then, we obtain the following.

Lemma 2 (normal form). *There is a polynomial-time translation that for every formula $\psi \in PathLogic$ returns a formula φ in normal form such that ψ is satisfiable iff φ is satisfiable.*

Proof: Given a formula φ we define the alphabet of the formula in normal form as the set of all subformulas of φ plus some extra element \perp . We denote these letters by a_ψ , for ψ a subformula of φ . Given a formula η let us call $tr(\eta)$ to the same formula where every instance of a path expression $[\psi_1] \cdots [\psi_n]$ (which does not appear nested inside another path expression) is replaced by $[a_{\psi_1}] \cdots [a_{\psi_n}]$.

To build the normal form, we define on the one hand $\varphi_1 = tr(\varphi)$. On the other hand, we build φ_2 as a conjunction of formulas

$$\neg \langle \overrightarrow{[a_\psi \wedge \neg tr(\psi)]} \rangle$$

for all subformulas ψ of φ . It follows that this translation is polynomial since the set of subformulas of φ is polynomial. ■

We fix once and for all \mathbb{A} as the finite alphabet and \mathbb{D} as the infinite domain we are going to work with. Given $\varphi = \varphi_1 \wedge \varphi_2$ in normal form, we call Ω_φ to all non-recursive subformulas of φ , and $\gamma_\varphi : \mathbb{A} \rightarrow \Omega_\varphi$ to the function such that $\gamma_\varphi(a) = \psi$ iff φ_2 contains $\neg \langle [a \wedge \neg \psi] \rangle$. Finally, we define $\Omega_\varphi^p = \{\alpha \in \Omega_\varphi \mid$

α is a path expression different from $\varepsilon\}$. We exclude ε from Ω_φ^p simply because we are interested in those path expressions that are ‘moving’, unlike ε .

C. The transition system

Given a formula φ in normal form, we show an abstraction of the important information that we need to keep track of in a model that satisfies φ . Each piece of information is called a **mosaic**. A mosaic maintains all necessary facts about some position i of a data word \mathbf{w} such that $\mathbf{w} \models \varphi$, namely

- the letter and datum of the current position,
- the data values that are to be found to the right of i , and with which paths of Ω_φ^p are reachable,
- and likewise for the data values to the left.

We say that a mosaic is the φ -**abstraction** of a position i of a data word \mathbf{w} if it contains all the aforementioned details of the position. A mosaic that abstracts a position i of \mathbf{w} contains enough information to test whether ψ is satisfied at position i (i.e., $\mathbf{w}, i \models \psi$) for any subformula ψ of φ . In order to be a valid mosaic, the mosaic must also comply with $\gamma_\varphi(a)$ where $a \in \mathbb{A}$ is the label of the mosaic.

Given two mosaics μ_1, μ_2 , we can check if they match, that one can be next to the other, that there are no incompatibilities. For example, suppose that μ_1 is the φ -abstraction of a position i . If μ_1 asserts that there is a data value d different from the current one that can be reached with $[a]$ navigating to the right, then the mosaic that abstracts $i + 1$ must also contain the information that d can be reached with $[a]$. Otherwise these two mosaics would not match. We will write \rightarrow for this matching relation, and we call a \rightarrow **sequence** to any sequence $\mu_1 \rightarrow \cdots \rightarrow \mu_n$. (In general, we call a \rightarrow sequence to $\mu_1 \rightarrow \cdots \rightarrow \mu_n$ for any transition \rightarrow .)

Indeed, the φ -abstraction of the positions of a data word w and formula φ such that $w \models \varphi$ verify that they are a \rightarrow sequence. Moreover, the leftmost mosaic must not contain information of data values to the left simply because there are no more elements to the left. For the same reasons, the rightmost mosaic must not contain information on the data values to the right. We call such mosaics **left-complete** and **right-complete** respectively. A \rightarrow sequence is **complete** if the leftmost element is left-complete and the rightmost element is right-complete. We show that from any complete \rightarrow sequence for φ whose first element verifies φ , we can build a data word satisfying φ . And in turn, the sequence of φ -abstractions of a word that satisfies φ is indeed a complete \rightarrow sequence for φ . Now the satisfiability problem reduces to the problem of testing if there is a complete \rightarrow sequence whose first element satisfies φ . We call this the **derivation problem**.

Definitions: Let us now turn to the basic definitions. An important object we will work with is that of a *profile*. A profile of a data value d is a subset of Ω_φ^p describing all the possible ways to reach d in a data word \mathbf{w} from a position $i \in \text{pos}(\mathbf{w})$ by going in one direction (either to the right or to the left). We use $\pi, \rho \subseteq \Omega_\varphi^p$ as names for profiles.

We define some useful operations on profiles. There is a constant profile σ_a for every $a \in \mathbb{A}$, which consists in all the path expressions that can be satisfied locally at a certain node, assuming that the node's label is a . (Note that σ_a is a constant whereas π, ρ are names.)

$$\sigma_a \stackrel{\text{def}}{=} \{\alpha \in \Omega_\varphi^p \mid \alpha = [\psi_1] \cdot \dots \cdot [\psi_n], \forall i : a \models \psi_i\} \quad a \in \mathbb{A}$$

where $a \models \psi_i$ means that ψ_i is verified by the label a —remember that $\psi_i \in BC(\text{labels})$. Also, given a profile π and a label $a \in \mathbb{A}$, we write $a\pi$ for the concatenation of the profiles of σ_a and π .

$$a\pi \stackrel{\text{def}}{=} \pi \cup \{\alpha \cdot \beta \in \Omega_\varphi \mid \alpha \in \sigma_a, \beta \in \pi\}$$

Finally, we write $\pi + a$ for the union of the profiles π and σ_a .

$$\pi + a \stackrel{\text{def}}{=} \pi \cup \sigma_a$$

Remark 1. These are idempotent and commutative operations: $a(a\pi) = a\pi$; $(\pi + a) + a = \pi + a$; $a(\pi + a) = a\pi + a$.

Definition 1 (mosaic). A *mosaic* over φ is a tuple $\mu = (a, d, \overleftarrow{\chi}, \overrightarrow{\chi})$ where $a \in \mathbb{A}$, $d \in \mathbb{D}$, and $\overleftarrow{\chi}, \overrightarrow{\chi} \in \wp_{<\infty}(\mathbb{D} \times \Omega_\varphi^p)$. The idea is that $\overleftarrow{\chi}$ contains the profiles of the data values that can be found to the left and $\overrightarrow{\chi}$ of those to the right.

We use letters μ, ν for mosaics and we fix the nomenclature $\mu_i = (a_i, d_i, \overleftarrow{\chi}_i, \overrightarrow{\chi}_i)$ for every i . We adopt the following notation for $\chi \in \{\overleftarrow{\chi}, \overrightarrow{\chi}\}$, $\alpha \in \Omega_\varphi^p$, and $e \in \mathbb{D}$.

$$\begin{aligned} \chi(e) &\stackrel{\text{def}}{=} \{\alpha \in \Omega_\varphi^p \mid (e, \alpha) \in \chi\} & \mu(e) &\stackrel{\text{def}}{=} (\overleftarrow{\chi}(e), \overrightarrow{\chi}(e)) \\ \chi(\alpha) &\stackrel{\text{def}}{=} \{e \in \mathbb{D} \mid (e, \alpha) \in \chi\} \end{aligned}$$

We call $\mu(e)$ the *profile* of e in μ . Note that we call a profile either to a subset of Ω_φ^p or to a pair of subsets of Ω_φ^p . The intended meaning will be clear from the context.

Given a formula $\varphi = \varphi_1 \wedge \varphi_2$ in normal form, a data word \mathbf{w} and a position i , we say that a mosaic μ over φ is the φ -**abstraction** of (\mathbf{w}, i) if it is defined in the expected way by taking into account the data values of \mathbf{w} and the path expressions Ω_φ^p . That is, if $\mathbf{w} = \mathbf{a} \otimes \mathbf{d}$ and $\mu = (a, d, \overleftarrow{\chi}, \overrightarrow{\chi})$ we define $(a, d) = \mathbf{w}(i)$, $\overleftarrow{\chi} = \{(e, [\psi_1] \cdot \dots \cdot [\psi_n]) \in \mathbb{D} \times \Omega_\varphi^p \mid \text{there are } i \leq j_1 \leq \dots \leq j_n \leq |\mathbf{w}| \text{ where } \mathbf{a}(j_k) \models \psi_k \text{ for all } k \text{ and } \mathbf{d}(j_n) = e\}$. Likewise for $\overrightarrow{\chi}_i$.

For every atomic expression $\langle \alpha = \beta \rangle$, $\langle \alpha \neq \beta \rangle$ or a from Ω_φ one can decide if the formula holds at position i of \mathbf{w} simply by inspecting the φ -abstraction of (\mathbf{w}, i) . Hence, one can define a satisfaction relation \models , where the following holds.

Lemma 3. For every formula φ in normal form and subformula $\psi \in \Omega_\varphi$ and for every data word \mathbf{w} and position i , $\mathbf{w}, i \models \varphi$ iff $\mu \models \psi$, where μ is the φ -abstraction of (\mathbf{w}, i) .

Not all mosaics are consistent with the φ -abstraction. For example, a mosaic over φ with a label $a \in \Omega_\varphi$ such that $\overrightarrow{\chi}([a]) = \emptyset$ cannot be the φ -abstraction of any (\mathbf{w}, i) : $\overrightarrow{\chi}([a])$ must contain the mosaic's datum.

Definition 2 (validity). A mosaic $(a, d, \overleftarrow{\chi}, \overrightarrow{\chi})$ is *valid* if it satisfies the following restrictions for all $\chi \in \{\overleftarrow{\chi}, \overrightarrow{\chi}\}$

1. for every $d' \in \mathbb{D}$, $\chi(d')$ is suffix closed: if $\alpha \cdot \beta \in \chi(d')$ then $\beta \in \chi(d')$,
2. $\overrightarrow{\chi}(d) = a\pi + a$ and $\overleftarrow{\chi}(d) = a\rho + a$ for some $\pi, \rho \subseteq \Omega_\varphi^p$,
3. for every $d' \in \mathbb{D}$, $\chi(d') = a\pi$ for some $\pi \subseteq \Omega_\varphi^p$, and
4. $\mu \vdash \gamma_\varphi(a)$.

Finally, $\mathcal{M}(\varphi)$ is the set of all valid mosaics over φ . Henceforth we say *mosaic* to denote a *valid* mosaic unless otherwise stated. Observe that the φ -abstraction of a data word \mathbf{w} and position i results always in a valid mosaic.

Lemma 4. The φ -abstraction of (\mathbf{w}, i) is a valid mosaic of $\mathcal{M}(\varphi)$, assuming $\mathbf{w} \models \varphi$ and $i \in \text{pos}(\mathbf{w})$.

The importance of the normal form presented before is that it allows us to work separately on each data value. The left profile of a data value at position $i + 1$ depends solely on the profile of the data value at position i , and the labels of these positions. In fact, all the interaction between the profiles of different data values is encapsulated in the restrictions imposed by the function γ_φ . This is due to the fact that all path expressions contained in Ω_φ^p depend solely on the label of the mosaic. We can hence define a relation that tests whether it is possible that two profiles of a data value can abstract two consecutive positions.

We define the relations $\overrightarrow{\mathbf{m}}, \overleftarrow{\mathbf{m}}$ that evidence this dependency of profiles that we discussed for a given data value d . Intuitively, the first three components of the relation specify the situation for d on the left position: the label of the node, whether d is equal (1) or not (0) to the current data value, and its profile ($\overleftarrow{\chi}(d)$ or $\overrightarrow{\chi}(d)$ respectively for $\overrightarrow{\mathbf{m}}, \overleftarrow{\mathbf{m}}$). The other three components contain a similar description for the same data value on the right position. The idea is that $\overrightarrow{\mathbf{m}}$ and $\overleftarrow{\mathbf{m}}$ hold if it is possible to have this situation in a data word.

$$\overrightarrow{\mathbf{m}}(a, 1, \pi_1, b, j, \pi_2) \stackrel{\text{def}}{\iff} \pi_1 = a\pi_2 + a \quad (1)$$

$$\overrightarrow{\mathbf{m}}(a, 0, \pi_1, b, j, \pi_2) \stackrel{\text{def}}{\iff} \pi_1 = a\pi_2 \quad (2)$$

$$\overleftarrow{\mathbf{m}}(a, i, \rho_1, b, j, \rho_2) \stackrel{\text{def}}{\iff} \overrightarrow{\mathbf{m}}(b, j, \rho_2, a, i, \rho_1) \quad (3)$$

We define $\mathbf{m}(a, i, (\rho_1, \pi_1), b, j, (\rho_2, \pi_2))$ iff $\overrightarrow{\mathbf{m}}(a, i, \pi_1, b, j, \pi_2)$ and $\overleftarrow{\mathbf{m}}(a, i, \rho_1, b, j, \rho_2)$. We use the expression $d \stackrel{?}{=} d'$ to denote 1 if $d = d'$ or 0 otherwise. Now we can define the transition system \mapsto :

Definition 3 (\mapsto). We say that μ_1 and μ_2 *match*, and we write it $\mu_1 \mapsto \mu_2$, iff for all $d \in \mathbb{D}$

$$\mathbf{m}(a_1, d \stackrel{?}{=} d, \mu_1(d), a_2, d \stackrel{?}{=} d, \mu_2(d)) .$$

We define the necessary properties that the abstractions of the borders (i.e., of the first and last elements) should admit.

