

Solutions in XML Data Exchange

Mikołaj Bojańczyk
University of Warsaw
bojan@mimuw.edu.pl

Leszek A. Kołodziejczyk
University of Warsaw
lak@mimuw.edu.pl

Filip Murlak
University of Warsaw
fmurlak@mimuw.edu.pl

ABSTRACT

The task of XML data exchange is to restructure a document conforming to a source schema under a target schema according to certain mapping rules. The rules are typically expressed as source-to-target dependencies using various kinds of patterns, involving horizontal and vertical navigation, as well as data comparisons. The target schema imposes complex conditions on the structure of solutions, possibly inconsistent with the mapping rules. In consequence, for some source documents there may be no solutions.

We investigate three problems: deciding if all documents of the source schema can be mapped to a document of the target schema (absolute consistency), deciding if a given document of the source schema can be mapped (solution existence), and constructing a solution for a given source document (solution building).

We show that the complexity of absolute consistency is rather high in general, but within the polynomial hierarchy for bounded depth schemas. The combined complexity of solution existence and solution building behaves similarly, but the data complexity turns out to be very low.

In addition to this we show that even for much more expressive mapping rules, based on MSO definable queries, absolute consistency is decidable and data complexity of solution existence is polynomial.

Categories and Subject Descriptors

H.2.5 [Database Management]: Heterogeneous Databases—*Data translation*; I.7.2 [Document and Text Processing]: Document Preparation—*XML*

General Terms

Theory, Algorithms, Languages

Keywords

XML data exchange, regular queries, patterns, absolute consistency, solution building, solution existence

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

ICDT 2011, March 21–23, 2011, Uppsala, Sweden.
Copyright 2011 ACM 978-1-4503-0529-7/11/0003 ...\$10.00

1. INTRODUCTION

One of the main challenges of modern data management is dealing with heterogeneous data. A typical scenario is that of data exchange, where one needs to restructure data stored in a source database under a target database schema, following a specification. The specification is given by a so-called *schema mapping*, a collection of logical formulas describing dependencies between the source schema and the target schema. The produced instance of data is called a *solution* for the source data with respect to the schema mapping.

Studies on relational data exchange and schema mappings were initiated several years ago [11, 12], and since then have been a major topic (see recent surveys [5, 6, 15]).

In the XML context, the target schema, a DTD or XML Schema, imposes complex conditions on the structure of the solution. One of the consequences is that in practice schema mappings are often *over-specified*, in the sense that for some source instances there is no solution satisfying the specification and conforming to the target schema. One needs to be able to check if a mapping admits a solution or not.

Nowadays, schema mapping design is usually assisted by specialized software [16]. The ability to issue warnings that a mapping does not admit solutions for some source documents would be a most welcome feature, especially if the warning could also point out the source of the problem.

When a mapping is already there, one has to materialize the appropriate solution for the given source data. This requires checking if a solution exists, and constructing it.

We concentrate on three problems related to solutions:

- *Absolute consistency*. Is there a solution for every possible source document?
- *Solution existence*. Is there a solution for a given source document?
- *Solution building*. Find a solution for a given source document.

So far, absolute consistency was only investigated for very simple mappings, based on so-called tree patterns using only the child and descendant relations [1]. Solution building was investigated for even simpler mappings allowing only a restricted class of DTDs [3]. More practical mapping languages involve sibling order and data comparisons. For such mappings only the *consistency* problem was considered, i.e., whether *some* source document has a solution [1, 3]. Consistency of mappings, though interesting theoretically, seems to have less practical importance as it gives no information on how the mapping will perform on a given source document.

We work with mappings using general XML schemas, formalized as tree automata, DAG-shaped patterns using horizontal and vertical navigation, as well as data comparisons. We also consider a restricted case: *bounded-depth* mappings, where schemas only admit trees of bounded height. As most real-life schemas have small depth, this case is particularly important in practice.

Our main findings for pattern-based mappings can be summarized as follows:

- Absolute consistency is $\Pi_2\text{EXP}$ -complete¹ in general, and $\Pi_4\text{P}$ -complete in the bounded depth case.
- Data complexity of solution existence is in LOGSPACE , while combined complexity is NEXP -complete in general, and $\Sigma_3\text{P}$ -complete for the bounded depth case.
- Solution building can be done in EXPSpace in general, in PSPACE in the bounded depth case, and in P if the mapping is fixed.

The high combined complexity of the general case is disturbing, but the practically relevant bounded-depth case brings it down to acceptable low levels of polynomial hierarchy, and the polynomial data complexity of solution existence and solution building gives real hope for applications.

We then go on to generalize these results even further. We show that for mappings based on queries definable in MSO , absolute consistency remains decidable, and data complexity of solution existence is still polynomial.

The paper is organized as follows. After recalling the basic notions (Sect. 2), we study the absolute consistency problem for pattern-based mappings (Sect. 3) and its restriction to mappings of bounded depth (Sect. 4). We finish the first part of the paper with the examination of solution building and solution existence (Sect. 5). The second part starts with a definition of regular queries and then develops a representation designed to facilitate deciding absolute consistency and solution existence of mappings using regular queries instead of patterns (Sect. 6). Next, we give the decision procedures (Sect. 7), and finish with some suggestions for future work (Sect. 8).

Due to space limitations some proofs are omitted.

2. PRELIMINARIES

Data trees, DTDs, automata. The abstraction of XML documents we use is *data trees*: unranked labelled trees storing in each node a natural number called *data value*. Formally, a *data tree* over a finite labelling alphabet Γ is a structure $\langle T, \downarrow, \downarrow^*, \rightarrow, \rightarrow^*, \ell_T, \rho_T \rangle$, where

- the set T is an unranked tree domain, i.e., a prefix-closed subset of \mathbb{N}^* such that $n \cdot i \in T$ implies $n \cdot j \in T$ for all $j < i$;
- the binary relations \downarrow and \rightarrow are the child relation ($n \downarrow n \cdot i$) and the next-sibling relation ($n \cdot i \rightarrow n \cdot (i+1)$);
- \downarrow^* and \rightarrow^* are transitive closures of \downarrow and \rightarrow ;
- the function ℓ_T is a labelling from T to Γ ;
- ρ_T is a function from T to \mathbb{N} . We say that a node $s \in T$ stores the value v when $\rho_T(s) = v$.

¹ $\Pi_2\text{EXP}$ is the second level of the exponential hierarchy.

Most often, when the interpretations of $\downarrow, \rightarrow, \ell_T$, and ρ_T are understood, we write just T to refer to a data tree. We use the terms “tree” and “data tree” interchangeably.² We write $|T|$ to denote the number of nodes of T .

The principal schema language we use are tree automata, abstracting Relax NG [17, 18].

A *nondeterministic word automaton* can be presented as a tuple $\mathcal{B} = \langle \Gamma, Q, q_I, \delta, F \rangle$, where Γ is a finite input alphabet, Q is a finite set of states, $q_I \in Q$ is the initial state, $F \subseteq Q$ is the set of final states, and $\delta \subseteq Q \times \Gamma \times Q$ is the transition relation. A run on a word $\sigma_1\sigma_2 \cdots \sigma_n \in \Gamma^*$ is a sequence $q_1q_2 \cdots q_{n+1} \in Q^*$ where q_1 is the initial state q_I , and $(q_i, \sigma_i, q_{i+1}) \in \delta$ for $i = 1, 2, \dots, n$. A run is accepting if it ends in a final state, $q_{n+1} \in F$. A word is accepted if there is an accepting run for it. The set of accepted words is denoted by $L(\mathcal{B})$.

A *nondeterministic tree automaton* can be presented as a tuple $\mathcal{A} = \langle \Gamma, Q, \delta, F \rangle$ where Γ is the input alphabet, Q is a finite set of states, $F \subseteq Q$ is the set of final states, and δ is a function $Q \times \Gamma \rightarrow 2^{Q^*}$ such that $\delta(p, \sigma)$ is a regular language over Q for every $p \in Q, \sigma \in \Gamma$. A run of \mathcal{A} on a data tree T is a labelling $\lambda : T \rightarrow Q$ such that for every $v \in T$ with n children,

$$\lambda(v_0)\lambda(v_1) \cdots \lambda(v_{n-1}) \in \delta(\lambda(v), \ell_T(v)).$$

If v is a leaf, the condition reduces to $\varepsilon \in \delta(\lambda(v), \ell_T(v))$, which explains the lack of initial states in \mathcal{A} . A run λ is *accepting*, if the root is labelled with a final state, $\lambda(\varepsilon) \in F$. A data tree T is accepted by \mathcal{A} if there is an accepting run for it. The set of trees accepted by \mathcal{A} is denoted by $L(\mathcal{A})$.