Definition 4. We define the left and right demands of a mosaic $\mu = (a, d, \overleftarrow{\chi}, \overrightarrow{\chi})$. A left- (resp. right-) demand consists of the

data values that necessarily have to be found to the left (resp. to the right).

$$\overleftarrow{\text{dem}}(\mu) \stackrel{\text{def}}{=} \{(e, \rho) \in \overleftarrow{\chi} \mid (e \neq d \wedge \rho \neq \emptyset) \vee (e = d \wedge \rho \neq \sigma_a)\}$$

$$\overrightarrow{\text{dem}}(\mu) \stackrel{\text{def}}{=} \{(e, \pi) \in \overrightarrow{\chi} \mid (e \neq d \wedge \pi \neq \emptyset) \vee (e = d \wedge \pi \neq \sigma_a)\}$$

We say that a mosaic is left-complete (resp. right-complete) if it does not have left-demands (resp. right-demands).

Definition 5. A \succrightarrow sequence is a sequence of mosaics μ_1, \dots, μ_n such that $\mu_i \succrightarrow \mu_{i+1}$ for all $i < n$. We say that the sequence is **complete** if $\overleftarrow{\text{dem}}(\mu_1) = \overrightarrow{\text{dem}}(\mu_n) = \emptyset$.

We state the derivation problem for a transition system \rightarrow .

The derivation problem for \rightarrow	
INPUT:	A formula $\varphi = \varphi_1 \wedge \varphi_2$ in normal form.
OUTPUT:	Is there a complete \rightarrow sequence over $\mathcal{M}(\varphi)$ whose leftmost mosaic μ verifies $\mu \vdash \varphi_1$?

D. From the satisfiability problem to the derivation problem

We can reduce the satisfiability problem for φ into a the derivation problem for \succrightarrow over $\mathcal{M}(\varphi)$.

Proposition 1. A PathLogic formula $\varphi = \varphi_1 \wedge \varphi_2$ in normal form is satisfiable iff the derivation problem for \succrightarrow with input φ has a positive answer.

IV. THE DERIVATION PROBLEM FOR \succrightarrow

In this section we show how to solve the derivation problem for \succrightarrow . The general idea is that we introduce a quasi-order \leq_K on the set of mosaics, parametrized by a set of data values K . The parameter of this relation relies on the concept of *rigid values*, which are data values with a specific role in an \succrightarrow sequence. We show that \leq_K is (upward) monotone with respect to \succrightarrow . This monotonicity allows us to define a transition relation \rightarrow_K that has an equivalent derivation problem. We postpone the solving of the derivation problem for \rightarrow_K to Section V.

We fix once and for all a formula φ of PathLogic in normal form, and $\mathcal{M}(\varphi)$ as the domain of the transition relation \succrightarrow .

A. Rigid and flexible data values

Not all data values are the same in a data word satisfying a formula, different data values have different roles. We introduce 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 word. Let us give some more precise intuition. We use the logic PathLogic to make this intuition clear, but we will then state the definitions in terms of our transition system.

Given a data word \mathbf{w} such that $\mathbf{w} \models \varphi$, suppose there is a data value d such that if d can be accessed through a path α from some position i , then d is the *only* value that can be accessed through α at position i . When there is such a d we call it a **rigid value** for i , 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 of \mathbf{w} we say that d is rigid for \mathbf{w} . All the remaining data values (which are the **flexible values**) of \mathbf{w} play the role of assuring that “there are at least two data values reachable through α from position i ” for some α and i . As such, its importance is only relative. For example, if \mathbf{w} is a data word satisfying φ and containing d as a flexible value, then consider \mathbf{w}' the result of replacing, for some fresh data value d' , every appearance of (a, d) in \mathbf{w} by $(a, d) (a, d')$. \mathbf{w}' will continue to satisfy φ . But this is not necessarily true if d is a rigid value. The same notions hold for our \succrightarrow sequences. We say that a mosaic μ is **K -compliant** for $K \subseteq D$ if all its rigid values are inside K .

Definition 6 (K -compliant). Given a finite set $K \subseteq \mathbb{D}$ we say that a mosaic $\mu = (a, d, \overleftarrow{\chi}, \overrightarrow{\chi})$ is K -compliant iff for all $\alpha \in \Omega_\varphi^p$ and $\chi \in \{\overleftarrow{\chi}, \overrightarrow{\chi}\}$

$$\chi(\alpha) = \{d'\} \text{ implies } d' \in K.$$

This notion of K -compliance will be central for our second reduction to the more abstract transition system. Observe that for every \succrightarrow sequence there is a trivial K such that all its mosaics are K -compliant: the set of all data values of the data word. Later on, we will see that we can effectively bound K as a function of φ . This bound is independent of the length of the \succrightarrow sequence.

Based on this idea of rigid values, we define a relation \leq_K specially suited for mosaics that are K -compliant. In this definition, if two mosaics are in the relation, then all the profiles from K must remain untouched. For any other flexible value we are only interested in having at least as many data values with a certain profile in the bigger profile. We also preserve the existence or non-existence of a flexible data value with a given profile.

Definition 7 (\leq_K). $\mu_1 \leq_K \mu_2$ iff:

1. they have the same label
2. $\forall d \in K, \mu_1(d) = \mu_2(d)$,
3. $\mu_1(d_1) = \mu_2(d_2)$ for d_i the data value of μ_i ,
4. $\forall (\rho, \pi) \in \wp(\Omega_\varphi^p) \times \wp(\Omega_\varphi^p)$,

$$|\mu_1^{-1}(\rho, \pi) \setminus K| \leq |\mu_2^{-1}(\rho, \pi) \setminus K|$$

where $n \leq m$ iff $0 = n = m$ or $1 \leq n \leq m$.

Lemma 5. Given $\mu_1 \leq_K \mu_2$ we have

1. if both mosaics are valid, then μ_1 is left-complete iff μ_2 is left-complete, and likewise for right-complete,
2. if μ_1 is K -compliant, then μ_2 is K -compliant,
3. if μ_1 is valid, then μ_2 is valid,
4. if both are K -compliant, for all $\psi \in \Omega_\varphi$, $\mu_1 \vdash \psi$ iff $\mu_2 \vdash \psi$.

B. A more abstract transition system

We define a transition system \rightarrow_K over the set of K -compliant mosaics, which is also parametrized by K . We then reduce the derivation problem for \succrightarrow to the derivation problem for \rightarrow_K .

Definition 8 (\rightarrow_K). $\mu_1 \rightarrow_K \mu_2$ iff μ_1 and μ_2 are K -compliant and $\mu_1 \leq_K \mu'_1 \succrightarrow \mu'_2 \geq_K \mu_2$ for some mosaics μ'_1, μ'_2 .

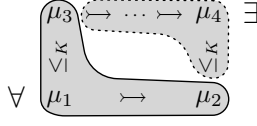


Fig. 2. The forward monotonicity Lemma.

Note that \rightarrow_K is a looser relation than \succrightarrow . This will enable us to easily bound the domain. But we leave this issue to the next Section V. Our present goal is to show that there is indeed a reduction.

Proposition 2. *There exists a K -compliant complete \rightarrow_K sequence over $\mathcal{M}(\varphi)$ whose first element μ_1 verifies $\mu_1 \vdash \varphi_1$ iff there exists a K -compliant complete \succrightarrow sequence over $\mathcal{M}(\varphi)$ whose first element μ'_1 verifies $\mu'_1 \vdash \varphi_1$.*

To prove the above proposition, first we need to show a monotonicity property of the transition system \succrightarrow . The next Lemma asserts that \succrightarrow is upward compatible with respect to \leq_K , as in Figure 2. This is a crucial technical Lemma of our proof.

Lemma 6 (forward monotonicity). *For every K -compliant mosaics $\mu_3 \geq_K \mu_1 \succrightarrow \mu_2$ there is a K -compliant mosaic μ_4 such that $\mu_3 \succrightarrow^+ \mu_4 \geq_K \mu_2$.*

Since all our definitions are symmetrical, by the proof of Lemma 7 we will also obtain a similar result for backward monotonicity.

Lemma 7 (backward monotonicity). *For every K -compliant mosaics $\mu_1 \succrightarrow \mu_2 \leq_K \mu_4$ there is a K -compliant mosaic μ_3 such that $\mu_1 \leq_K \mu_3 \succrightarrow^+ \mu_4$.*

Before going into the details of the proof, we state some simple facts.

Lemma 8 (about \mathbf{m}). *For any two valid mosaics $\mu_1, \mu_2 \in \mathcal{M}(\varphi)$ and $d \in \mathbb{D}$, suppose that $a_1 = a_2 = a$ for some $a \in \mathbb{A}$. Let $(\rho_1, \pi_1) = \mu_1(d)$, $(\rho_2, \pi_2) = \mu_2(d)$, $i = (d \stackrel{?}{=} d_1)$, $j = (d \stackrel{?}{=} d_2)$. Then $\vec{\mathbf{m}}(a, i, \pi_1, a, j, \pi_2)$ holds iff either*

$$i = 1, j = 0 \text{ and } \pi_1 = \pi_2 + a, \quad (4)$$

$$i = 0, j = 0 \text{ and } \pi_1 = \pi_2, \quad (5)$$

$$i = 1, j = 1 \text{ and } \pi_1 = \pi_2, \text{ or} \quad (6)$$

$$i = 0, j = 1 \text{ and } \pi_1 = \pi_2. \quad (7)$$

Analogously, $\overleftarrow{\mathbf{m}}(a, i, \pi_1, a, j, \pi_2)$ holds iff either

$$i = 1, j = 0 \text{ and } \rho_1 = \rho_2, \quad (8)$$

$$i = 0, j = 0 \text{ and } \rho_1 = \rho_2, \quad (9)$$

$$i = 1, j = 1 \text{ and } \rho_1 = \rho_2, \text{ or} \quad (10)$$

$$i = 0, j = 1 \text{ and } \rho_2 = \rho_1 + a. \quad (11)$$

Proof of Lemma 7 (monotonicity): Let $\mu_i = (a_i, d_i, \chi_i, \chi_i)$ for all $i \in \{1, 2, 3, 4\}$. In this context we will call *flexible value* to a value $d \in \mathbb{D} \setminus K$ and a *rigid* one to a value $d \in K$. Notice that $a_1 = a_3$ by definition of $\mu_1 \leq_K \mu_3$. By $\mu_1 \succrightarrow \mu_2$ we can assume $\mu_1(d_1) = (\rho_1, a_1\pi_1 + a_1)$ and

$\mu_2(d_1) = (a_2\rho_1, \pi_1)$ for some ρ_1 of the form $\rho_1 = a_1\rho' + a_1$. Similarly, we assume $\mu_2(d_2) = (a_2\rho_2 + a_2, \pi_2)$ and $\mu_1(d_2) = (\rho_2, a_1\pi_2)$ for some π_2 of the form $\pi_2 = a_2\pi' + a_2$. Hence, if $d_1 = d_2$ we have that $\rho_2 = \rho_1$ and $\pi_2 = \pi_1$.

We divide this proof into two parts. The first one groups all the easy cases, where the simulation can be done in one step: $\mu_3 \succrightarrow \mu_4$. For the second part, we will need to use more than one step to arrive to such μ_4 mosaic.

Easy situation. The easy situation is when all data values in μ_3 can be ‘simulated’ by the behavior of data values of μ_1 . We say that a flexible data value d of μ_3 can *simulate* another flexible data value d' of μ_1 if they have the same profile and one is equal to d_3 iff the other one is equal to d_1 . If this holds, d can have the same profile as d' by making a one-step transition, reading the same letter as μ_2 . If all flexible data values of μ_3 simulate some data value of μ_1 and in turn every flexible data value from μ_1 is simulated by at least one data value from μ_3 , we can arrive to a mosaic $\mu_4 \geq_K \mu_2$; but there is one caveat. There cannot be more than one data value that is chosen to simulate d_2 , since only one of them can become the data value of μ_4 . In this case the simulation would fail. But if there were another flexible data value d' with the same profile as d_2 in μ_1 , then we can choose only one data value to simulate d_2 , and all the rest to simulate d' . We call this last restriction *the condition B*. By definition of \leq_K it is a fact that for every data value d such that $\mu_3(d) \neq \mu_3(d_3) (= \mu_1(d_1))$ there is at least one data value in μ_1 that can simulate. Then, in order to assure that every data value of μ_3 can be simulated, it suffices to check that either d_1 is rigid (and hence also d_3), or if d_1 is flexible, that there are at least another flexible data value with the same profile as d_1 . We call this *the condition A*. If this would not hold, there could happen that there is only d_1 with the profile $\mu_1(d_1)$ in μ_1 , and that there is a flexible data value $d \neq d_3$ with $\mu_3(d) = \mu_1(d_1)$. This data value d cannot simulate any value, and the simulation would then fail.

We formalize the properties we have discussed.

- (A) If d_1 is a flexible value, then there are other flexible values with the same profile. That is, there exists $d' \neq d_1$ with $\mu_1(d') = \mu_1(d_1)$.
- (B) The same applies to d_2 : if it is a flexible value then there exists another flexible value $d' \neq d_2$ with $\mu_1(d') = \mu_1(d_2)$.

If A and B hold, by definition of \leq_K then there exists a surjective function $\xi : \mathbb{D} \rightarrow \mathbb{D}$ such that for every $d \in \mathbb{D}$,

1. $d = d_3$ iff $\xi(d) = d_1$,
2. $|\xi^{-1}(d_2)| = 1$.
3. $\mu_3(d) = \mu_1(\xi(d))$
4. if $d \in K$ then $\xi(d) = d$,

In other words, ξ behaves as the identity on K , sends data values from μ_3 to data values with the same profile in μ_1 , and sends only one data value to d_1 (namely d_3) and to d_2 . We call ξ a *simulating function* between μ_3 and μ_1 . If ξ is defined only for subsets $D_1, D_2 \subseteq \mathbb{D}$, we call $\xi : D_1 \rightarrow D_2$ a *partial simulating function* on D_1, D_2 between μ_3 and μ_1 , and in this case conditions 1 and 2 are replaced by

- 1'. if $d_3 \in D_1$, $d = d_3$ iff $\xi(d) = d_1$,
 2'. if $d_2 \in D_2$, $|\xi^{-1}(d_2)| = 1$.