Throughout the paper we assume that all states are reachable, i.e., for each state $p \in Q$ there is a tree which evaluates to p . This can be guaranteed by polynomial preprocessing.

In the decision problems involving tree automata we assume that the regular languages $\delta(p, \sigma)$ are given by nondeterministic word automata $\mathcal{B}_{p,\sigma}$, and write $\|\mathcal{A}\|$ for the size of \mathcal{A} including the automata $\mathcal{B}_{p,\sigma}$.

In this case, an *extended run* of \mathcal{A} on T is a run λ together with a labelling κ of T such that whenever $\ell_T(v) = \sigma$ and $\lambda(v) = p$, and v 's children are v_1, v_2, \dots, v_k , it holds that $q_I^{p,\sigma}, \kappa(v_1), \kappa(v_2), \dots, \kappa(v_k)$ is an accepting run of $\mathcal{B}_{p,\sigma}$ on $\lambda(v_1), \lambda(v_2), \dots, \lambda(v_k)$. The κ -label of the root is irrelevant, it can be any state of any $\mathcal{B}_{p,\sigma}$.

A simpler schema language is provided by DTDs. A *document type definition* (DTD) over a labelling alphabet Γ is a pair $D = \langle r, P_D \rangle$, where

- $r \in \Gamma$ is a distinguished root symbol;
- P_D is a function assigning regular expressions over $\Gamma - \{r\}$ to the elements of Γ , usually written as $\sigma \rightarrow e$, if $P_D(\sigma) = e$.

A data tree T *conforms to* a DTD D , denoted $T \models D$, if its root is labelled with r and for each node $s \in T$ the sequence of labels of children of s is in the language of $P_D(\ell_T(s))$. The set of data trees conforming to D is denoted $L(D)$.

Both DTDs and tree automata define languages of data trees. A DTD $D = \langle r, P_D \rangle$ over Γ can be viewed as a tree automaton $\mathcal{A}_D = \langle \Gamma, \delta, \{r\} \rangle$, with $\delta(\sigma, \sigma) = P_D(\sigma)$ and $\delta(\sigma, \tau) = \emptyset$ for $\sigma \neq \tau$. We can think of DTDs as restricted tree automata.

²A different abstraction allows several *attributes* in each node, each attribute storing one data value [1, 3]. Attributes can be easily modelled with additional children.

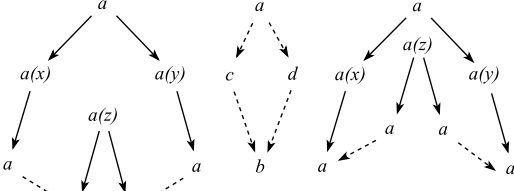


Figure 1: Typical patterns

Patterns. Patterns were originally invented as convenient syntax for conjunctive queries on trees [7, 8, 13]. While most schema mapping research has concentrated on tree-shaped patterns, definable with an XPath-like syntax [1, 3], the full expressive power of conjunctive queries is only guaranteed by DAG-shaped patterns [7, 8]. Indeed, a tree shaped pattern cannot express “there is a path from a to b via c and d ”, but a DAG-shaped pattern or a conjunctive query can (see the middle pattern in Fig. 1). We base our mappings on DAG-shaped patterns, extending the setting used previously.

A pattern π over Γ can be presented as

$$\pi = \langle V, E_c, E_d, E_n, E_f, \ell_\pi, \xi_\pi \rangle$$

where $\langle V, E_c \cup E_d \cup E_n \cup E_f \rangle$ is a finite DAG whose edges are split into child edges E_c , descendant edges E_d , next sibling edges E_n , and following sibling edges E_f , ℓ_π is a partial function from V to Γ , and ξ_π is a partial function from V to the set of variables. The range of ξ_π , denoted $\text{Rg } \xi_\pi$, is the set of variables used by π . By $\|\pi\|$ we mean the size of the underlying DAG.

A data tree $\langle T, \downarrow, \downarrow^*, \rightarrow, \rightarrow^*, \ell_T, \rho_T \rangle$ satisfies a pattern $\pi = \langle V, E_c, E_d, E_n, E_f, \ell_\pi, \xi_\pi \rangle$ under a valuation $\theta : \text{Rg } \xi_\pi \rightarrow \mathbb{N}$, denoted $T \models \pi\theta$, if there exists a homomorphism

$h : \langle V, E_c, E_d, E_n, E_f, \ell_\pi, \xi_\pi \circ \theta \rangle \rightarrow \langle T, \downarrow, \downarrow^*, \rightarrow, \rightarrow^*, \ell_T, \rho_T \rangle$, i.e., a function $h : V \rightarrow T$ such that

- $h : \langle V, E_c, E_d, E_n, E_f \rangle \rightarrow \langle T, \downarrow, \downarrow^*, \rightarrow, \rightarrow^* \rangle$ is a homomorphism of relational structures;
- $\ell_T(h(v)) = \ell_\pi(v)$ for all $v \in \text{Dom } \ell_\pi$; and
- $\rho_T(h(u)) = \theta(\xi_\pi(u))$ for all $u \in \text{Dom } \xi_\pi$.

We write $\pi(\bar{x})$ to express that $\text{Rg } \xi_\pi \subseteq \bar{x}$. For $\pi(\bar{x})$, instead of $\pi\theta$ we usually write $\pi(\bar{a})$, where $\bar{a} = \theta(\bar{x})$. We say that T satisfies π , denoted $T \models \pi$, if $T \models \pi\theta$ for some θ .

Note that we use the usual non-injective semantics, where different vertices of the pattern can be witnessed by the same tree node, as opposed to injective semantics, where each vertex is mapped to a different tree node [10]. Under the adopted semantics patterns are closed under conjunction: $\pi_1 \wedge \pi_2$ can be expressed by the disjoint union of π_1 and π_2 .

Without loss of generality we assume throughout the paper that *siblings have a common parent*, i.e., whenever two vertices of a pattern are connected by a next-sibling or a following sibling edge, there is a vertex connected to both of them by a child edge.

Examples of patterns are given in Fig. 1. Solid and dashed arrows represent E_c and E_d respectively. Figure 2 shows examples of homomorphisms witnessing that the left pattern is satisfied for $x = y = z = 1$ and the right one is satisfied for $x = 3, y = 4, z = 0$. Note that in the left pattern, x and y always take the same value.

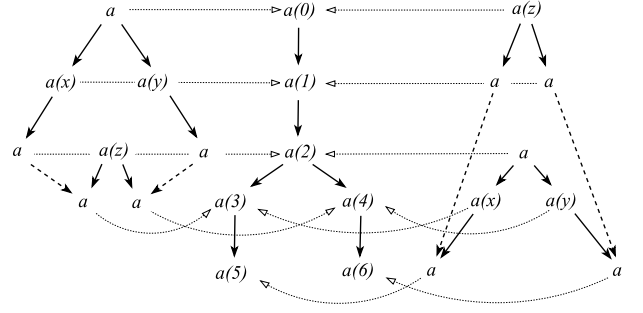


Figure 2: Homomorphisms witness satisfaction

Pattern-based mappings. A (pattern-based) schema mapping $\mathcal{M} = \langle \mathcal{S}_s, \mathcal{S}_t, \Sigma \rangle$ consists of a source schema \mathcal{S}_s , a target schema \mathcal{S}_t , and a set Σ of *source-to-target dependencies* (stds) that relate source and target instances. The source and target schemas are given as tree automata or as DTDs. StdS are expressions of the form:

$$\pi(\bar{x}), \eta(\bar{x}) \rightarrow \pi'(\bar{x}, \bar{y}), \eta'(\bar{x}, \bar{y}), \quad (1)$$

where π, π' are patterns, and η, η' are conjunctions of equalities and inequalities among \bar{x} and \bar{x}, \bar{y} respectively. We assume the usual safety condition: each variable used in η is also used in π .