We build $\mu_4 = (a_4, d_4, \overleftarrow{\chi}, \overrightarrow{\chi})$ where $a_4 = a_2$, $\{d_4\} = \xi^{-1}(d_2)$ and we define $\mu_4(d) = \mu_2(\xi(d))$ for all $d \in \mathbb{D}$. It is straightforward to check that for every $d \in \mathbb{D}$

$$\mathbf{m}(a_1, d \stackrel{?}{=} d_3, \mu_3(d), a_2, d \stackrel{?}{=} d_4, \mu_4(d)) \text{ iff } \\ \mathbf{m}(a_1, \xi(d) \stackrel{?}{=} d_1, \mu_1(\xi(d)), a_2, \xi(d) \stackrel{?}{=} d_2, \mu_2(\xi(d))) .$$

Complex situation 1. We suppose first that the condition A just seen does not hold, but that condition B holds. In other words, d_1 is the only flexible value with profile $\mu_1(d_1)$ in μ_1 . In particular, $d_1 \neq d_2$. Otherwise we would have a contradiction, since A and B would express the same property, and we assume that one holds but the other does not.

Note that in this case we cannot define a function ξ just as before, because there may be a flexible data value $d \neq d_3$ with $\mu_3(d) = \mu_1(d_1)$. (If there is no such data value, we proceed just as in the easy case.) We basically do not know what to do with this data value d , since it cannot simulate any other data value from μ_1 . Let P be the set of these problematic data values,

$$P = \{d \in \mathbb{D} \setminus K \mid d \neq d_3, \mu_3(d) = \mu_1(d_1)\}$$

If $P = \{e_1, \dots, e_n\}$, we will create n intermediary mosaics μ'_1, \dots, μ'_n with

$$\mu_3 \succ \mu'_1 \succ \dots \succ \mu'_n \succ \mu_4 \geq_K \mu_2$$

where μ'_i has the same letter a_1 as μ_3 and e_i as data value. Further, all data values except $\{d_3, e_1, \dots, e_n\}$ will preserve the same profile they have in μ_3 all along μ'_1, \dots, μ'_n . This is just a consequence of (5). Each mosaic μ'_i is defined as follows

- a_1 is the label, and e_i is the data value,
- d_3, e_1, \dots, e_{i-1} have profile $(\rho_1, a_1\pi_1)$,
- e_i, \dots, e_n have profile $(\rho_1, a_1\pi_1 + a_1)$, and
- any other data value $d \notin \{d_3, e_1, \dots, e_n\}$ has profile $\mu_3(d)$.

These are indeed legal mosaics.

Claim 1. For all $i \in [n]$, μ'_i is valid and K -compliant.

Claim 2. $\mu_3 \succ \mu'_1$, and for all $i \in [n-1]$, $\mu'_i \succ \mu'_{i+1}$.

Proof: To simplify notation, let us define $\mu'_0 = \mu_3$ and $e_0 = d_3$. Take an arbitrary data value $d \notin \{e_0, e_1, \dots, e_n\}$ and take any pair of mosaics μ'_i, μ'_{i+1} with $0 \leq i < n$. Then we have that: d is not the data value of μ'_i nor μ'_{i+1} ; and μ'_i and μ'_{i+1} carry the same letter a_1 . We can thus apply conditions (5) and (9) and obtain that $\mu_i(d) = \mu'_{i+1}(d)$ verifies \mathbf{m} . If $d \in \{e_0, \dots, e_{i-1}, e_{i+2}, \dots, e_n\}$, the same applies, since d is not the data value of μ'_i nor μ'_{i+1} , and then $\mu'_i(d) = \mu'_{i+1}(d)$ verifies \mathbf{m} . If $d = e_{i+1}$ we have that $\mu'_i(e_{i+1}) = \mu'_{i+1}(e_{i+1}) = (\rho_1, a_1\pi_1 + a_1)$. $a_1\pi_1 + a_1$ verifies $\overleftarrow{\mathbf{m}}$ by condition (7), and ρ_1 verifies $\overrightarrow{\mathbf{m}}$ by condition (11), since $\rho_1 = a_1\rho' + a_1$ and then $\rho_1 + a_1 = \rho_1$ by idempotence of $+a_1$. Finally, if $d = e_i$ we

have that $\mu'_i(d) = (\rho_1, a_1\pi_1 + a_1)$ and $\mu'_{i+1}(d) = (\rho_1, a_1\pi_1)$. The fact that ρ_1 is preserved is verified by $\overleftarrow{\mathbf{m}}$ by condition (8). And on the other hand, by (4) $a_1\pi_1 + a_1$ and $a_1\pi_1$ verify $\overrightarrow{\mathbf{m}}$. ■

Let $D_1 = \mathbb{D} \setminus \{d_3, e_1, \dots, e_n\}$, $D_2 = \mathbb{D} \setminus \{d_1\}$ and $\xi : D_1 \rightarrow D_2$ be a partial simulating function on D_1, D_2 between μ_3 and μ_1 . It must exist since condition B holds.

Claim 3. There exists a partial simulating function $\xi : D_1 \rightarrow D_2$ between μ_3 and μ_1 .

We then define μ_4 as $\mu_4 = (a_4, d_4, \overleftarrow{\chi}_4, \overrightarrow{\chi}_4)$ with $a_4 = a_2$, $\{d_4\} = \xi^{-1}(d_2)$, $\mu_4(d) = (a_2\rho_1, \pi_1)$ for all $d \in \{d_3, e_1, \dots, e_n\}$ and $\mu_4(d) = \mu_2(\xi(d))$ for all other $d \in \mathbb{D}$. It then follows that $\mu'_n \succ \mu_4$ and that $\mu_4 \geq_K \mu_2$.

Claim 4. $\mu'_n \succ \mu_4$

Proof: For every data value d that is simulated by ξ we have that $\mathbf{m}(a_1, d \stackrel{?}{=} e_n, \mu'_n(d), a_2, d \stackrel{?}{=} d_4, \mu_4(d))$. Hence, since ξ does all the work for us for all data values except $\{d_3, e_1, \dots, e_n\}$, we concentrate on these values. First, take any $e \in \{d_3, e_1, \dots, e_{n-1}\}$. We have that $\mu'_n(e) = (\rho_1, a_1\pi_1)$ and $\mu_4(e) = (a_2\rho_1, \pi_1)$. It is immediate that $\overrightarrow{\mathbf{m}}(a_1, 0, a_1\pi_1, a_2, 0, \pi_1)$ and $\overleftarrow{\mathbf{m}}(a_1, 0, \rho_1, a_2, 0, a_2\rho_1)$ and hence that e behaves correctly. Now take e_n , where we have $\mu'_n(e_n) = (\rho_1, a_1\pi_1 + a_1)$ and $\mu_4(e_n) = (a_2\rho_1, \pi_1)$. We simply need then to verify $\overrightarrow{\mathbf{m}}(a_1, 1, a_1\pi_1 + a_1, a_2, 0, \pi_1)$ and $\overleftarrow{\mathbf{m}}(a_1, 1, \rho_1, a_2, 0, a_2\rho_1)$, which are clearly true. ■

Claim 5. $\mu_4 \geq_K \mu_2$

There are two other complex situations: when neither conditions A nor B hold, and when condition A holds but B does not. These conditions are treated in an analogous way. Indeed, if condition A holds but B does not, we use a symmetrical strategy, creating intermediate mosaics but this time with the same label as μ_2 . If neither of the conditions hold, we combine both approaches. The details are in the Appendix. ■

We can now show that the reduction from the derivation problem for \succ to the derivation problem for \rightarrow_K is sound and complete by the following two lemmas.

Lemma 9 (complete). For every K -compliant complete \rightarrow sequence $\mu_1 \succ \dots \succ \mu_n$ over $\mathcal{M}(\varphi)$ there exists a K -compliant complete sequence $\mu'_1 \rightarrow_K \dots \rightarrow_K \mu'_n$ over $\mathcal{M}(\varphi)$.

Proof: Trivial, since \succ over K -compliant mosaics is a particular case of \rightarrow_K by reflexivity of \leq_K . ■

Lemma 10 (sound). For every K -compliant sequence $\mu_1 \rightarrow_K \dots \rightarrow_K \mu_n$ over $\mathcal{M}(\varphi)$ there exists a K -compliant \succ sequence $\mu'_1 \succ \dots \succ \mu'_m$ over $\mathcal{M}(\varphi)$ with $\mu_1 \leq_K \mu'_1$ and $\mu_n \leq_K \mu'_m$.

Proof: We proceed by induction on n . The base case is immediate. For the inductive step, suppose we have a sequence $\mu_1 \rightarrow_K \dots \rightarrow_K \mu_n \rightarrow_K \mu_{n+1}$. By inductive hypothesis we obtain $\nu_1 \succ \dots \succ \nu_m$ with $\nu_1 \geq_K \mu_1$ and $\nu_m \geq_K \mu_n$, as

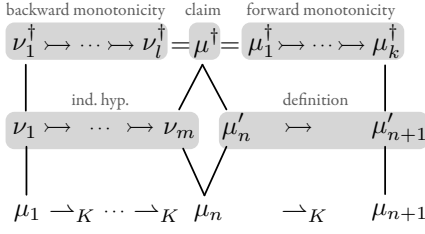


Fig. 3. General argument of Lemma 10. Solid lines correspond to the relation \leq_K , where the upper element in the diagram is the bigger.

shown in Figure 3. Since $\mu_n \rightarrow_K \mu_{n+1}$, then

$$\mu_n \leq_K \mu_n' \rightarrow \mu_{n+1}' \geq_K \mu_{n+1}$$

for some μ_n' and μ_{n+1}' by definition of \rightarrow_K . Let μ^\dagger be such that $\mu^\dagger \geq_K \nu_m$ and $\mu^\dagger \geq_K \mu_n'$. It always exists such μ^\dagger .

Claim 6. For any three mosaics μ_1, μ_2, μ_3 with $\mu_1 \geq_K \mu_3$ and $\mu_2 \geq_K \mu_3$, there exists a mosaic ν such that $\mu_1 \leq_K \nu$ and $\mu_2 \leq_K \nu$.

By applying the monotonicity lemma we obtain on the one hand a sequence

$$\nu_1^\dagger \rightarrow \dots \rightarrow \nu_l^\dagger$$

where $\nu_l^\dagger = \mu^\dagger$ and $\nu_1 \leq_K \nu_1^\dagger$, and by the backwards monotonicity lemma that

$$\mu_1^\dagger \rightarrow \dots \rightarrow \mu_k^\dagger$$

where $\mu_1^\dagger = \mu^\dagger$ and $\mu_k^\dagger \geq_K \mu_{n+1}'$. This is depicted in Figure 3. Further, these \rightarrow sequences are K -compliant. Then,

$$\nu_1^\dagger, \dots, \nu_{l-1}^\dagger, \mu_1^\dagger, \mu_2^\dagger, \dots, \mu_k^\dagger$$

is the the desired \rightarrow sequence. \blacksquare

By Lemmas 9 and 10 the Proposition 2 follows. The completeness and the fact that the first mosaic verifies φ_1 follow from Lemma 5.

V. THE DERIVATION PROBLEM FOR \rightarrow_K

In this Section we solve the derivation problem for \rightarrow_K .

Since the transition system \rightarrow_K is parametrized by K , in fact the derivation problem for \rightarrow_K is in fact: “is there a finite set $K \subseteq D$ such that the derivation problem for \rightarrow_K over $\mathcal{M}(\varphi)$ has a positive solution?”.

Note that we are dealing with many sources of infinity here: the choice of K , the domain of the transition system (i.e., the set of all mosaics), and the definition of \rightarrow_K , which involves making a potentially infinite number of tests. Nevertheless, they can all be tamed. This is the objective of this section.

We now show that the object that we are dealing with can be effectively bounded, to end up with finite objects that can be treated by an algorithm.

A. A bound on the rigid values

We will show that for every complete \rightarrow sequence over $\mathcal{M}(\varphi)$ there are only linearly bounded many rigid values. This Lemma will allow us to bound the number of minimal \leq_K classes of equivalence.

Lemma 11. For every complete sequence $\mu_1 \rightarrow \dots \rightarrow \mu_n$ over $\mathcal{M}(\varphi)$ there is a set $K \subseteq \mathbb{D}$ such that:

1. $|K| \leq 2 \cdot |\Omega_\varphi^p|$
2. μ_1, \dots, μ_n are K -compliant.

Corollary 1. If there is a finite set $K' \subseteq \mathbb{D}$ such that the derivation problem for $\rightarrow_{K'}$ has a solution, then the derivation problem for \rightarrow_K has also a solution, where K is any fixed subset of \mathbb{D} of $2 \cdot |\varphi|$ elements.

Thus, henceforth we assume without any loss of generality that K is any fixed subset of \mathbb{D} of $2 \cdot |\Omega_\varphi^p|$ elements.

B. A bound on the domain

Let us write $\text{MIN}_{\mathcal{M}(\varphi)}$ for the set of mosaics of $\mathcal{M}(\varphi)$ that are minimal with respect to \leq_K . This set contains only one element (arbitrarily chosen) for every minimal class of equivalence of \leq_K .

Remark 2. Notice that if $\mu_1 \rightarrow_K \mu_2$ then $\mu_1' \rightarrow_K \mu_2'$ for any $\mu_1' \leq_K \mu_1$ and $\mu_2' \leq_K \mu_2$. Hence, we can restrict the domain of \rightarrow_K to $\text{MIN}_{\mathcal{M}(\varphi)}$. That is, to mosaics that are minimal with respect to \leq_K .

Lemma 12. The number of elements of $\text{MIN}_{\mathcal{M}(\varphi)}$ is double exponential in $|\varphi|$.

From the previous statements we obtain the following.