A pair of trees (T, T') satisfies the std (1) if whenever $T \models \pi(\bar{a})$ and $\eta(\bar{a})$ holds, there is \bar{b} such that $T' \models \pi'(\bar{a}, \bar{b})$ and $\eta'(\bar{a}, \bar{b})$ holds. Given a source $T \in L(\mathcal{S}_s)$, a target $T' \in L(\mathcal{S}_t)$ is called a *solution* for T under \mathcal{M} if (T, T') satisfies all the stds in Σ . We let $\mathcal{M}(T)$ stand for the set of all solutions for T .

Suppose that the source schema is given as a DTD $r \rightarrow a^*b^*$, and the target schema is $r \rightarrow ab^*$. A mapping might be given, e.g. by a single std

$$a(x_1) \rightarrow b \text{ } _ (x_2) \rightarrow a(x_1) \dots \rightarrow b(x_2)$$

where solid and dashed arrows represent E_n and E_f , respectively. Observe that for this mapping every source tree has a solution. On the other hand, if we replace the source DTD with $r \rightarrow (a + b)^*$, some source trees will have no solutions.

In the complexity analysis of computational problems we use various restrictions on mappings. These include imposing tree-structure on patterns (instead of DAG structure), forbidding some axes in patterns, forbidding equality or inequality in stds, and restricting schemas to DTDs.

Evaluation and satisfiability. The evaluation problem is: given a pattern $\pi(\bar{x})$, a tuple \bar{a} , and a tree T , decide if $T \models \pi(\bar{a})$. The complexity of this problem is folklore.

PROPOSITION 2.1. *Data complexity of evaluating patterns is in LOGSPACE and combined complexity is in P.*

PROOF. A pattern is an FO query using relations $\rightarrow^*, \downarrow^*$, which can be computed in LOGSPACE (for \downarrow^* move against the arrows). In consequence, a fixed pattern can be evaluated in LOGSPACE. Dynamic programming gives an algorithm with polynomial combined complexity. \square

A pattern π is satisfiable with respect to an automaton \mathcal{A} if there is a tree $T \in L(\mathcal{A})$ such that $T \models \pi$. The *satisfiability problem* is: given a pattern π and an automaton

\mathcal{A} , decide if π is satisfiable with respect to \mathcal{A} . This problem was shown to be NP-complete for many variants of patterns [2, 4, 8, 14]. As we need the details of the NP-algorithm in some arguments, we sketch it here briefly.

Let h be a homomorphism from π to T . The *support* of h , denoted $\text{supp } h$, is the subtree of T obtained by removing all nodes that cannot be reached from Rgh by going up, left, and right.

LEMMA 2.2. *For each pattern π satisfiable wrt an automaton \mathcal{A} , there exists $T \in L(\mathcal{A})$ and a homomorphism $h : \pi \rightarrow T$ with $|\text{supp } h| \leq 12\|\pi\| \cdot \|\mathcal{A}\|^2$.*

PROOF. Take a tree $T \in L(\mathcal{A})$ satisfying π and let h be a homomorphism from π to T . Divide the nodes of $\text{supp } h$ into four categories: the nodes from the image of h are *red*, the nodes that are not red and have more than one child that is an ancestor of a red node (or is red itself) are *green*, the others are *yellow* if they are ancestors of red nodes, and *blue* otherwise. Let N_{red} , N_{green} , N_{yellow} , and N_{blue} be the numbers of red, green, yellow, and blue nodes.

By definition, $N_{\text{red}} \leq \|\pi\|$. Also $N_{\text{green}} \leq \|\pi\|$: when going bottom-up, each green node decreases the number of subtrees containing a red node by at least one, and since in the root we arrive with one subtree containing a red node, $N_{\text{green}} \leq N_{\text{red}}$. By a pumping argument we may assume that all yellow \downarrow -paths and all blue \rightarrow -paths in $\text{supp } h$ are not longer than $\|\mathcal{A}\|$. The number of (maximal) yellow \downarrow -paths is at most $N_{\text{red}} + N_{\text{green}}$. Hence there are at most $2\|\pi\| \cdot \|\mathcal{A}\|$ yellow nodes. Since all blue nodes are siblings of nodes of other colours, the number of (maximal) blue \rightarrow -paths is at most $2(N_{\text{red}} + N_{\text{green}} + N_{\text{yellow}}) \leq 4\|\pi\| \cdot (\|\mathcal{A}\| + 1)$ and so $N_{\text{blue}} \leq 4\|\pi\| \cdot (\|\mathcal{A}\| + 1)\|\mathcal{A}\|$. Altogether we have at most $2\|\pi\|(\|\mathcal{A}\| + 1)(2\|\mathcal{A}\| + 1) \leq 12\|\pi\| \cdot \|\mathcal{A}\|^2$ nodes. \square

PROPOSITION 2.3. *The satisfiability of patterns is in NP.*

PROOF. By Lemma 2.2, we can guess a homomorphism into a polynomial tree T together with an “almost” run on T : in every leaf we additionally guess a sequence of states to which the missing subtrees should evaluate (by the pumping lemma, the sequence can be linear in the size of the automaton). Since all states are reachable, T can be extended to a tree accepted by the automaton. Checking correctness of the “almost” run, and of the homomorphism is polynomial in the size of the pattern and the tree T , hence it is polynomial. \square

3. ABSOLUTE CONSISTENCY FOR PATTERN-BASED MAPPINGS

As we have mentioned, a schema mapping $\mathcal{M} = \langle \mathcal{S}_s, \mathcal{S}_t, \Sigma \rangle$ is called *absolutely consistent* if every source tree has a solution, i.e., $\mathcal{M}(T) \neq \emptyset$ for every $T \in L(\mathcal{S}_s)$. We are interested in the following decision problem.

PROBLEM: ABCONS
INPUT: Mapping $\mathcal{M} = \langle \mathcal{A}_s, \mathcal{A}_t, \Sigma \rangle$
QUESTION: Is \mathcal{M} absolutely consistent?

Assume for a while that the mapping \mathcal{M} contains a single dependency $\pi(\bar{x}) \rightarrow \pi'(\bar{a}, \bar{b})$. The logical structure of the condition we need to check is: *for every* source tree T *there exists* a target tree T' such that *for every* \bar{a} satisfying $T \models$

$\pi(\bar{a})$ *there exists* \bar{b} satisfying $T' \models \pi'(\bar{a}, \bar{b})$. To turn this into an algorithm we would need to show a bound on the size of trees that need to be considered.

Instead, we will try to change the order of quantifiers to the following: “for every T and every \bar{a} , there exists T' and some witnessing \bar{b} ”. The modified condition can be checked easily in Π_2P . Indeed, what really matters is the equality type of \bar{a} and \bar{b} , so it is enough to chose their entries from a fixed set of linear size. Furthermore, one does not need to guess T and T' explicitly, it is enough to witness their existence. By Lemma 2.2, there exists a polynomial witness.

To justify the reordering of the quantifiers, we would need to show that for every two target trees T_1 and T_2 there is a π' -union tree T , such that whenever $T_1 \models \pi'(\bar{a}, \bar{b})$ or $T_2 \models \pi'(\bar{a}, \bar{b})$, then $T \models \pi'(\bar{a}, \bar{b})$. Without a target schema a π' -union can be obtained by combining T_1 and T_2 under a fresh root node. In the presence of a target schema, however, a π' -union need not exist. For instance, consider the schema $r \rightarrow ab^*$, two trees $r[a(1), b(1), b(2)]$ and $r[a(2), b(3), b(4)]$, and $\pi'(x_1, x_2)$ saying “there is an a -node storing x_1 with a following sibling storing x_2 ”. The first tree satisfies $\pi'(1, 2)$, and the second satisfies $\pi'(2, 4)$, but clearly no tree conforming to the schema can satisfy both $\pi'(1, 2)$ and $\pi'(2, 4)$.

We resolve this difficulty by showing that one can split every schema-defined language of data trees into subsets which are closed under π' -union, and have a description exponential in the size of the schema and π' . Those descriptions will be called *kinds*. In the example above, let L_d consist of all trees conforming to the schema $r \rightarrow ab^*$, which store d in the unique a -node. It is easy to see that each L_d is closed under π' -union: the union tree simply needs to include all the b -nodes of T_1 and T_2 .