Lemma 13. If there is a complete K -compliant \rightarrow_K sequence $\mu_1 \rightarrow_K \dots \rightarrow_K \mu_n$ over $\mathcal{M}(\varphi)$, then there is a complete K -compliant sequence $\mu_1' \rightarrow_K \dots \rightarrow_K \mu_m'$ where

- $\mu_1' \leq_K \mu_1$, $\mu_m' \leq_K \mu_n$,
- $m \leq |\text{MIN}_{\mathcal{M}(\varphi)}|$, and
- $\mu_i' \in \text{MIN}_{\mathcal{M}(\varphi)}$ for all $1 \leq i \leq m$.

C. A bound on the relation

The lemmas above show that we have that the problem now reduces to find a complete small \rightarrow_K sequence for a small K . In fact, we can test the relation \rightarrow_K between any two elements in polynomial space.

We define the space of a mosaic μ that we denote it by $|\mu|$ as the space needed to write the multiset of all non-empty profiles of the mosaic. It is not hard to see that \rightarrow_K can be in fact checked in polynomial space.

Lemma 14. Given two K -compliant mosaics μ_1, μ_2 , we can test $\mu_1 \rightarrow_K \mu_2$ in space polynomial in $|\mu_1|$ and $|\mu_2|$.

D. The algorithm

With all these ingredients we can now build the algorithm.

Proposition 3. The derivation problem is in EXPSpace.

VI. LOWER BOUND

There is also a lower bound of EXPSPACE for the satisfiability problem of PathLogic. This bound also translates into a EXPSPACE-hardness for XPath($*\leftarrow, \rightarrow^*, =$) since there is a PTIME translation from PathLogic to XPath($*\leftarrow, \rightarrow^*, =$), as shown in the next Section.

Proposition 4. *The satisfiability problem for PathLogic is EXPSPACE-hard.*

This can be shown by coding a solution for an instance of the 2^n corridor tiling problem. The coding, although it is not quite straightforward, utilizes several coding techniques developed in [FS09]. In particular, the coding of [FS09] was introduced to show non-primitive recursiveness of the logic XPath($\rightarrow^+, =$), whereas here we work with a logic which is expressive-equivalent to XPath($\rightarrow^*, =$). Although the precise coding of [FS09] fails for XPath($\rightarrow^*, =$), we reuse the same kind of techniques. The main idea is in the Appendix.

Since the coding for this problem need only to use expressions in one direction, we also obtain the EXPSPACE-hardness for XPath($\rightarrow^*, =$).

VII. COROLLARIES ON XPATH AND XML DOCUMENTS

Let us call *plain*-XPath($*\leftarrow, \rightarrow^*, =$) to the fragment of XPath such that that

- all path expressions do not mix $*\leftarrow$ and \rightarrow^* , i.e., each path expressions uses only $*\leftarrow$ or only \rightarrow^* , and
- all test node expressions $\langle \alpha = \beta \rangle$ and $\langle \alpha \neq \beta \rangle$ are such that α uses only $*\leftarrow$ and β only \rightarrow^* .

It is immediate that this basically corresponds to PathLogic.

Lemma 15. *There are PTIME translations such that for every node expression of plain-XPath($*\leftarrow, \rightarrow^*, =$) (resp. of PathLogic) return an equivalent node expression PathLogic (resp. plain-XPath($*\leftarrow, \rightarrow^*, =$)).*

Corollary 2. *The satisfiability problem for plain-XPath($*\leftarrow, \rightarrow^*, =$) on data words is EXPSPACE-complete.*

In fact, plain-XPath($*\leftarrow, \rightarrow^*, =$) (and hence PathLogic) is expressive-equivalent to XPath($*\leftarrow, \rightarrow^*, =$) through an EXP-TIME translation.

Lemma 16. *There exists a computable EXPTIME translation that for every node expression φ of XPath($*\leftarrow, \rightarrow^*, =$) returns an equivalent node expression ψ of plain-XPath($*\leftarrow, \rightarrow^*, =$).*

VIII. DISCUSSION

We believe that the proof we give here allows to extend existing decidability results of XPath fragments on XML with navigation through reflexive-transitive axes. However, obtaining such results is non-trivial, since one has to combine the techniques of such results with the ones contained in our main proof. We propose that the following are candidates of (non-trivial) corollaries of the work contained in this paper and known results on XPath.

Conjecture 1. *Satisfiability of XPath($\downarrow, \downarrow^*, *\leftarrow, \rightarrow^*, =$) on XML documents is decidable, with elementary complexity.*

This would be an extension of the result of [Fig09] stating that XPath($\downarrow, \downarrow^*, =$) is EXPTIME-complete, and in contrast to the fact that XPath($\downarrow, \downarrow^*, *\leftarrow, \rightarrow^+, =$) is undecidable and XPath($\downarrow, \downarrow^*, \rightarrow^+, =$) has non-primitive recursive complexity [FS09]. In fact, we believe that this can be pushed to show the following.

Conjecture 2. *Satisfiability of XPath($\downarrow, \downarrow^*, \uparrow^*, *\leftarrow, \rightarrow^*, =$) on XML documents is decidable, with elementary complexity.*

Furthermore, it is plausible that the following fragment of XPath is decidable.

Conjecture 3. *Satisfiability of XPath($\downarrow, \downarrow^*, \uparrow, \uparrow^*, *\leftarrow, \rightarrow^*, =$) on XML documents is decidable.*

This would be an extension of a recent result stating that XPath($\downarrow, \downarrow^*, \uparrow, \uparrow^*, =$) is decidable [FS11]. However, even XPath($\downarrow, \downarrow^*, \uparrow, \uparrow^*, \rightarrow^+, =$) or XPath($\downarrow, \downarrow^*, \uparrow, \uparrow^*, \rightarrow, =$) are undecidable [FS09].

Future work

We conclude by mentioning some open questions.

1. *Decidability of XPath fragments.* Is any of the aforementioned conjectures true?
2. *Infinite words.* Is PathLogic decidable on infinite data words? Observe that PathLogic does not enjoy the finite model property.
3. *Linear order.* As a possible extension, one may add comparisons between data values with respect to a linear order.

REFERENCES

- [BMS⁺06] Mikolaj Bojańczyk, Anca Muscholl, Thomas Schwentick, Luc Segoufin, and Claire David. Two-variable logic on words with data. In *Annual IEEE Symposium on Logic in Computer Science (LICS'06)*, pages 7–16. IEEE Computer Society Press, 2006.
- [CD99] James Clark and Steve DeRose. XML path language (XPath). Website, 1999. W3C Recommendation. <http://www.w3.org/TR/xpath>.
- [Ch186] Bogdan S. Chlebus. Domino-tiling games. *J. Comput. Syst. Sci.*, 32(3):374–392, 1986.
- [DL09] Stéphane Demri and Ranko Lazić. LTL with the freeze quantifier and register automata. *ACM Transactions on Computational Logic*, 10(3), 2009.
- [Fig09] 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.
- [FS09] 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.
- [FS11] 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. Leibniz-Zentrum für Informatik, 2011. To appear.
- [KSZ10] Ahmet Kara, Thomas Schwentick, and Thomas Zeume. Temporal logics on words with multiple data values. In *IARCS Annual Conference on Foundations of Software Technology and Theoretical Computer Science (FSTTCS'10)*, 2010.
- [SZ10] Thomas Schwentick and Thomas Zeume. Two-variable logic with two order relations. In *EACSL Annual Conference on Computer Science Logic (CSL'10)*, 2010.

APPENDIX

APPENDIX TO SECTION III

Proof of Proposition 1: We want to show that a PathLogic formula $\varphi = \varphi_1 \wedge \varphi_2$ in normal form is satisfiable iff there exists a complete \rightarrow sequence of mosaics

$$\mu_1 \rightarrow \mu_2 \rightarrow \dots \rightarrow \mu_n$$

over $\mathcal{M}(\varphi)$ such that $\mu_1 \vdash \varphi_1$. Throughout the proof we will only work with the $\vec{\chi}$ part of the mosaic, since the situation for $\overleftarrow{\chi}$ is symmetrical.

(if) Now, we show the other direction. Suppose we have a complete \rightarrow sequence

$$\mu_1 \rightarrow \dots \rightarrow \mu_n$$

over $\mathcal{M}(\varphi)$ such that $\mu_1 \vdash \varphi_1$. We build a data word \mathbf{w} of length n such that μ_i is the φ -abstraction of (\mathbf{w}, i) . The data word is defined as $\mathbf{w}(i) = (a_i, d_i)$ where (a_i, d_i) is the label and data value of μ_i . We claim that for every i , μ_i is the φ -abstraction of (\mathbf{w}, i) . In the following proof we only work with the $\vec{\chi}$ part of the mosaics. The proof for $\overleftarrow{\chi}$ will follow by symmetry.

(sound) If from a position i we can access a data value d by going to the right with a path expression $\alpha \in \Omega_\varphi^p$, we want to show that $(d, \alpha) \in \vec{\chi}$. We proceed by combined induction on the size of α and the index $n - i$. The base cases for $i = n$ are immediate since everything is satisfied locally. If $\alpha = [\psi]$ and $d_i = d$, then the data value can be accessed locally. By item 2 of validity of μ_i , $\pi_i = a_i \pi_i + a_i$, where $\pi_i = \vec{\chi}_i(d) \ni \alpha$. Since $\alpha \in \sigma_{a_i}$ then $\alpha \in \pi_i$. If on the other hand $d_i \neq d$, this means that there is a future position $j > i$ with data value $d_j = d$ such that $a_j \models \psi$. This means that d can be accessed through α from position $i + 1$. We then apply the inductive hypothesis and we obtain that $(d, \alpha) \in \vec{\chi}_{i+1}$. By Eq. (2) of $\vec{\mathbf{m}}$, $\pi_i = a_i \pi_{i+1}$ where $\pi_{i+1} = \vec{\chi}_{i+1}(d)$. Since $\alpha \in \pi_{i+1}$, then $\alpha \in \pi_i$. This concludes the base case. For the inductive case, suppose $\alpha = [\psi] \cdot \alpha'$. Suppose first that $a_i \models \psi$. Observe that if we can access d through α we can also access d through α' . By inductive hypothesis we have that $\alpha' \in \vec{\chi}_i(d)$. By item 3 of validity of μ_i we have that $\vec{\chi}_i(d) = a_i \pi$ for some π . since $[\psi] \in \sigma_{a_i}$, then $[\psi] \cdot \alpha' \in \vec{\chi}_i(d)$. Suppose now that $a_i \not\models \psi$. This means that d can be accessed actually from position $i + 1$ through α . We then apply the inductive hypothesis and obtain that $\alpha \in \vec{\chi}_{i+1}(d)$. By Eq. (1) or (2) (depending on whether $d = d_i$ or not) we obtain in both cases that α must also be in $\vec{\chi}_i(d)$.

(complete) Suppose that $(d, \alpha) \in \vec{\chi}_i$, we want to show that from position i we can access d by going to the right with a path expression α . We proceed first by induction on the length of α and then by induction on the index $n - i$. Suppose first that $\alpha = [\psi]$. If $a_i \models \psi$ and $d_i = d$, then α is satisfied locally

and we are done. If $d_i \neq d$ then μ_i is not right-complete, this means that there is a mosaic μ_{i+1} where

$$\vec{\mathbf{m}}(a_i, 0, \pi_i, a_{i+1}, d_{i+1} \stackrel{?}{=} d, \pi_{i+1})$$

where $\pi_i = \vec{\chi}_i(d) \ni \alpha$ and $\pi_{i+1} = \vec{\chi}_{i+1}(d)$. By Eq. (2) of $\vec{\mathbf{m}}$, $\pi_i = a_i \pi_{i+1}$. Since α has only one element, by definition of $a_i \pi_{i+1}$, $\alpha \in \pi_2$. We can then apply the inductive hypothesis and obtain that d can be accessed through α from position $(\mathbf{w}, i + 1)$, which implies that it can also be accessed from position (\mathbf{w}, i) .

Finally, if $d_i = d$ but $a_i \not\models \psi$ then μ_i is not right-complete and there must be a mosaic μ_{i+1} where

$$\vec{\mathbf{m}}(a_i, 1, \pi_i, a_{i+1}, d_{i+1} \stackrel{?}{=} d, \pi_{i+1})$$

where $\pi_i = \vec{\chi}_i(d) \ni \alpha$ and $\pi_{i+1} = \vec{\chi}_{i+1}(d)$. By Eq. (1) of $\vec{\mathbf{m}}$, $\pi_i = a_i \pi_{i+1} + a_i$. Since $a_i \not\models \psi$ we must have that $\alpha \in \pi_{i+1}$. Again, we apply the inductive hypothesis and obtain that d can be accessed through α from position $(\mathbf{w}, i + 1)$, which implies that it can also be accessed from position (\mathbf{w}, i) . This concludes the base case. For the inductive case, suppose $\alpha = [\psi] \cdot \alpha'$. If $a_i \models \psi$ then by condition 1 of a valid mosaic, since $\alpha \in \vec{\chi}_i(d)$ then $\alpha' \in \vec{\chi}_i(d)$. By applying inductive hypothesis we obtain that d can be reached from position i through α' and since $a_i \models \psi$ it can also be reached through $[\psi] \cdot \alpha'$. If on the other hand $a_i \not\models \psi$, then by equations (1) or (2) (depending on whether $d_i = d$ or not) we obtain that $\pi_i = a_i \pi_{i+1}$ or $\pi_i = a_i \pi_{i+1} + a_i$. In any case, since $\psi \notin a_i$, we obtain that $\alpha \in \pi_{i+1}$ and we apply the inductive hypothesis just as in the base case.

(only if) This direction is easier. Given a data words we have to show that the φ -abstraction of each position results in a complete \rightarrow sequence whose first element verifies φ_1 .

Given a data word $\mathbf{w} = \mathbf{a} \otimes \mathbf{d}$ such that $\mathbf{w} \models \varphi$ we can build a \rightarrow sequence $\mu_1, \dots, \mu_{|\mathbf{w}|}$ over $\mathcal{M}(\varphi)$ that is complete, such that $\mu_1 \vdash \varphi$. We define $\mu_i = (a_i, d_i, \overleftarrow{\chi}_i, \vec{\chi}_i)$ as the φ -abstraction of (\mathbf{w}, i) .

Let us first show that $\mu_1, \dots, \mu_{|\mathbf{w}|}$ are valid mosaics of $\mathcal{M}(\varphi)$. It is easy to see that the path expressions of $\vec{\chi}(d)$ are suffix closed since if $i \in \llbracket \alpha \cdot \beta \rrbracket_{\vec{\chi}}^{\mathbf{w}}$ then $i \in \llbracket \beta \rrbracket_{\vec{\chi}}^{\mathbf{w}}$. Hence, item 1 holds. Items 2 and 3 are also very easy to verify. Finally, we need to show that all mosaics with a letter a verify $\gamma_\varphi(a)$. But this must be true by the formula φ_2 holding at the leftmost element. We then have that item 4 of the validity conditions is also met.

Secondly, we show that $\mu_1, \dots, \mu_{|\mathbf{w}|}$ is a \rightarrow sequence. Let us analyze a pair of arbitrary positions $i, i + 1$ and a data value $d \in \mathbb{D}$. Let us show that

$$\vec{\mathbf{m}}(a_i, d_i \stackrel{?}{=} d, \vec{\chi}_i(d), a_{i+1}, d_{i+1} \stackrel{?}{=} d, \vec{\chi}_{i+1}(d))$$

holds true. Take an arbitrary path expression $\alpha = [\psi_1] \cdot \dots \cdot [\psi_n] \in \Omega_\varphi^p$. If $(d, \alpha) \in \vec{\chi}_{i+1}$ then there are n ordered positions greater or equal to $i + 1$ satisfying respectively ψ_1, \dots, ψ_n by definition of φ -abstraction. Then in particular they are greater or equal to i , and hence $(d, \alpha) \in \vec{\chi}_i$. Thus,

$\vec{\chi}_{i+1}(d) \subseteq \vec{\chi}_i(d)$. On the other hand, suppose there is a prefix of α , say $[\psi_1] \cdots [\psi_k]$, such that $a_i \models \psi_j$ for all $j \in [1..k]$, and that $(d, [\psi_{k+1}] \cdots [\psi_n]) \in \vec{\chi}_{i+1}(d)$. Then, there are witness positions for $\psi_{k+1}, \dots, \psi_n$ which are greater or equal to $i+1$. Also, there is a witness for ψ_1, \dots, ψ_k at position i . Then, d can be accessed through α from position i , and then $\alpha \in \vec{\chi}_i(d)$ by definition of φ -abstraction. We just showed that $a_i \vec{\chi}_{i+1}(d) \subseteq \vec{\chi}_i(d)$, that will be useful next. We separate the two different cases from the definition of \vec{m} : either $d_i = d$ or $d_i \neq d$. Suppose first $d_i \neq d$. Any path $\alpha \in \vec{\chi}_i(d)$ has a (possibly empty) prefix which is satisfied at position i by letter a_i , and a suffix that is satisfied at positions bigger or equal to $i+1$. But this suffix cannot be empty, since this would mean that $d_i = d$. Then, we have that $\vec{\chi}_i(d) \subseteq a_i \vec{\chi}_{i+1}(d)$ and we conclude that $a_i \vec{\chi}_{i+1}(d) = \vec{\chi}_i(d)$. On the other hand, if $d_i = d$, we can show that in this case $a_i \vec{\chi}_{i+1} + a_i = \vec{\chi}_i$ proceeding in a similar way. We showed that $\vec{m}(a_i, d_i \stackrel{?}{=} d, \vec{\chi}_i, a_{i+1}, d_{i+1} \stackrel{?}{=} d, \vec{\chi}_{i+1})$ holds. The proof for $\vec{m}(a_i, d_i \stackrel{?}{=} d, \vec{\chi}_i, a_{i+1}, d_{i+1} \stackrel{?}{=} d, \vec{\chi}_{i+1})$ is completely analogous. We conclude that the sequence $\mu_1, \dots, \mu_{|\mathbf{w}|}$ is a \rightarrow sequence over $\mathcal{M}(\varphi)$.

Now, we show that it is *complete*. By way of contradiction, suppose that μ_1 is not left-complete. Let $(d, \pi) \in \overleftarrow{\text{dem}}(\mu_1)$. Then:

- If $d = d_1$ and $\pi \neq \sigma_a$, this means that there is $\alpha \in \Omega_\varphi^p$ such that $(d, \alpha) \in \overleftarrow{\chi}(\mu_1)$ and $a \not\models \psi$ where $[\psi]$ is a particle of α . By construction of μ_1 and the fact that $(d, \alpha) \in \overleftarrow{\chi}(\mu_1)$, this means that there is a position to the left of the first position that can reach the data value d through α . But since $a \not\models \psi$ this position must be *strictly* to the left. This is a contradiction, since there are no positions to the left of the first one.
- Otherwise, $d \neq d_1$ and $(d, \alpha) \in \overleftarrow{\chi}(\mu_1)$. Again, we conclude that there must be a position strictly to the left of the first one, which is a contradiction.

Then, we have that necessarily μ_1 is left-complete. Proceeding in a similar way, we prove that $\mu_{|\mathbf{w}|}$ is right-complete.

Finally, we show that $\mu_1 \vdash \varphi_1$. But this is immediate by Lemma 3 from the fact that $\mathbf{w}, 1 \models \mu_1$. ■

APPENDIX TO SECTION IV

Proof of Lemma 5: We show, one by one, all conditions of the Lemma.

Condition 1: If μ_1 is not right-complete, then either:

- There is a path expression $\alpha \in \sigma_{a_1}$ and a data value $e \neq d_1$ such that $(e, \alpha) \in \vec{\chi}_1$. Then for the profile $(\rho, \pi) = \mu_1(e)$ there must be at least one data value e' such that $\mu_2(e') = (\rho, \pi)$. Hence, μ_2 is not right-complete either.
- There is a path expression $\alpha \notin \sigma_{a_1}$ such that $(d_1, \alpha) \in \vec{\chi}_1$. Since $\mu_1(d_1) = \mu_2(d_2)$, then $(d_2, \alpha) \in \vec{\chi}_2$ and μ_2 is not right-complete.

- There is a path expression $\alpha \in \sigma_{a_1}$ such that $(d_1, \alpha) \notin \vec{\chi}_1$. But in this case μ_1 is not valid, since d_1 must contain at least all the path expressions from σ_{a_1} .

The back condition is identical, and the condition for left-complete is completely symmetrical.

Condition 2: Suppose that μ_2 is not K -compliant. Then suppose without any loss of generality that there is $\alpha \in \Omega_\varphi^p$ such that $\vec{\chi}_2(\alpha) = \{e\}$ for $e \notin K$. Let $\mu_2(e) = (\rho, \pi)$, and hence $\mu_2^{-1}(\rho, \pi) = \{e\}$. By $\mu_1 \leq_K \mu_2$, $|\mu_1^{-1}(\rho, \pi) \setminus K| \leq 1$ $|\mu_2^{-1}(\rho, \pi) \setminus K|$, and by definition of \leq_1 we have that $|\mu_2^{-1}(\rho, \pi) \setminus K| = 1$. Then there is also a data value e' such that $\vec{\chi}_1(\alpha) = \{e'\}$. Hence, μ_1 is not K -compliant.

Condition 4: Suppose that $\mu_1 \vdash \psi$. Let $\psi = \langle \overleftarrow{\alpha} = \overrightarrow{\beta} \rangle$ for some path expressions α, β . Suppose first that either α is equal to ε . Then the test amounts to having $(d_1, \beta) \in \vec{\chi}_1$. Since $\mu_1(d_1) = \mu_2(d_2)$ we have that $(d_2, \beta) \in \vec{\chi}_2$ and hence that $\mu_2 \vdash \psi$. If none of α, β are ε then there must be a data value $e \in \mathbb{D}$ such that $(e, \alpha) \in \overleftarrow{\chi}_1$ and $(e, \beta) \in \vec{\chi}_1$. Consider the profile of this data value $(\rho, \pi) = \mu_1(e)$. By $\mu_1 \leq_K \mu_2$ there must be a data value e' such that $\mu_2(e') = (\rho, \pi)$. Then, $\mu_2 \vdash \langle \overleftarrow{\alpha} = \overrightarrow{\beta} \rangle$. If $\psi = \langle \overleftarrow{\alpha} \neq \overrightarrow{\beta} \rangle$ we proceed in a similar way. If $\psi = \neg \langle \overleftarrow{\alpha} = \overrightarrow{\beta} \rangle$ this means that $\mu_1(\alpha) \cap \mu_1(\beta) = \emptyset$. By way of contradiction, suppose that $\mu_2(\alpha) \cap \mu_2(\beta) \ni e$. Let $(\rho, \pi) = \mu_2(e)$. Then we have that by $\mu_1 \leq_K \mu_2$ there must be a data value e' such that $\mu_1(e) = (\rho, \pi)$. Hence $\mu_1(\alpha) \cap \mu_1(\beta) \ni e'$ which is a contradiction. Finally, if $\psi = \neg \langle \overleftarrow{\alpha} \neq \overrightarrow{\beta} \rangle$, then either

- $\vec{\chi}_1(\alpha) = \emptyset$,
- $\vec{\chi}_1(\beta) = \emptyset$, or
- $\vec{\chi}_1(\alpha) = \vec{\chi}_1(\beta) = \{e\}$ for some $e \in \mathbb{D}$.

We can see that in each of this cases a similar property holds for μ_2 . In particular for the last one, this condition necessarily means that $e \in K$ since μ_1 is K -compliant. Let $(\rho, \pi) = \mu_1(e)$. Then we have that $|\mu_1^{-1}(\rho, \pi) \setminus K| = 0$ and hence $|\mu_2^{-1}(\rho, \pi) \setminus K| = 0$. This, added to the fact that $\mu_2(e) = (\rho, \pi)$ yields that $\vec{\chi}_2(\alpha) = \vec{\chi}_2(\beta) = \{e\}$.

Condition 3: This condition is immediate from Condition 4 and the fact that there is a data value e with $\mu_1(e) = (\rho, \pi)$ iff there is a data value e' with $\mu_2(e') = (\rho, \pi)$. ■

Proof of Lemma 8: We only prove conditions (4), (5), (6) and (7), since the second set of conditions is completely symmetrical.

If $i = 1, j = 0$, then $d = d_1, d \neq d_2$. By (1) $\pi_1 = a\pi_2 + a$. On the other hand, since $a_2 = a$, by condition 3 of validity applied to μ_2 , $\pi_2 = a\pi'_2$ for some π'_2 . By idempotence we have that $a\pi_2 = \pi_2$ and hence that $\pi_1 = a\pi_2 + a = \pi_2 + a$.

If $i = 0, j = 0$ by (2) $\pi_1 = a\pi_2$ and $a\pi_2 = \pi_2$ by the same reason as before.

If $i = 1, j = 1$, by (1) we have that $\pi_1 = a\pi_2 + a$ and $\pi_2 = a\pi'_2 + a$ by condition 2 of validity of μ_2 , since $d = d_2$. Again by idempotence we obtain $\pi_1 = a\pi_2 + a = \pi_2$.

Finally, if $i = 0, j = 1$, by (2) we have that $\pi_1 = a\pi_2$ and by condition by condition 2 of validity of μ_2 we have $\pi_2 = a\pi'_2 + a$. By idempotence, $\pi_1 = a\pi_2 + a = \pi_2$. ■

Proof of Claim 1: The only thing that changes between μ'_{i-1} and μ'_i (or between μ_3 and μ'_1) is that e_{i-1} is perhaps no longer reachable by some $\alpha \in \sigma_{a_1}$. There is however always at least one data value (namely e_i) that can reach α . In fact there are always two data values that can reach α except, perhaps, in μ'_n , where e_n may be the only (flexible) data value that can reach α .

So the only possibility to falsify the validity property 4 of a mosaic is that $\gamma_\varphi(a_1)$ demands that there exists at least two data values reachable by α , and that this is not true in μ'_n . Suppose ad absurdum that α is accessible by only one data value in μ'_n . First, note that α is accessible by at least two data values in μ_1 , namely d_1 and d' for some $d' \neq d_1$. If only d_1 could access α , then d_1 would be a rigid value, contradicting our assumption.

If $\mu_1(d') = (\rho_1, a_1\pi_1 + a_1)$ then d' must be rigid (because d_1 is the only flexible value with that profile), and $\mu'_n(d') = (\rho_1, a_1\pi_1 + a_1)$ and in this case there are two data values (e_n and d') that can access α .

If $\mu_1(d') \neq (\rho_1, a_1\pi_1 + a_1)$ then there is a data value d'' such that $\mu_3(d'') = \mu_1(d')$ by definition of $\mu_1 \leq_K \mu_3$ where, of course, $d'' \notin \{d_3, e_1, \dots, e_n\}$. Since μ'_n and μ_3 coincide in the profiles of all data values except perhaps in $\{d_3, e_1, \dots, e_n\}$ we obtain that $\mu'_n(d'') = \mu_3(d'') = \mu_1(d')$. Hence, there are two data values (namely e_n and d'') that can access α at μ'_n .

The properties 1, 2, and 3 of validity are immediate to check.