A single kind usually cannot provide solutions to all source trees. For instance, in the example from the previous section the source schema is $r \rightarrow a^*b^*$ and $\pi(x_1, x_2)$ says “there is an a -node storing x_1 whose next sibling is a b -node, and some following sibling stores x_2 ”. Clearly, a tree from L_d can only be solution to source trees that store d in the last a -node. Thus source documents have to be split into subsets admitting solutions of a single kind. It turns out that the latter is guaranteed by closure under pattern unions, which means that we can also use kinds to split the source documents.

Based on this we reformulate absolute consistency condition as “for every source kind K there exists a target kind K' such that for every T of kind K and every \bar{a} , there exists T' of kind K' and witnessing \bar{b} ”, which ultimately leads to a $\Pi_2\text{EXP}$ -algorithm.

Kinds. Recall the schema $r \rightarrow ab^*$. As we have noticed, the a -node is the critical area of every tree conforming to this schema. One can easily collect b -nodes from two trees in one tree, but there is always only one a -node. The main idea behind kinds is to distinguish the critical and the non-critical areas of trees. A kind will specify the critical areas entirely, including data values. The non-critical areas will be represented as “holes” in the tree, associated with the corresponding parts of the schema.

We find it most convenient to incorporate into each kind an extended run of the automaton representing the original schema. The run will have holes, just like the underlying tree. Each hole will be associated with a strongly connected component of the automaton, reflecting the intuition that the non-critical areas cover those parts of trees that can be

repeated.

With a tree automaton $\mathcal{A} = \langle \Gamma, Q, \delta, F \rangle$ we associate a graph $G_{\mathcal{A}} = (Q, E)$, where $(p, p') \in E$ iff $Q^* p Q^* \cap \delta(p', \sigma)$ is nonempty for some $\sigma \in \Gamma$. We speak of strongly connected components (SCCs) of \mathcal{A} meaning SCCs of $G_{\mathcal{A}}$. We say that an SCC is *non-trivial* if it contains an edge (it might have only one vertex though). A non-trivial SCC X is *branching* if there exist $p, p_1, p_2 \in X$ such $Q^* p_1 Q^* p_2 Q^* \cap \delta(p, \sigma)$ is nonempty for some σ . If this is not the case, X is *non-branching*. We also work with SCCs of word automata, defined in the natural way.

Let $\mathcal{A} = \langle \Gamma, Q, \delta, F \rangle$ and let the automaton recognizing $\delta(p, \sigma)$ be $\mathcal{B}_{p, \sigma} = \langle Q^{\mathcal{A}}, Q^{p, \sigma}, q_I^{p, \sigma}, \delta^{p, \sigma}, F^{p, \sigma} \rangle$ for $p \in Q, \sigma \in \Gamma$. Let $\mathcal{C}_{p, \sigma}$ and \mathcal{D} be the sets of nontrivial SCCs of $\mathcal{B}_{p, \sigma}$ and \mathcal{A} , respectively. An \mathcal{A} -*kind* is a tree K labelled with elements of

$$\mathcal{D} \cup \bigcup_{p, \sigma} \mathcal{C}_{p, \sigma} \cup \left(Q \times \Gamma \times \mathbb{N} \times \bigcup_{p, \sigma} Q^{p, \sigma} \right)$$

according to certain rules. The rules simply say that the partial run information is consistent. For technical convenience we also demand that \mathcal{D} -labelled nodes have no siblings, and SCC-labelled nodes are leaves, unless their label is a non-branching SCC of \mathcal{A} . The formal rules are as follows:

1. the root is labelled with (p, σ, a, q) such that $p \in F$ (the state q is irrelevant);
2. each internal node v is either labelled with (p, σ, a, q) and has a single child, labelled with an SCC $X \in \mathcal{D}$ containing p , or its sequence of children is labelled with

$$\ell_1, \ell_2, \dots, \ell_k \in \mathcal{C}_{p, \sigma} \cup (Q \times \Gamma \times \mathbb{N} \times Q^{p, \sigma})$$

for some p and σ where

- (a) v is labelled with (p, σ, a, q) for some a, q or with a non-branching $X \in \mathcal{D}$ such that $p \in X$ and $\ell_i = (p_i, \sigma_i, a_i, q_i)$ with $p_i \in X$ for some $1 \leq i \leq k$,
 - (b) $\ell_1 = (p_1, \sigma_1, a_1, q_1)$ and $(q_I^{p, \sigma}, p_1, q_1) \in \delta^{p, \sigma}$,
 - (c) for $1 \leq i < k$ either $\ell_i = (p_i, \sigma_i, a_i, q_i)$, $\ell_{i+1} = (p_{i+1}, \sigma_{i+1}, a_{i+1}, q_{i+1})$ and $(q_i, p_{i+1}, q_{i+1}) \in \delta^{p, \sigma}$, or one of ℓ_i, ℓ_{i+1} is $X \in \mathcal{C}_{p, \sigma}$ and the other is (p', σ', a', q') with $q' \in X$,
 - (d) $\ell_k = (p_k, \sigma_k, a_k, q_k)$ and q_k is accepting;
3. each leaf is labelled with a branching SCC from \mathcal{D} , an element of $\bigcup_{p, \sigma} \mathcal{C}_{p, \sigma}$, or (p, σ, a, q) such that $\mathcal{B}_{p, \sigma}$ accepts the empty word.

To gain better understanding of this notion, let us examine a special case of DTDs admitting only trees of height 1. Recall that we can view DTDs as automata, whose state space is equal to the labelling alphabet Γ , the only final state is the root label, r , and only $\delta(\sigma, \sigma)$ are nonempty. If only trees of height 1 are admitted, for each $\sigma \neq r$, $\delta(\sigma, \sigma) = \{\varepsilon\}$. The only interesting transition is $\delta(r, r) = L(e)$, where $r \rightarrow e$ is the production for the root symbol in the DTD. Since no internal node can have label r , the tree automaton has only trivial SCCs. In consequence, each kind is a tree of height 1, whose root is labelled with (r, r, d, p) for some meaningless p , and some $d \in \mathbb{N}$. The sequence of labels of root's children in K is of the form

$$\alpha_0 X_1 \alpha_1 X_2 \alpha_2 X_3 \alpha_3 \dots X_n \alpha_n,$$

where X_i are SCCs of a word automaton \mathcal{B} recognizing $\delta(r, r)$, α_i are essentially data words decorated with runs of \mathcal{B} starting in X_i (or the initial state if $i = 0$) and finishing in X_{i+1} (or a final state if $i = n$).

The purpose of kinds is to define languages of data trees. A data tree T *agrees with* K if there exists an extended run λ, κ of \mathcal{A} on T and a surjection $f: T \rightarrow K$ such that

- T 's root is mapped to K 's root
- if $v \rightarrow v'$, then $f(v) \rightarrow f(v')$ or $f(v) = f(v') = w$ and w is labelled with an element of $\mathcal{D} \cup \bigcup_{p, \sigma} \mathcal{C}_{p, \sigma}$, and analogously for $v \downarrow v'$,
- $f(v)$ is labelled with $(\lambda(v), \ell_T(v), \rho_T(v), \kappa(v))$, or with an element of $\mathcal{D} \cup \bigcup_{p, \sigma} \mathcal{C}_{p, \sigma}$.

This definition formalizes the intuition of filling the holes: for v labelled with an SCC X , $f^{-1}(v)$ is the area of the tree T filling the hole represented by v in K . If X is a branching SCC of \mathcal{A} or an SCC of $\mathcal{B}_{p, \sigma}$, this area is a subforest of T consisting of a set of subtrees rooted at subsequent children of the same node (the inverse image of the parent of v). If X is a non-branching SCC of \mathcal{A} , then $f^{-1}(v)$ is a subforest like above, minus the subforest formed by the inverse image of v 's descendants.

Let $L(K)$ denote the set of trees agreeing with K . From the definition above and the consistency conditions satisfied by the labelling of each kind it follows that $L(K) \subseteq L(\mathcal{A})$.