Finally, we prove by a similar argument that all μ'_i are K -compliant. Suppose ad absurdum that there is $\alpha \in \Omega_\varphi^p$ accessible only by one value that is flexible in μ'_i . Then it must be a data value from $\{e_1, \dots, e_n\}$ since the other values do not change of profile with respect to μ_3 , which is K -compliant. Further, it must be at μ'_n since otherwise α can be accessed both by e_{n-1} and e_n . As we already showed, if α is accessible by only one data value in μ'_n then it is also accessible by only one data value in μ_1 and it is then rigid, by K -compliance of μ_1 , which is in contradiction with our hypothesis. ■

Proof of Claim 3: For any $d \in D_1 \cap K$, we define $\xi(d) = d$. For any $d \in D_1 \setminus K$ such that $\mu_1(d) \neq \mu_1(d_2)$, we define $\xi(d) = d' \in D_2 \setminus K$ such that $\mu_1(d') = \mu_3(d)$, making sure that ξ is surjective onto all flexible values of $D_2 \setminus K$ that do not have profile $\mu_1(d_2)$, this can be done by definition of $\mu_1 \leq_K \mu_3$. If $d_2 \in K$, we are done with the definition of ξ . Otherwise, suppose $\mu_1(d_2) = (\rho, \pi)$ and $\mu_1^{-1}(\rho, \pi) \setminus K = \{d_2\} \cup D'_2$ such that $|D'_2| > 0$ by condition B. Let $d \in D_1$ be a flexible data value from D_1 with $\mu_3(d) = (\rho, \pi)$ (there must exist by \leq_K). We define $\xi(d) = d_2$, and for every other $d' \in D_1 \setminus \{d\}$ with the same profile we define $\xi(d') = d''$ for some $d'' \in D'_2$ (where there exists at least one element), again making sure that ξ is surjective onto D'_2 . We can see that ξ satisfies conditions 3, 4, and 2'. (condition 1' is trivially satisfied). ■

Proof of Claim 5: Note that condition 1 of \leq_K is immediate by definition of μ_4 . We now check condition 2. Let $d \in K$. Then we have that $\mu_4(d) = \mu_2(\xi(d)) = \mu_2(d)$ since $\xi(d) = d$. Condition 3 also holds: $\mu_4(d_4) = \mu_2(\xi(d_4)) = \mu_2(d_2)$.

To check condition 4, first consider the profile $(a_2\rho_1, \pi_1)$. $\mu_2^{-1}(a_2\rho_1, \pi_1)$ contains d_1 , perhaps some rigid data values, and perhaps some flexible data values d (with $\mu_1(d) \neq \mu_1(d_1)$, since condition A does not hold). We know that $\mu_4^{-1}(a_2\rho_1, \pi_1)$ contains d_3, e_1, \dots, e_n . So both sets contain at least one element. We show that for every data value in $\mu_2^{-1}(a_2\rho_1, \pi_1)$ there is a distinct data value in $\mu_4^{-1}(a_2\rho_1, \pi_1)$ and hence that $|\mu_2^{-1}(a_2\rho_1, \pi_1)| \leq |\mu_4^{-1}(a_2\rho_1, \pi_1)|$. For any $d \in \mu_2^{-1}(a_2\rho_1, \pi_1)$ with $d \neq d_1$ there exists a distinct data value $d' \notin \{d_3, e_1, \dots, e_n\}$ such that $\xi(d') = d$ and hence such that $d' \in \mu_4^{-1}(a_2\rho_1, \pi_1)$ by definition of μ_4 . On the other hand, for $d_1 \in \mu_2^{-1}(a_2\rho_1, \pi_1)$ there is $d_3 \in \mu_4^{-1}(a_2\rho_1, \pi_1)$. Thus, we conclude that $|\mu_2^{-1}(a_2\rho_1, \pi_1)| \leq |\mu_4^{-1}(a_2\rho_1, \pi_1)|$.

Now consider any profile (ρ, π) different from $(a_2\rho_1, \pi_1)$. By ξ , we have that for every data value $d \in \mu_2^{-1}(\rho, \pi)$ there exists d' with $\xi(d') = d$ such that $d' \in \mu_4^{-1}(\rho, \pi)$ by definition of μ_4 . Hence, $|\mu_2^{-1}(\rho, \pi)| \leq |\mu_4^{-1}(\rho, \pi)|$. Further, if there is some $d \in \mu_4^{-1}(\rho, \pi)$, then $\xi(d) \in \mu_2^{-1}(\rho, \pi)$. Then, we conclude that $|\mu_2^{-1}(\rho, \pi)| \leq |\mu_4^{-1}(\rho, \pi)|$, and that condition 4 holds. We thus have that $\mu_4 \geq_K \mu_2$. ■

Continuation of the Proof of Lemma 7 (monotonicity):

Complex situation 2. Now, suppose that neither condition A nor B hold. Then we can generate μ'_1, \dots, μ'_n just as before, but now our problem is that from μ'_n we cannot generate μ_4 as before since condition B fails. Let us call

$$R = \{e'_1, \dots, e'_m\} = \{d \in \mathbb{D} \setminus K \mid \mu_3(d) = \mu_1(d_2)\}.$$

These are the data values that can only be simulated by d_2 (otherwise B would hold). Notice that d_2 is necessarily flexible by the negation of condition B, and that R is non-empty by $\mu_1 \leq_K \mu_3$. We will also be using the data values P and mosaic μ'_i 's as defined for the first complex situation. Observe that if $d_1 = d_2$, then $\{d_3\} \cup P = R$, and otherwise $(\{d_3\} \cup P) \cap R = \emptyset$.

The strategy we will take is similar as the previous one, we will create intermediary mosaics μ''_1, \dots, μ''_m , but this time with a_2 as letter, to obtain

$$\mu_3 \succ \mu'_1 \succ \dots \succ \mu'_n \succ \mu''_1 \succ \dots \succ \mu''_m = \mu_4.$$

Each of the data values e'_i of R will change its profile first from $(\rho_2, a_1\pi_2) = \mu'_n(e'_i)$ to $(a_2\rho_2, \pi_2) = \mu''_1(e'_i)$, and then to $(a_2\rho_2 + a_2, \pi_2) = \mu''_m(e'_i)$. The only exception is the case of e'_1 , that will only change once its profile, directly to $(a_2\rho_2 + a_2, \pi_2)$ in μ''_1 . Consider a simulating function $\xi : D_1 \rightarrow D_2$ with $D_1 = \mathbb{D} \setminus \{d_3, e_1, \dots, e_n, e'_1, \dots, e'_m\}$, $D_2 = \mathbb{D} \setminus \{d_1, d_2\}$ between μ_3 and μ_1 . It is easy to verify that it must exist.

Claim 7. *There exists a partial simulating function $\xi : D_1 \rightarrow D_2$ between μ_3 and μ_1 .*

ξ is also a partial simulating function between μ'_n and μ_1 because it only simulates data values whose profiles were not changed with respect to μ_3 . We build μ''_1, \dots, μ''_m where each μ''_i

- has letter a_2 and data value e'_i
- e'_1, \dots, e'_i have profile $(a_2\rho_2 + a_2, \pi_2)$
- e'_{i+1}, \dots, e'_m have profile $(a_2\rho_2, \pi_2)$
- all data value from $(\{d_3\} \cup P) \setminus R$ (if any) have profile $(a_2\rho_1, \pi_1)$,
- any other data value $d \in D_1$ has profile $\mu_2(\xi(d))$.

Finally, the role of μ_4 is played by μ''_m . As before, we must verify that everything closes.

Claim 8. $\mu''_m \geq_K \mu_2$.

Proof: Note that condition 1 of \leq_K is immediate by definition of μ_4 . We now check condition 2. Let $d \in K$. Then we have that $\mu_4(d) = \mu_2(\xi(d)) = \mu_2(d)$ since $\xi(d) = d$. Condition 3 also holds: $\mu_4(d_4) = \mu''_m(e'_m) = (a_2\rho_2 + a_2, \pi_2) = \mu_2(d_2)$.

To check condition 4, first we proceed as in the proof of Claim 5 to show that $|\mu_2^{-1}(a_2\rho_1, \pi_1)| \leq_1 |\mu_4^{-1}(a_2\rho_1, \pi_1)|$.

For $(a_2\rho_2 + a_2, \pi_2)$ note that $|\mu_2^{-1}(a_2\rho_2 + a_2, \pi_2)| \leq |\mu_4^{-1}(a_2\rho_2 + a_2, \pi_2)|$ for each element in $\mu_4^{-1}(a_2\rho_2 + a_2, \pi_2)$ there is a distinct element in $\mu_2^{-1}(a_2\rho_2 + a_2, \pi_2)$. Indeed, for d_2 we have d_4 and for any other data value $d \in \mu_4^{-1}(a_2\rho_2 + a_2, \pi_2)$ we can take any data value from $\xi^{-1}(d)$, which has the same profile. Moreover, $|\mu_2^{-1}(a_2\rho_2 + a_2, \pi_2)| \leq_1 |\mu_4^{-1}(a_2\rho_2 + a_2, \pi_2)|$ since both sets contain at least one element.

Observe that all data values from $\mathbb{D} \setminus D_1$ have one of the two μ_4 -profiles we have just treated. Then, for any other mosaic (ρ, π) we know that all its data values are simulated by data values of μ_2 , and we can then easily prove that $|\mu_2^{-1}(a_2\rho_2 + a_2, \pi_2)| \leq_1 |\mu_4^{-1}(a_2\rho_2 + a_2, \pi_2)|$. Then we conclude that condition 4 of $\mu_2 \leq_K \mu_4$ holds. ■

Claim 9. $\mu'_n \succrightarrow \mu''_1$.

Proof: For simplicity of notation let us call $e_0 = d_3$. Suppose first that $d_2 \neq d_1$, and hence that $(\{e_0\} \cup P) \cap R = \emptyset$, as already pointed out. Take any e_i . If $i \neq n$, then $\mu'_n(e_i) = (\rho_1, a_1\pi_1)$ and $\mu''_1(e_i) = (a_2\rho_1, \pi_1)$, which verifies $\mathbf{m}(a_1, 0, (\rho_1, a_1\pi_1), a_2, 0, (a_2\rho_1, \pi_1))$. If $i = n$, then $\mu'_n(e_i) = (\rho_1, a_1\pi_1 + a_1)$ and $\mu''_1(e_i) = (a_2\rho_1, \pi_1)$, which verifies $\mathbf{m}(a_1, 1, (\rho_1, a_1\pi_1 + a_1), a_2, 0, (a_2\rho_1, \pi_1))$. Take now any e'_j . Note that $\mu'_n(e'_j) = \mu_3(e'_j) = \mu_1(d_2) = (\rho_2, a_1\pi_2)$. If $j \neq 1$, then $\mu''_1(e'_j) = (a_2\rho_2, \pi_2)$, which verifies $\mathbf{m}(a_1, 0, (\rho_2, a_1\pi_2), a_2, 0, (a_2\rho_2, \pi_2))$. And if $j = 1$, then $\mu''_1(e'_j) = (a_2\rho_2 + a_2, \pi_2)$, which verifies $\mathbf{m}(a_1, 0, (\rho_2, a_1\pi_2), a_2, 1, (a_2\rho_2 + a_2, \pi_2))$.

The pathological case is when $d_2 = d_1$. As we already remarked, in this situation we have $\{e_0\} \cup P = R$. Also, we have $\rho_2 = \rho_1$ and $\pi_2 = \pi_1$. Then take any data value $e_i = e'_j$. If $i \neq n$ and $j \neq 1$, $\mu'_n(e_i) = (\rho_1, a_1\pi_1)$ and $\mu''_1(e_i) = \mu''_1(e'_j) = (a_2\rho_2, \pi_2) = (a_2\rho_1, \pi_1)$; it is immediate to verify that $\mathbf{m}(a_1, 0, (\rho_1, a_1\pi_1), a_2, 0, (a_2\rho_1, \pi_1))$ holds true. If $i = n$ and $j \neq 1$, $\mu'_n(e_i) = (\rho_1, a_1\pi_1 + a_1)$ and $\mu''_1(e_i) = (a_2\rho_1, \pi_1)$, and hence $\mathbf{m}(a_1, 1, (\rho_1, a_1\pi_1 +$

$a_1), a_2, 0, (a_2\rho_1, \pi_1))$. If $i \neq n$ and $j = 1$, $\mu'_n(e_i) = (\rho_1, a_1\pi_1)$, $\mu''_1(e_i) = (a_2\rho_2 + a_2, \pi_2) = (a_2\rho_1 + a_2, \pi_1)$, which verifies $\mathbf{m}(a_1, 0, (\rho_1, a_1\pi_1), a_2, 1, (a_2\rho_1 + a_2, \pi_1))$. Finally, if $i = n$ and $j = 1$, $\mu'_n(e_i) = (\rho_1, a_1\pi_1 + a_1)$, $\mu''_1(e_i) = (a_2\rho_1 + a_2, \pi_1)$ which trivially verifies $\mathbf{m}(a_1, 1, (\rho_1, a_1\pi_1 + a_1), a_2, 1, (a_2\rho_1 + a_2, \pi_1))$. ■

Claim 10. For every $i \in [m-1]$, $\mu''_i \succrightarrow \mu''_{i+1}$.

Proof: Since μ''_i and μ''_{i+1} have the same letter, it is immediate that the preservation of the profile for every data value other than those of R is compatible with the definition of \mathbf{m} . In other words, that $\mathbf{m}(a_2, 0, \mu''_i(d), a_2, 0, \mu''_{i+1}(d))$ for every $d \notin R$, since $\mu''_i(d) = \mu''_{i+1}(d)$. For e'_{i+1} we have to verify $\mathbf{m}(a_2, 0, (a_2\rho_2, \pi_2), a_2, 1, (a_2\rho_2 + a_2, \pi_2))$ which is true by (7) and (11). For e'_i we verify $\mathbf{m}(a_2, 1, (a_2\rho_2 + a_2, \pi_2), a_2, 0, (a_2\rho_2 + a_2, \pi_2))$. Notice that in this case, (4) is verified since π_2 of the form $\pi_2 = a_2\pi' + a_2$, and then we have that $\pi_2 + a_2 = \pi_2$ by idempotence of $+a_2$ and (8) is trivially verified. Finally, for e'_j with $j < i$ or $j > i+1$, we can apply (5) and (9) since they preserve the profile. ■

Claim 11. For all $i \in [n]$, μ''_i is valid and K -compliant.

Proof: This proof is completely equivalent and symmetrical to the proof of Claim 1. There, it was showed that μ'_1, \dots, μ'_n are valid and K -compliant from the fact that μ_3 is valid and K -compliant. The same argument works to show that $\mu''_{m-1}, \dots, \mu''_1$ are valid and K -compliant from the fact that μ''_m is valid and K -compliant. Observe that μ''_m is valid and K -compliant by Lemma 5 since $\mu''_m \geq_K \mu_2$. ■

Complex situation 3. Finally, if condition A hold true, it is the same procedure as before with $P = \emptyset$. ■

Proof of Claim 6: By definition of \leq_K we have that μ_1 and μ_2 verify the properties 1, 2 and 3 of \leq_K , but not necessarily 4. However, we still conserve the following property. For every profile (ρ, π) ,

$$|\mu_1^{-1}(\rho, \pi) \setminus K| = 0 \quad \text{iff} \quad |\mu_2^{-1}(\rho, \pi) \setminus K| = 0 \quad (12)$$

because this property must coincide with the mosaic μ_3 verifying $|\mu_3^{-1}(\rho, \pi) \setminus K| = 0$. This is a consequence of the property 4 of $\mu_3 \leq_K \mu_1$ and $\mu_3 \leq_K \mu_2$.

We define $\nu = (a, d, \overleftarrow{\chi}, \overrightarrow{\chi})$, where a and d are the label and data value of μ_1 (the label is also equal to that of μ_2). Consider the mosaic μ'_2 which is the result of renaming some data values of μ_2 in such a way that it has the same data value d as μ_1 , and the only data values shared between μ'_2 and μ_1 are a subset of $K \cup \{d\}$. We then define

- $\nu(e) = \mu_1(e)$ if $\mu_1(e) \neq \emptyset$, and
- $\nu(e) = \mu'_2(e)$ if $\mu_1(e) = \emptyset$.

Observe that this definition preserves all the profiles of the rigid values and of the data value d . Further, for every profile (ρ, π) both the flexible data values of $\mu_1^{-1}(\rho, \pi)$ and of $(\mu'_2)^{-1}(\rho, \pi)$ are subsets of those of $\nu^{-1}(\rho, \pi)$. This, combined with the property (12) ensures that the condition 4 of the

definition of \leq_K holds both for μ_1 and μ'_2 , and we hence have $\nu \geq_K \mu_1$, $\nu \geq_K \mu'_2$. Since μ'_2 and μ_2 are equal modulo renaming of data values from $\mathbb{D} \setminus K$, by definition of \leq_K we obtain $\nu \geq_K \mu'_2$, and we conclude the proof. ■

APPENDIX TO SECTION V

Proof of Lemma 12: Notice that, by definition of \leq_K , we need count the following

- the possible profiles of the current data value of the mosaic, which is bounded by $Prof \stackrel{\text{def}}{=} |\wp(\Omega_\varphi^p) \times \wp(\Omega_\varphi^p)|$
- the possible profiles of all the data values from K , bounded by $Prof^{|K|}$
- the possible configurations accounting for which profiles have one flexible data value and which are empty, bounded by $2^{|\wp(\Omega_\varphi^p)|}$,
- the possible labels, bounded by $|\mathbb{A}|$.

We then obtain that there are not more than $|\mathbb{A}| \cdot Prof^{|K|+1} \cdot 2^{|\wp(\Omega_\varphi^p)|}$. In fact, $|\mathbb{A}|$ can be considered to be polynomial in $|\varphi|$, since otherwise we can replace all letters not appearing in φ by a single one \perp . Hence the number of minimal elements of $\mathcal{M}(\varphi)$ with respect to \leq_K is double exponential in $|\varphi|$. ■

Proof of Lemma 13: Let $\hat{\mu}_i$ be a mosaic such that $\hat{\mu}_i \leq_K \mu_i$ and $\hat{\mu}_i \in \text{MIN}_{\mathcal{M}(\varphi)}$ for all $1 \leq i \leq n$ (it always exist). Then we have that $\hat{\mu}_1 \rightarrow_K \dots \rightarrow_K \hat{\mu}_n$, $\hat{\mu}_1 \leq_K \mu_1$, $\hat{\mu}_n \leq_K \mu_n$. Further, the sequence is complete, because \leq_K preserves the left/right completeness of a mosaic by Lemma 5. If there are two indices $i < j$ such that $\hat{\mu}_i = \hat{\mu}_j$, then we can remove all the mosaics between them obtaining still a complete \rightarrow_K sequence with the aforementioned properties. We can repeat this procedure until we obtain a \rightarrow_K sequence where all the elements are different. Hence, the sequence is bounded by $|\text{MIN}_{\mathcal{M}(\varphi)}|$. ■

A. Proof of Lemma 14

To prove Lemma 14, we first need to be able to compute \rightarrow_K . The following lemma shows how to compute such relation.

Lemma 17. *Given two K -compliant mosaics $\mu_1, \mu_2 \in \mathcal{M}(\varphi)$, $\mu_1 \rightarrow_K \mu_2$ iff for every data value $d \in \mathbb{D}$:*

1. if $d \in K$ then $\mathbf{m}(a_1, d \stackrel{?}{=} d_1, \mu_1(d), a_2, d_2 \stackrel{?}{=} d, \mu_2(d))$ holds,
2. if $d \in \mathbb{D} \setminus K$ then there is a bijection $f : \mathbb{D} \setminus K \rightarrow \mathbb{D} \setminus K$ such that
 - a) $\mathbf{m}(a_1, d \stackrel{?}{=} d_1, \mu_1(d), a_2, d_2 \stackrel{?}{=} f(d), \mu_2(f(d)))$, or
 - b) if $\mu_1(d) = (\emptyset, \emptyset) \neq \mu_2(f(d))$, then

$$\mathbf{m}(a_1, d \stackrel{?}{=} d_1, (\rho, \pi), a_2, d_2 \stackrel{?}{=} f(d), \mu_2(f(d)))$$

for some profile (ρ, π) such that $\mu_1^{-1}(\rho, \pi) \setminus K \neq \emptyset$, or

- c) if $\mu_1(d) \neq (\emptyset, \emptyset) = \mu_2(f(d))$, then

$$\mathbf{m}(a_1, d \stackrel{?}{=} d_1, \mu_1(d), a_2, d_2 \stackrel{?}{=} f(d), (\rho, \pi))$$

for some profile (ρ, π) such that $\mu_2^{-1}(\rho, \pi) \setminus K \neq \emptyset$.

Proof of Lemma 17: (if) Suppose we have two mosaics μ_1 and μ_2 that verify the conditions imposed by the Lemma through a witnessing bijection $f : \mathbb{D} \setminus K \rightarrow \mathbb{D} \setminus K$. Consider two mosaics μ'_1 and μ'_2 such that $a'_1 = a_1$, $d'_1 = d_1$, $a'_2 = a_2$ and $d'_2 = f^{-1}(d_2)$ if $d_2 \notin K$, or $d'_2 = d_2$ otherwise. For every $d \in K$ we define $\mu'_1(d) = \mu_1(d)$ and $\mu'_2(d) = \mu_2(d)$. For every $d \in \mathbb{D} \setminus K$ we define

- (i) if $\mu_1(d) = \mu_2(f(d)) = (\emptyset, \emptyset)$ then $\mu'_1(d) = \mu'_2(d) = (\emptyset, \emptyset)$
- (ii) if $\mu_1(d) = (\emptyset, \emptyset)$ but $\mu_2(f(d)) \neq (\emptyset, \emptyset)$, there is some (ρ_1, π_1) such that $\mathbf{m}(a_1, d_1 \stackrel{?}{=} d, (\rho_1, \pi_1), a_2, d_2 \stackrel{?}{=} f(d), \mu_2(f(d)))$ where $\mu_1^{-1}(\rho_1, \pi_1)$ contains at least one flexible data value e by condition 2b. We then define $\mu'_1(d) = (\rho_1, \pi_1)$, $\mu'_2(d) = \mu_2(f(d))$.
- (iii) if $\mu_1(d) \neq (\emptyset, \emptyset)$ but $\mu_2(f(d)) = (\emptyset, \emptyset)$, there must be (ρ_2, π_2) such that $\mathbf{m}(a_1, d_1 \stackrel{?}{=} d, \mu_1(d), a_2, d_2 \stackrel{?}{=} f(d), (\rho_2, \pi_2))$ where $\mu_2^{-1}(\rho_2, \pi_2)$ contains at least one flexible data value e by condition 2c. We then define $\mu'_1(d) = \mu_1(d)$, $\mu'_2(d) = (\rho_2, \pi_2)$.
- (iv) Finally, if $\mu_1(d) \neq (\emptyset, \emptyset)$ and $\mu_2(f(d)) \neq (\emptyset, \emptyset)$, we define $\mu'_1(d) = \mu_1(d)$, $\mu'_2(d) = \mu_2(f(d))$.

We need to show now that $\mu'_1 \geq_K \mu_1$, $\mu'_2 \geq_K \mu_2$ and $\mu'_1 \rightarrow_K \mu'_2$. Note that by Lemma 5 we will obtain that μ'_1 and μ'_2 are also valid and K -compliant.

$\mu'_1 \geq_K \mu_1$. The conditions 1, 2 of \leq_K are easily verified by construction of μ'_1 . Condition 3 also holds since because $d_1 = d'_1$ and because $\mu_1(d_1) \neq \emptyset$ it is defined by conditions (iii) or (iv) in which $\mu_1(d_1) = \mu'_1(d_1)$. To check that condition 4 is verified, take any profile (ρ, π) . First note that

$$\mu_1^{-1}(\rho, \pi) \subseteq (\mu'_1)^{-1}(\rho, \pi). \quad (13)$$

Further, if $\mu_1^{-1}(\rho, \pi)$ has no flexible values, this means that condition 2b of the Lemma is never true for the profile (ρ, π) . Note that the only possible step in the construction where μ'_1 is defined to be something else than the same profile as in μ_1 is in step ii, but this step is never reached with the profile (ρ, π) since it needs 2b to hold true. Then, we have

$$\text{if } \mu_1^{-1}(\rho, \pi) \setminus K = \emptyset \text{ then } (\mu'_1)^{-1}(\rho, \pi) \setminus K = \emptyset. \quad (14)$$

Hence, by (13) and (14), condition 4 of \leq_K is met and we have $\mu'_1 \geq_K \mu_1$.

$\mu'_2 \geq_K \mu_2$. The conditions 1, 2 of \leq_K are easily verified by construction of μ'_2 . If $d_2 \in K$, condition 3 is also easily verified. If $d_2 \notin K$, we have by construction that $d'_2 = f^{-1}(d_2)$. We need to show that $\mu_2(d_2) = \mu'_2(f^{-1}(d_2))$. For $f^{-1}(d_2)$ condition (i) cannot be true, since it cannot be true that $\mu_2(f(f^{-1}(d_2))) = \mu_2(d_2) = (\emptyset, \emptyset)$ by definition of mosaic. Condition (iii) cannot either be true for the same reason. Then, $\mu'_2(f^{-1}(d_2))$ has to be built according to condition (ii) or (iv), which define $\mu'_2(f^{-1}(d_2)) = \mu_2(f(f^{-1}(d_2))) = \mu_2(d_2)$. Then, condition 3 of $\mu'_2 \geq_K \mu_2$ holds. To check that condition 4 is verified, take any profile (ρ, π) . First note that

$$\mu_2^{-1}(\rho, \pi) \setminus K \subseteq f((\mu'_2)^{-1}(\rho, \pi)) \setminus K. \quad (15)$$

Further, if $\mu_2^{-1}(\rho, \pi)$ has no flexible values, this means that condition 2c of the Lemma is never true for the profile (ρ, π) . As before the only possible step in the construction where $\mu'_2(d)$ is defined to be something different to $\mu_2(f(d))$ is in step (iii), but this step is never reached with the profile (ρ, π) since it needs 2c to hold true. Then, we have

$$\text{if } \mu_2^{-1}(\rho, \pi) \setminus K = \emptyset \text{ then } (\mu'_2)^{-1}(\rho, \pi) \setminus K = \emptyset. \quad (16)$$

Hence, by (15) and (16), condition 4 of \leq_K is met and we have $\mu'_2 \geq_K \mu_2$.

$\mu'_1 \succ \mu'_2$. Take any data value $d \in \mathbb{D}$. If $d \in K$ then by condition 1 of the Lemma we have $\mathbf{m}(a_1, d \stackrel{?}{=} d_1, \mu_1(d), a_2, d_2 \stackrel{?}{=} d, \mu_2(d))$ and by construction $\mu_1(d) = \mu'_1(d)$, $\mu_2(d) = \mu'_2(d)$. Hence,

$$\mathbf{m}(a_1, d \stackrel{?}{=} d_1, \mu_1(d), a_2, d_2 \stackrel{?}{=} d, \mu_2(d)).$$

Suppose then that $d \notin K$. If $\mu_1(d) = \mu_2(f(d)) = (\emptyset, \emptyset)$ we have that $\mathbf{m}(a_1, d \stackrel{?}{=} d_1, \mu_1(d), a_2, d_2 \stackrel{?}{=} f(d), \mu_2(f(d)))$ by item 2a. By step (i) of the construction, we have that $\mu'_1(d) = \mu'_2(d) = (\emptyset, \emptyset)$ and then $\mathbf{m}(a'_1, d \stackrel{?}{=} d'_1, \mu'_1(d), a'_2, d'_2 \stackrel{?}{=} d, \mu'_2(d))$. If $\mu_1(d) \neq (\emptyset, \emptyset)$ and $\mu_2(f(d)) \neq (\emptyset, \emptyset)$ we proceed in a similar way. By step (iv) of the construction we have that $\mu'_1(d) = \mu_1(d)$ and $\mu'_2(d) = \mu_2(f(d))$ and we then conclude $\mathbf{m}(a'_1, d \stackrel{?}{=} d'_1, \mu'_1(d), a'_2, d'_2 \stackrel{?}{=} d, \mu'_2(d))$. Otherwise, if only one of $\mu_1(d)$, $\mu_2(f(d))$ is empty, either by step (ii) or (iii) we will define $\mu'_1(d)$ or $\mu'_2(d)$ as the profile necessary to verify $\mathbf{m}(a'_1, d \stackrel{?}{=} d'_1, \mu'_1(d), a'_2, d'_2 \stackrel{?}{=} d, \mu'_2(d))$, which exists by condition 2b or 2c.