Recall once again the example discussed in the previous subsection. The languages L_d cannot be expressed by kinds, but a similar effect can be obtained by taking kinds K_{d_1, d_2, d_3} with the root labelled with (r, r, d_1, q_I) and the root's children labelled with

$$(a, a, d_2, q) \{q\} (b, b, d_3, q),$$

where q_I, q are the states of the two-state automaton recognizing $a b^*$ (transitions are $q_I, a \rightarrow q$ and $q, b \rightarrow q$, the only final state is q).

The language $L(K_{d_1, d_2, d_3})$ consists of trees that store d_1 in the root, d_2 in the a -node, d_3 in the rightmost b -node, and have at least two b -nodes (something has to be mapped to the node labelled with $\{q\}$ in K_{d_1, d_2, d_3}). In order to cover all trees conforming to $r \rightarrow a b^*$, we need to add kinds whose root's children are labelled with $(a, a, d_2, q)(b, b, d_3, q)$ or (a, a, d_2, q) .

The outline of the algorithm presented at the end of the previous subsection suggests that we should be able to provide a small object witnessing that a pattern is satisfied in a tree agreeing with a kind. An argument similar to the proof of Lemma 2.2 and Proposition 2.3 leads to the following bound.

LEMMA 3.1. *For a kind K , a pattern π , and a tuple \bar{a} , satisfiability of $\pi(\bar{a})$ in a tree agreeing with K can be witnessed by an object polynomial in the size of π , the height of K , and the branching of K .*

Recall that our aim is to show that each regular language can be covered by exponential kinds, each defining a language closed under π' -unions. Obviously, the choice of kinds has to take π' into account. For instance, in our running example the kinds with the root's children labelled with $(a, a, d_2, q) \{q\} (b, b, d_3, q)$, $(a, a, d_2, q)(b, b, d_3, q)$ or (a, a, d_2, q) are closed under union for $\pi'(x_1, x_2)$ saying "there is an a -node storing x_1 with a following sibling storing x_2 ", but not

for $\pi''(x_1, x_2)$ saying “there is an a -node storing x_1 whose next sibling stores x_2 ”. For π'' one should use kinds given by

$$\begin{aligned} &(a, a, d_2, q)(b, b, d_3, q)\{q\}(b, b, d_4, q), \\ &(a, a, d_2, q)(b, b, d_3, q)(b, b, d_4, q), \\ &(a, a, d_2, q)(b, b, d_3, q), \text{ and } (a, a, d_2, q), \end{aligned}$$

where d_2, d_3, d_4 range over \mathbb{N} . Intuitively, to ensure closure under π' -union, we need to specify large enough areas of the tree around the holes. This is formalized by the following notion.

For $m \in \mathbb{N}$, we say that K has *horizontal margins of size M* if whenever a node v is labelled with $X \in \mathcal{C}_{p,\sigma}$, there is a sequence of siblings $v_{-M}, v_{-M+1}, \dots, v_M$ with $v_0 = v$ such that for all $0 < |j| \leq M$, v_j is labelled with $(p_j, \sigma_j, a_j, q_j)$ such that $q_j \in X$.

Similarly, K has *vertical margins of size M* if

- whenever v is labelled with a non-branching SCC X , there is a \downarrow -path $v_{-M}, v_{-M+1}, \dots, v_M$ with $v_0 = v$ such that for all $0 < |j| \leq M$, v_j is labelled with $(p_j, \sigma_j, a_j, q_j)$ where $q_j \in X$,
- whenever v is labelled with a branching SCC X , there is a \downarrow -path $v_{-M}, v_{-M+1}, \dots, v_0$ with $v_0 = v$ such that for all $-M \leq j < 0$, v_j is labelled with $(p_j, \sigma_j, a_j, q_j)$ where $q_j \in X$,

Note that margins of size 1 are enforced already by the definition of kinds.

LEMMA 3.2. *For all \mathcal{A} , $T \in L(\mathcal{A})$ and $m, n \in \mathbb{N}$ there is an \mathcal{A} -kind K with vertical margins of size m , horizontal margins of size n , height at most $(2m+1)\|\mathcal{A}\|$, branching at most $(2n+1)\|\mathcal{A}\|$, such that $T \in L(K)$.*

PROOF. Fix an extended accepting run λ, κ of \mathcal{A} on T . Let \tilde{T} be a tree obtained from T by labelling each node v with $(\lambda(v), \ell_T(v), \rho_T(v), \kappa(v))$. Prune \tilde{T} putting SCC labelled nodes as stubs according to the following rules.

First process the SCCs of $\mathcal{B}_{p,\sigma}$. Each X in which the run stays for at least $2n+1$ steps becomes a $\mathcal{C}_{p,\sigma}$ -node, and the first and last n steps of the run spent in X become the horizontal margins around this node. The subtrees rooted at the margin nodes are kept, but the subtrees rooted at the remaining nodes in the middle are lost together with their roots.

Next, deal with the SCCs of \mathcal{A} . Each maximal path P on which the automaton stays in a branching SCC X is cut off at depth $m+1$; under the freshly obtained leaf put a single child labelled with X .

For non-branching X only consider paths of length at least $2m+1$. Cut off at depth $m+1$, add a fresh X node, and under it put all subtrees originally rooted at the children of the $(m+1)$ st node of the path, counting from the bottom.

Let K be the resulting tree. By construction $T \in L(K)$. The bounds on the height and branching of K follow from the observation that no SCC occurs more than once in a sequence of children, or a branch. \square

The purpose of margins is to isolate the regions of the tree filling the holes from each other, and from the fixed part beyond the margins. In this way realizations of patterns touching the filling regions can be rearranged according to our needs. We use this in the lemma below, which constitutes

the technical core of the upper bounds for pattern-based mappings.

LEMMA 3.3. *Let $\pi(\bar{x})$ be a pattern, and let K be an \mathcal{A} -kind with horizontal and vertical margins of length $\|\pi\|$. For all $T_1, T_2, \dots, T_n \in L(K)$ there is $T \in L(K)$ such that $T \models \pi(\bar{a})$ whenever $T_i \models \pi(\bar{a})$ for some $i \in \{1, 2, \dots, n\}$.*

The proof of this lemma can be found in the on-line version of the paper [9]. To have a sample of the techniques used, let us consider a simple word example. Fix a word automaton \mathcal{A} with the state space Q . We work with data words consistently decorated with states of \mathcal{A} . Abusing notation, we write $\alpha \in L(K)$ and $\alpha \models \pi(\bar{a})$ for $\alpha \in (\Gamma \times \mathbb{N} \times Q)^*$, when in fact those relations hold for the underlying data words.

Let K be an \mathcal{A} -kind with margins of length $\|\pi\|$, and just one occurrence of an SCC, say X . Hence, $K = \alpha\beta X\beta'\alpha'$ where $\beta, \beta' \in (\Gamma \times \mathbb{N} \times X)^{\|\pi\|}$. Take $\gamma_1, \gamma_2 \in (\Gamma \times \mathbb{N} \times Q)^*$ such that $\gamma_1, \gamma_2 \in L(K)$. Then $\gamma_i = \alpha\beta\delta_i\beta'\alpha'$ for some $\delta_i \in (\Gamma \times \mathbb{N} \times X)^*$. It is easy to find $\eta \in (\Gamma \times \mathbb{N} \times X)^*$ such that $\gamma = \alpha\beta\delta_1\beta'\eta\beta\delta_2\beta'\alpha'$ is in $L(K)$, where $\tilde{\beta}$ is β with the first letter chopped off.

Assume that $\gamma_i \models \pi(\bar{a})$ and let h be a witnessing homomorphism from $\pi(\bar{a})$ to γ_i . If $h(\pi)$ is disjoint from δ_i , partition the vertices of π into those mapped to $\alpha\beta$ and those mapped to $\beta'\alpha'$. Observe that only E_f edges can exist between those two parts of π . It follows easily that $\gamma \models \pi(\bar{a})$. Suppose that $h(\pi)$ is not disjoint from δ_i . Since $|\beta| = |\beta'| = \|\pi\|$, by the pigeonhole principle there exist positions j, j' of the words β, β' which are not in $h(\pi)$. Partition the vertices of π into those mapped to the left of j , between j and j' , and to the right of j' . Again, only E_f edges can connect different parts. It follows that $\gamma \models \pi(\bar{a})$. For example, for $i = 1$ a witnessing homomorphism is obtained by composing h with the injection $g : \alpha\beta\gamma_1\beta'\alpha' \rightarrow \alpha\beta\gamma_1\beta'\eta\gamma_2\beta'\alpha'$ which maps $\alpha\beta\gamma_1\beta'_{\leq j'}$ to the prefix $\alpha\beta\gamma_1\beta'_{\leq j'}$ and $\beta'_{> j'}$ to the suffix $\beta'_{> j'}$, where $\beta'_{\leq j'}$ is the prefix of β' of length j' and $\beta'_{> j'}$ is the remaining suffix.

Algorithm. We start with a simple observation that lets us concentrate on mappings with a single std.

LEMMA 3.4. *Let $\mathcal{M} = \langle \mathcal{S}_s, \mathcal{S}_t, \{\pi_i, \eta_i \rightarrow \pi'_i, \eta'_i \mid i = 1, 2, \dots, n\} \rangle$ such that each std uses different variables. For $I \subseteq \{1, 2, \dots, n\}$, let*

$$\mathcal{M}_I = \left\langle \mathcal{S}_s, \mathcal{S}_t, \left\{ \bigwedge_{i \in I} \pi_i, \bigwedge_{i \in I} \eta_i \rightarrow \bigwedge_{i \in I} \pi'_i, \bigwedge_{i \in I} \eta'_i \right\} \right\rangle.$$

For all T, T' it holds that $T' \in \mathcal{M}(T)$ iff $T' \in \mathcal{M}_I(T)$ for all $I \subseteq \{1, 2, \dots, n\}$. \square

THEOREM 3.5. **ABCONS** is in $\Pi_2\text{EXP}$.

PROOF. We present the algorithm as a two-round game between two players, \forall and \exists . In each round, \forall moves first. Moves are made by the choice of an object of size exponential in $\|\mathcal{M}\|$ during the first round, and polynomial in $\|\mathcal{M}\|$ during the second round. The winning condition, a polynomial time property of the moves, is defined so that \exists wins exactly if \mathcal{M} is absolutely consistent. In the first round, \forall states what *kind* of tree (in the sense defined above) is a counterexample to absolute consistency, while \exists chooses the *kind* of tree that gives solutions to the purported counterexamples. In the

second round, \forall picks a tree of the declared kind and a tuple witnessing that the solutions fail, and \exists tries to respond with a tree and a tuple that would prove \forall wrong.

By Lemma 3.4 we can assume that \mathcal{M} has a single std: \forall can choose \mathcal{M}_I as a part of the first move. Let an instance of ABCONS be $\mathcal{M} = \langle \mathcal{A}_s, \mathcal{A}_t, \{\pi(\bar{x}), \eta(\bar{x}) \rightarrow \pi'(\bar{x}, \bar{y}), \eta'(\bar{x}, \bar{y})\} \rangle$, where $\mathcal{A}_s = \langle \Gamma, Q^s, \delta^s, F^s \rangle$ is the source automaton and $\mathcal{A}_t = \langle \Gamma, Q^t, \delta^t, F^t \rangle$ is the target automaton.

In the first round \forall plays an \mathcal{A}_s -kind K_\forall with horizontal and vertical margins $\|\pi\|$, height at most $(2\|\pi\| + 1)\|\mathcal{A}_s\|$, and branching at most $(2\|\pi\| + 1)\|\mathcal{A}_s\|$. The data values used in K are to represent an equality type, so it is enough to choose them from $\{1, 2, \dots, |K_\forall|\}$.

The response K_\exists of \exists is similar except that \mathcal{A}_s is replaced by \mathcal{A}_t , π is replaced by π' , and some of the nodes can store nulls taken from a fixed set $\{\perp_1, \perp_2, \dots, \perp_{|K_\exists|}\}$. Each null can appear more than once in K_\exists . It is intended to represent “data distinct from whatever appears in the source tree”. (Formally, we partition \mathbb{N} into two infinite sets, \mathbb{N}_s containing data values that are allowed in the source and \mathbb{N}_t containing the data values allowed only on the target side.)

In the second round, \forall chooses a data tuple \bar{a} (without nulls) such that $\eta(\bar{a})$ holds, together with a polynomial object witnessing that $\pi(\bar{a})$ can be realized in a tree agreeing with K_\forall (Lemma 3.1)

\exists then responds with a tuple \bar{b} (possibly including nulls) such that $\eta'(\bar{a}, \bar{b})$ holds, and a polynomial witness that $\pi'(\bar{a}, \bar{b})$ can be realized in a tree agreeing with K_\exists .

A player loses if he fails to make a move complying with the rules. If all moves are made, \exists wins.

It remains to show that \exists has a winning strategy if and only if \mathcal{M} is absolutely consistent.

If \mathcal{M} is not absolutely consistent, \forall 's strategy in the first round is to choose a data tree T for which no solution exists and play K_\forall such that $T \in L(K_\forall)$ (see Lemma 3.2).

If $|\mathcal{M}|$ is absolutely consistent, \exists 's strategy in the first round is to choose K_\exists so that for every tree agreeing with K_\forall there exists a solution agreeing with K_\exists . If such a K_\exists cannot be produced, then there exists a tree agreeing with K_\forall for which there is no solution at all, contradicting absolute consistency. To see this, reason as follows. For each possible response K to K_\forall , let T_K be a tree agreeing with K_\forall for which there is no solution agreeing with K . By Lemma 3.3 there is a tree $T \in L(K_\forall)$ that satisfies $\pi(\bar{a})$ whenever one of the T_K 's satisfies $\pi(\bar{a})$. Since every $T' \in L(\mathcal{A}_t)$ agrees with one of K 's (Lemma 3.2), there is no solution for T .

In the second round, if \mathcal{M} is not absolutely consistent, there is some T agreeing with K_\forall for which there is no solution. \forall 's strategy now is to choose a tuple \bar{a} such that $T \models \pi(\bar{a}), \eta(\bar{a})$, but $\pi'(\bar{a}, \bar{y}), \eta'(\bar{a}, \bar{y})$ cannot be realized in a tree agreeing with K_\exists . Some suitable \bar{a} exists, as otherwise a solution for T could be obtained by an argument similar to the one above.

If $|\mathcal{M}|$ is absolutely consistent, then whatever \bar{a} was played, $\pi'(\bar{a}, \bar{b}), \eta'(\bar{a}, \bar{b})$ can be realized in a tree from $L(K_\exists)$. \exists 's strategy is simply to choose suitable \bar{b} and a witness. \square

Hardness. Consider the following problem:

PROBLEM:	2^n -UNIVERSALITY
INPUT:	nondeterministic Turing machine M , number n in unary
QUESTION:	Does M accept every word of length 2^n in at most 2^n steps?

The problem is obviously in $\Pi_2\text{EXP}$. It can be also shown that it is hard for this class.

LEMMA 3.6. 2^n -UNIVERSALITY is $\Pi_2\text{EXP}$ -complete.

PROOF. Take a $\Pi_2\text{EXP}$ language L . There is a nondeterministic machine M and $k \leq l$ such that $x \in L$ iff for all y such that $|y| \leq 2^{|x|^k}$, M accepts $\langle x, y \rangle$ (M always stops after at most $2^{|x|^l}$ steps). Let $M_{x,k}$ write x on the tape, mark the first $2^{|x|^k}$ positions of the input as relevant, and simulate M on x and the relevant part of the input. Since $k \leq l$, the machine stops after at most $2^{C|x|^l}$ steps for some C independent of x . Hence, $x \in L$ iff $\langle M_{x,k}, 1^{C|x|^l} \rangle$ is a positive instance of 2^n -UNIVERSALITY. \square

We give a reduction from 2^n -UNIVERSALITY to ABCONS for a restricted class of mappings, based on tree-shaped patterns using only vertical navigation. For such patterns we use the following syntax [3]:

$$\begin{array}{ll} \pi ::= \sigma(x)[\lambda] \mid \sigma[\lambda] & \text{patterns} \\ \lambda ::= \varepsilon \mid \pi \mid \|\pi \mid \lambda, \lambda & \text{lists} \end{array} \quad (2)$$

where $\sigma \in \Gamma \cup \{-\}$. That is, a tree pattern is given by its root node and a listing of its subtrees. A subtree can be rooted at a child of the root (corresponding to π in the definition of λ), or its descendant (corresponding to $\|\pi$). The wildcard symbol $-$ is used to denote a vertex without a label. We also write σ for $\sigma[\]$, σ/π for $\sigma[\pi]$, $\sigma\|\pi$ for $\sigma[\|\pi]$, $\sigma[\lambda, \bigwedge_{i=1}^k \pi_i, \lambda']$ for $\sigma[\lambda, \pi_1, \pi_2, \dots, \pi_k, \lambda']$, and similarly for $\sigma(x)$.