(only if) If we have two mosaics $\mu_1 \rightarrow_K \mu_2$ then there must exist $\mu_1 \leq_K \mu'_1 \succ \mu'_2 \geq_K \mu_2$. By $\mu_1 \leq_K \mu'_1$ there is a bijection of data values $f_1 : \mathbb{D} \setminus K \rightarrow \mathbb{D} \setminus K$ that sends each data value of μ_1 with non-empty profile to a data value in μ_2 with the same profile. Similarly we have a bijection f_2 that relates the data values between μ_2 and μ'_2 . We can then define $f = f_1 \circ f_2^{-1}$ and it is easy to check that it verifies all the conditions defined in the Lemma. ■

Proof of Lemma 14: By Lemma 17, we only need to test the conditions specified there. We disregard the space needed to encode the data values, because we can always abstract them away in classes of equivalence. The algorithm simply needs to first guess a bijection f only restricted to the data values of μ_1 and μ_2 that have non-empty profiles. That is, we are only interested in pairs of data values (d, d') such that $\mu_1(d) \neq \emptyset$ or $\mu_2(f(d)) \neq \emptyset$. Then, this bijection f needs only polynomial space. Then, we verify the conditions for each data value in the bijection, each condition taking polynomial space to test. ■

Proof of Lemma 11: Consider any complete \succ sequence $\mu_1 \succ \dots \succ \mu_n$ over $\mathcal{M}(\varphi)$. Consider also a pair of functions $f_l, f_r : \mathbb{D} \rightarrow \wp(\Omega_\varphi^p)$ that assign to each data value a set of path expressions such that the following holds for every $d \in D$

- $\alpha \in f_l(d)$ iff there is a position $1 \leq i \leq n$ s.t. $\bar{\chi}_i(\alpha) = \{d\}$.

- $\alpha \in f_r(d)$ iff there is a position $1 \leq i \leq n$ s.t. $\bar{\chi}_i(\alpha) = \{d\}$.

We claim that for every two data values $d \neq d'$, $f_l(d) \cap f_l(d') = \emptyset$, and $f_r(d) \cap f_r(d') = \emptyset$. By way of contradiction, suppose that $\alpha \in f_l(d) \cap f_l(d')$, and assume that the witnessing positions for this membership are i and i' . Suppose that $i \leq i'$ without any loss of generality. This means that from position i we can access the data value d by going to the left through a path expression α . Since we can reach it from i we can also reach it from i' , since $i \leq i'$. Then we can reach both d and d' from i' , which is in contradiction with i' being the witnessing position for $\alpha \in f_l(d')$.

Then, there are not more than $|\Omega_\varphi^p|$ data values d such that $f_l(d) \neq \emptyset$, and likewise for f_r . We can build K by collecting these data values, and we will have that $|K| \leq 2 \cdot |\Omega_\varphi^p|$. ■

Proof of Proposition 3: For simplicity we build a NEXPSpace algorithm. Since $\text{NEXPSpace} = \text{EXPSpace}$ (Savitch's Theorem) this proves the proposition. The algorithm tries to build progressively a complete sequence $\mu'_1 \rightarrow_K \dots \rightarrow_K \mu'_m$ such that $\mu'_i \in \text{MIN}_{\mathcal{M}(\varphi)}$ for every i and $\mu'_1 \vdash \varphi_1$. It answers 'yes' iff there is such a sequence. By Lemma 13 we know that if there is some complete \rightarrow_K sequence then there is one over minimal elements of size bounded by $|\text{MIN}_{\mathcal{M}(\varphi)}|$. The algorithm will count the length of the sequence it is building. Once the counter reaches $|\text{MIN}_{\mathcal{M}(\varphi)}|$ it fails, answering that there is no complete sequence.

The algorithm maintains two variables:

- one counter c to count from 1 up to $|\text{MIN}_{\mathcal{M}(\varphi)}|$, and
- the last mosaic it has treated μ .

Both variables require exponential space by Lemma 12. In the first step it initializes the counter $c = 1$, and the variable μ with a mosaic from $\text{MIN}_{\mathcal{M}(\varphi)}$ that it guesses. Next, it checks that it is left-complete and that verifies φ_1 (otherwise it answers 'no'). Then it enters in a loop. At each step it checks if μ is complete to the right. If it is right-complete, it answers 'yes'. Otherwise it performs the following actions.

- If $c < |\text{MIN}_{\mathcal{M}(\varphi)}|$ it increments c by one, otherwise it answers 'no'.
- It guesses a mosaic $\mu' \in \text{MIN}_{\mathcal{M}(\varphi)}$,
- It checks that $\mu \rightarrow_K \mu'$ using exponential space by Lemma 14. If it does not hold it answers 'no'.
- Finally, it overwrites μ with μ' .

If there is a complete sequence as the one described in Lemma 13, it will first guess μ'_1 , then μ'_2 up to μ'_m and it will answer 'yes'. If there is no complete sequence, it will answer 'no' because for every possible guessed sequence $\mu'_1 \dots \mu'_m$ either

- μ'_1 is not left-complete, or
- $\mu'_1 \not\vdash \varphi_1$, or
- μ'_m is not right-complete
- there is some μ'_i and μ'_{i+1} that are not in the \rightarrow_K relation, or
- $m > |\text{MIN}_{\mathcal{M}(\varphi)}|$.

In each one of these cases the algorithm answers 'no'. ■

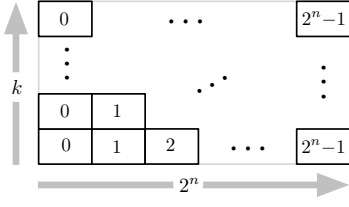


Fig. 4. 2^n corridor tiling. Each depicted box contains a tile from T . Each number inside a box is the *index* of the position.

APPENDIX TO SECTION VI

Proof of Proposition 4: We show this by coding an instance of the 2^n corridor tiling problem. An instance of this problem consists in a size of the corridor n (encoded in unary), a set of tiles $T = \{T_1, \dots, T_s\}$, a special final tile T_s , an initial tile T_1 , and the horizontal and vertical tiling relations $H, V \subseteq T \times T$. The problem consists in answering whether there is $k \in \mathbb{N}$ and a function $f : 2^n \times k \rightarrow T$ such that $f(1, 1) = T_0$, $f(2^n, k) = T_s$, and $H(f(i, j), f(i+1, j))$, $V(f(i, j), f(i, j+1))$ for all possible i, j . In this case we say that the problem has a solution. This problem is known to be complete for EXPSpace [Chl86]. We will see f as a matrix of 2^n columns and k rows, whose indices grow to the left and to the top, like in Figure 4.

Groups: Since PathLogic cannot avoid having repeated consecutive nodes with the same symbol, in this coding we will always have to work with *groups* of nodes with the same symbol. Unfortunately this makes the coding somewhat more cumbersome. For the purpose of this poof, we consider a data values as elements of $(\mathbb{A} \times \mathbb{D})^*$, as opposed to the functional definition given in Section II. (Of course, both definitions are equivalent.) This will simplify the description of the structure in our coding.

Let us fix an instance of the 2^n corridor tiling problem. We will code a solution like the one in Figure 4 in a data word by coding each of the position of the matrix, starting with the first row from left to right, then with the second row, etc.

The particular strategy to link consecutive columns and rows follows a similar coding idea as the one developed in [FS09]. Let

$$A = \{b_1, \dots, b_n, \mathbf{N}_{\rightarrow}, \mathbf{N}_{\uparrow}, \#, @, \mathbf{beg}, \mathbf{end}, T_1, \dots, T_s\}.$$

We will use data values as *pointers* to future elements, in the sense that “a node (a, d) points to another node (b, e) ” if (b, e) appears to the right of (a, d) and $d = e$. The coding of a run consists in a succession of *blocks*. Each block corresponds to a position of the solution matrix. We first describe the general structure of the block in terms of the labels it must contain, and then we describe more precise properties that this structure must verify.

The shape of each block is of the form

$$(\mathbf{beg}, d_1) \cdots (\mathbf{beg}, d_n) \ w \ (\mathbf{end}, e_1) \cdots (\mathbf{end}, e_m)$$

where $w \in (A \times D)^*$, and such that for every d_i there exists e_j with $d_i = e_j$. Further, every data value of w is not equal

to any d_i nor e_j . Also, w is of the form

$$\text{Tile Num } \text{Next}_{\rightarrow} \ \text{Next}_{\uparrow}$$

such that

- $\text{Tile} \in (T \times \mathbb{D})^+$ is a group of nodes with the same symbol T_i , for some $i \leq s$ (the data values will play no role here).
- $\text{Next}_{\rightarrow} \in (\{\mathbf{N}_{\rightarrow}\} \times \mathbb{D})^*$, such that each element $(\mathbf{N}_{\rightarrow}, d)$ points to an element (\mathbf{beg}, d) of a future block. (The idea will be that every one of these nodes point to the next block of the same row and next column.)
- $\text{Next}_{\uparrow} \in (\{\mathbf{N}_{\uparrow}\} \times \mathbb{D})^*$, such that each element $(\mathbf{N}_{\uparrow}, d)$ points to an element (\mathbf{beg}, d) of a future block. (The idea is that these nodes point to the block of the same column and next row.)
- $\text{Num} \in (\{b_1, \dots, b_n\} \times \mathbb{D})^*$ is the coding of a binary number of n bits (the data values play no role here). We interpret that the i -th bit is 1 iff b_i is in Num . From now on we call the *index* of the block, to the number coded in Num . Figure 4 shows the indices for each position of the matrix.

Now we constrain this general structure to ensure that it encodes a valid solution of the problem. In the sequel we say that a block B points to a block B' through some symbol S iff B contains an element (S, d) and B' contains an element (\mathbf{beg}, d') such that $d = d'$.

- The first block has the tile T_1 .
- The last block has the tile T_s .
- If the $\text{Next}_{\rightarrow}$ part of a block B is the empty string, then B is the last block.
- The first block of the data word has index 0.
- For every pair of blocks B, B' such that B points to B' through \mathbf{N}_{\rightarrow} then
 - the index n of B and n' of B' are such that $n' = n + 1 \pmod{2^n}$
 - B' is the next immediate block of B
 - the tile t_i of B and t_j of B' are in H -relation.
- If B points to B' through \mathbf{N}_{\uparrow} then
 - the index of B equals to the index of B'
 - there is exactly one block B'' of index 0 between B and B' .
 - the tile t_i of B and t_j of B' are in V -relation.

It is easy to see that there exists a model with these properties iff the 2^n corridor tiling problem instance has a solution.

It is possible to code these conditions using similar techniques as in [FS09]. The cited work uses a different logic that is expressive equivalent $\text{XPath}(\rightarrow^+, =)$ to code a much harder (non-primitive recursive) problem which involves properties that we cannot express with PathLogic. However, in order to code these properties we do not need the *strict* transitive relation \rightarrow^+ , we can use \rightarrow^* . ■

Remark 3. The above coding of EXPSpace-hardness uses only formulas that move in one direction.

APPENDIX TO SECTION VII

$\text{XPath}(*\leftarrow, \rightarrow^*, =)$ and $\text{plain-XPath}(*\leftarrow, \rightarrow^*, =)$ are expressive-equivalent. There is a translation that transforms each node expression of $\text{XPath}(*\leftarrow, \rightarrow^*, =)$ in an exponential disjunction of node expressions of $\text{plain-XPath}(*\leftarrow, \rightarrow^*, =)$.

Proof idea of Lemma 16: We can factorize any atomic node expression like $\langle \alpha = \beta \rangle$ or $\langle \alpha \neq \beta \rangle$ uses path expressions with two directions into a disjunction of expressions, such that each one of these is in $\text{plain-XPath}(*\leftarrow, \rightarrow^*, =)$. We proceed by induction on the nesting depth of the subformula, defined as the maximum number of nested path expressions it contains. The base case is immediate. Take an expression $\langle \alpha = \beta \rangle$ where N_α and N_β are the sets of node expressions $[\psi]$ that use respectively α and β , where each of these is a formula of $\text{plain-XPath}(*\leftarrow, \rightarrow^*, =)$ by inductive hypothesis. The translation consists in checking from all the possible linearizations of $N_\alpha \cup N_\beta \cup \{curr\}$ those that are consistent with α, β . Here, $curr$ represents the current point of evaluation. For example, the linearization $[\psi_1][\psi_3]curr[\eta][\psi_2]$ is consistent with $\alpha = *\leftarrow[\psi_1]\rightarrow^*[\psi_2]*\leftarrow[\psi_3]$, $\beta = \rightarrow^*[\eta]$. However, $[\psi_1][\psi_2][\psi_3]curr[\eta]$ is not. For each consistent linearization the translation produces one disjunct $\langle \alpha' = \beta' \rangle$ with the aforementioned properties. For our example of linearization it would produce

$$\langle [\psi_1]*\leftarrow[\psi_3]*\leftarrow = \rightarrow^*[\eta]\rightarrow^*[\psi_2] \rangle.$$

The procedure for formulas $\langle \alpha \neq \beta \rangle$ is similar. This translation is exponential. ■