THEOREM 3.7. ABCONS is $\Pi_2\text{EXP}$ -hard, even if schemas are non-recursive DTDs, and the only relations available are child and equality on the target side.

PROOF. Let an instance of 2^n -UNIVERSALITY be n , and a Turing machine M with states q_0, q_1, \dots, q_f . W.l.o.g. we assume that q_f is the only final accepting state. For simplicity we assume that the tape alphabet is $0, 1, \flat$.

The idea of the reduction is to encode the input word in the source tree, and the run in the target tree. The run will be encoded as a sequence of 2^n configurations of length 2^n . Additionally, in the source tree we store a linear order of length 2^n , which will be used to address configurations and their cells, and in the target tree we store a specially preprocessed transition relation of M .

The source DTD is given as

$$\begin{array}{l} r \rightarrow \text{ord } q_0 q_1 \dots q_f \perp \text{ zero one blank}, \\ \text{ord} \rightarrow a_1 b_1, \\ a_i, b_i \rightarrow a_{i+1} b_{i+1}, \\ a_n, b_n \rightarrow c, \\ c \rightarrow \text{zero} \mid \text{one} \end{array}$$

with $i = 1, 2, \dots, n-1$ and the target DTD given as

$$\begin{array}{l} r \rightarrow a_1 b_1 (tr)^d, \\ tr \rightarrow st_1 \text{ sym}_1 st_2 \text{ sym}_2 \dots st_6 \text{ sym}_6, \\ a_i, b_i \rightarrow a_{i+1} b_{i+1}, \\ a_{2n}, b_{2n} \rightarrow \text{confnum cellnum st sym}, \end{array}$$

where $i = 1, 2, \dots, 2n - 1$, and d will be defined shortly. Under a_{2n} and b_{2n} nodes we store a configuration number, a cell number, a state, and a tape symbol. The tr nodes store $\hat{\delta}$, the *extended transition relation* of M , describing possible transitions in a window of three consecutive tape cells. Formally, $\hat{\delta} \subseteq (\{q_0, q_1, \dots, q_f, \perp\} \times \{0, 1, b\})^6$, where \perp means “the head is elsewhere”, and $(p_1, \sigma_1, p_2, \sigma_2, \dots, p_6, \sigma_6) \in \hat{\delta}$ iff at most one of p_1, p_2, p_3 is not equal to \perp , and $p_4\sigma_4p_5\sigma_5p_6\sigma_6$ is obtained from $p_1\sigma_1p_2\sigma_2p_3\sigma_3$ by performing a transition of M . In particular, if $p_1 = p_2 = p_3 = \perp$, it is possible that $p_4 \neq \perp$ or $p_6 \neq \perp$. Note that $\hat{\delta}$ can be computed in P. The constant d used in the target DTD is equal to $|\hat{\delta}|$.

Assume for a while that we only need to handle source trees in which all data values are distinct. The c -nodes encode a linear order of length 2^n , and their leaves labelled with *zero* or *one* encode the input word (the data values they store are ignored). The children of the root, q_i, \perp , *zero*, *one*, and *blank*, store data values encoding the states and tape symbols of M . To ensure that the extended transition relation is stored properly on the target side, for each $(q_1, \sigma_1, q_2, \sigma_2, \dots, p_6, \sigma_6) \in \hat{\delta}$ add an std

$$r[q_1(x_1), \hat{\sigma}_1(y_1), q_2(x_2), \hat{\sigma}_2(y_2), \dots, q_6(x_6), \hat{\sigma}_6(y_6)] \longrightarrow \\ \longrightarrow r/tr \left[\bigwedge_{i=1}^6 st_i(x_i), \bigwedge_{i=1}^6 sym_i(y_i) \right]$$

where $\hat{0} = \textit{zero}$, $\hat{1} = \textit{one}$, and $\hat{b} = \textit{blank}$. Note that d different stds are introduced, so all tr -nodes in the target tree are filled according to $\hat{\delta}$.

Now we need to ensure that the target tree encodes an accepting run of M . In the stds we use auxiliary patterns $Begin(x)$, $Begin_\sigma(x)$, $End(x)$, $Succ(x, y)$, and $Succ_3(x, y, z)$, which we define after describing the stds, as well as

$$Cell(x, y, u, v) = _ [confnum(x), cellnum(y), st(u), sym(v)].$$

We use $//$ as abbreviation only. Each use can be replaced with a sequence of $_$ and $/$ operators of suitable length.

To build the first configuration we copy the input word with the head in state q_0 over the first cell,

$$r[Begin_\sigma(x), q_0(u), \sigma(v)] \longrightarrow r// Cell(x, x, u, v), \\ r[Begin(x), //c(y)/\sigma, \perp(u), \sigma(v)] \longrightarrow r// Cell(x, y, u, v),$$

with $\sigma \in \{\textit{zero}, \textit{one}\}$. Transition correctness is ensured by

$$r[Succ(x_0, x_1), Succ_3(y_1, y_2, y_3)] \longrightarrow \\ \longrightarrow r \left[\bigwedge_{i,j} //Cell(x_i, y_j, u_{3i+j}, v_{3i+j}), tr \left[\bigwedge_{i=1}^6 st_i(u_i), \bigwedge_{i=1}^6 sym_i(v_i) \right] \right].$$

Finally, the last configuration needs to be accepting (w.l.o.g. we assume that the accepting state is looping),

$$r[End(x), q_f(u)] \longrightarrow r// Cell(x, y, u, v).$$

It remains to define the auxiliary patterns. With every c -node v we associate the sequence of a 's and b 's on the path leading from the root to v . This sequence is interpreted as a binary number (a read as 0, b read as 1), which is the position of v in the order. The first three patterns are defined easily:

$$Begin(x) = ord/a_1/a_2/\dots/a_n/c(x), \\ Begin_\sigma(x) = ord/a_1/a_2/\dots/a_n/c(x)/\sigma, \\ End(x) = ord/b_1/b_2/\dots/b_n/c(x).$$

The remaining two cannot be defined as single patterns, but can be expressed as disjunctions of patterns. As we only use the auxiliary patterns on the source side of stds, disjunction can be easily eliminated at the cost of multiplying stds.

$$Succ(x, y) = \bigvee_{i=1}^n // _ [a_i/b_{i+1}/b_{i+2}/\dots/b_n/c(x), \\ b_i/a_{i+1}/a_{i+2}/\dots/a_n/c(y)],$$

$$Succ_3(x, y, z) = \\ \bigvee_{i=1}^{n-1} // _ [a_i/b_{i+1}/b_{i+2}/\dots/b_{n-1} [a_n/c(x), b_n/c(y)], \\ b_i/a_{i+1}/a_{i+2}/\dots/a_n/c(z)] \vee \\ \bigvee_{i=1}^{n-1} // _ [a_i/b_{i+1}/b_{i+2}/\dots/a_n/c(x), \\ b_i/a_{i+1}/a_{i+2}/\dots/a_{n-1} [a_n/c(y), b_n/c(z)]] .$$

We claim that the mapping we have just defined is absolutely consistent iff the answer to n -UNIVERSALITY is “yes”. Assume that the mapping is absolutely consistent. Every input word w can be encoded in a source tree using distinct data values. An inductive argument shows that a solution to such a tree encodes an accepting run of M on w . Conversely, if the answer is “yes”, for each source tree S using distinct data values, a solution is obtained from the run of M on the word encoded in the sequence of *zero* and *one* leaves of S . What if S uses some data values more than once? For a function $h : \mathbb{N} \rightarrow \mathbb{N}$ and a tree U , let $h(U)$ be the tree obtained from U by replacing each data value a with $h(a)$. Now, let S' be a tree with the structure identical as S , but using distinct data values, and let h be a function on data values such that $h(S') = S$. By the previously considered case, there is a solution T' for S' . Since our mapping does not use inequality on the target side, nor equality on the source side, $h(T')$ is a solution for $h(S') = S$. \square

Remark. In the reduction above we can remove disjunction from the DTDs at a cost of relaxing restrictions on patterns. We only need to modify the encoding of the 0s and 1s of the input word. One way is to set $c \rightarrow c_1 c_2$, and replace c/\textit{zero} with $c[c_1(z), c_2(z')]$, $z \neq z'$ and c/\textit{one} with $c[c_1(z), c_2(z)]$. Another way is to set $c \rightarrow c_1 c_2^* c_3$ and use next-sibling to distinguish between $c_1 c_2$ and $c_1 c_2^+ c_3$.

4. BOUNDED DEPTH MAPPINGS

We have seen that absolute consistency is highly untractable even for tree-shaped patterns using only vertical axes and non-recursive DTDs with very simple productions. In this section we show that the complexity can be lowered substantially if the height of trees is bounded by a constant. We say that a mapping \mathcal{M} has *depth at most d* if the source and target schema only admit trees of height at most d .

THEOREM 4.1. *ABCONS for mappings of bounded depth is in $\Pi_4\text{P}$.*

PROOF. We claim that the general algorithm presented in Theorem 3.5 has the desired complexity for mappings of bounded depth.

Assume that some \mathcal{A} only accepts trees of height at most d and let K be an \mathcal{A} -kind. Obviously K has height at most

d. (In fact, K contains no nodes labelled with SCCs of \mathcal{A} , as otherwise arbitrarily high trees would agree with K .)

In consequence, the kinds played in the first round have polynomial branching, and bounded depth, hence are polynomial. In the second round, polynomial objects are played. As the correctness of the moves is polynomial, this gives a $\Pi_4\text{P}$ -algorithm. \square

Remark. A small modification of our techniques makes it possible to prove Theorem 4.1 for a more general definition of boundedness: \mathcal{M} has depth at most d if every pattern it uses can only be realized within the initial d levels of every tree conforming to the schema. This includes mappings using patterns starting at the root, that do not use descendant, nor E_c paths of length greater than d .

We now show that the absolute consistency is $\Pi_4\text{P}$ -hard even for word schema mappings, which can be viewed as depth 1 tree mappings. In the word case schemas are regular expressions (or word automata), and stds only use word patterns, i.e., patterns with $E_c = E_d = \emptyset$. In fact, we only need patterns that are disjoint unions of E_n -paths. We write them as sequences of data words with variables, e.g., $a(x)b(y)$, $a(x)$ is a pattern consisting of two E_n -paths.

THEOREM 4.2. *ABCONS is $\Pi_4\text{P}$ -hard for bounded-depth mappings, even if depth is 1, and the only relations available are next sibling, and equality on the target side.*

PROOF. We provide a reduction from TAUTOLOGY for Π_4 quantified propositional formulas. Let

$$\varphi = \forall x_1, x_2, \dots, x_n \exists y_1, y_2, \dots, y_n \bigwedge_{i=1}^m X_i \vee Y_i \vee Z_i$$

with $X_i, Y_i, Z_i \in \{x_j, y_j, u_j, v_j, \bar{x}_j, \bar{y}_j, \bar{u}_j, \bar{v}_j \mid j = 1, \dots, n\}$. Let the source and target schemas be

$$\begin{aligned} \mathcal{S}_s &= (a_1|a'_1)(a_2|a'_2) \dots (a_n|a'_n)ee, \\ \mathcal{S}_t &= eee a_1 a_2 a_2 \dots a_n a_n b_1 b_1 b_2 b_2 \dots b_n b_n (\#ggg)^7. \end{aligned}$$

The source word encodes a valuation of x_1, x_2, \dots, x_n , a_i means that x_i is *true*, a'_i means it is *false*. In e -positions we store values representing *true* and *false*. On the target side, we want to keep a copy of the valuation of x_i 's and a guessed valuation of y_i 's, except this time we use a different coding. The first position labelled with a_i stores the value of x_i , *true* or *false*, and the following position stores the value of the negation of x_i . Similarly, b_i 's store values of y_i . We also want in the target word two copies of *true* and a copy of *false* arranged so as to enable nondeterministic choice between a pair (*true*, *false*) or (*false*, *true*), as well as all triples with at least one entry *true*, which will help us to check that each clause of φ is satisfied.

Let us now describe the stds. First we make sure that values representing true and false are copied properly,

$$e(x)e(y) \longrightarrow e(x)e(y)e(x),$$

for each i translate the a_i/a'_i coding of values of x_i into true/false coding,

$$\begin{aligned} a_i, e(t)e(f) &\longrightarrow a_i(t)a_i(f), \\ a'_i, e(t)e(f) &\longrightarrow a_i(f)a_i(t), \end{aligned}$$

and enforce in the target word all triples with at least one entry *true*,

$$e(t)e(f) \longrightarrow \#g(f)g(f)g(t)\#g(f)g(t)g(f)\# \dots \#g(t)g(t)g(t).$$

Next, we guess a value of y_i for each i ,

$$e(t)e(f) \longrightarrow b_i(t), b_i(f),$$

and ensure that it makes the internal Π_2^P part of φ true for x_1, x_2, \dots, x_n encoded in the source word:

$$\begin{aligned} e(u_1), e(u_2), \dots, e(u_n), e(t)e(f) &\longrightarrow \\ e(u_1)e(\bar{u}_1), e(u_2)e(\bar{u}_2), \dots, e(u_n)e(\bar{u}_n), \\ e(v_1)e(\bar{v}_1), e(v_2)e(\bar{v}_2), \dots, e(v_n)e(\bar{v}_n), \\ a_1(x_1)a_1(\bar{x}_1)a_2(x_2)a_2(\bar{x}_2) \dots a_n(x_n)a_n(\bar{x}_n), \\ b_1(y_1)b_1(\bar{y}_1)b_2(y_2)b_2(\bar{y}_2) \dots b_n(y_n)b_n(\bar{y}_n), \\ g(X_1)g(Y_1)g(Z_1), g(X_2)g(Y_2)g(Z_2), \dots, g(X_m)g(Y_m)g(Z_m) \end{aligned}$$

(the literals X_j, Y_j, Z_j are taken from φ).

The obtained mapping is absolutely consistent iff φ is a tautology. Indeed, if the mapping is absolutely consistent, in particular it has a solution for each source word that uses two different data values in e -positions. By construction, such words have solutions iff φ is a tautology. If the data values in e -positions are equal, the stds are satisfied trivially. \square

Remark. Disjunction can be eliminated from the source schema at a cost of allowing data comparisons on the source side. To achieve this, replace $(a_i|a'_i)$ with $a_i a_i$, and encode the truth value as (in)equality of the two data values.

5. BUILDING SOLUTIONS

In previous sections we examined existence of solutions for all source trees. The task of this section is to build a solution for a given tree. Techniques used in previous sections give us a procedure with polynomial data complexity; the combined complexity is higher.

THEOREM 5.1. *For a mapping \mathcal{M} and a source tree T one can build a solution (or determine that it does not exist) in EXPSPACE in general, in PSPACE if the mapping has bounded depth, and in P if the mapping is fixed.*

PROOF. As a first step we remove stds whose source sides are not satisfied in T . (Checking if π is satisfied in T can be done in NP, and for fixed π in P.) The remaining stds can be merged into one just like in Lemma 3.4.

By Lemma 3.2, if there is a solution to T , it agrees with one of the possible kinds which \exists can play in the first round of the game described in the proof of Theorem 3.5. Note that the size of those kinds only depends on \mathcal{M} (single exponential and polynomial for bounded depth), and so can be viewed as fixed. The kinds only use data values from T and nulls from a set independent of T (single exponential in $\|\mathcal{M}\|$, polynomial for bounded depth). Therefore, the number of possible kinds is polynomial in $|T|$ and the problem amounts to finding a solution agreeing with a given kind K .

First, for each tuple \bar{a} such that $\eta(\bar{a})$ holds and $T \models \pi(\bar{a})$, we need to find a tuple \bar{b} satisfying $\eta'(\bar{a}, \bar{b})$, whose entries are data values used in T or K , or nulls from $\{\perp'_1, \perp'_2, \dots, \perp'_{|\bar{y}|}\}$, and a tree $T_{\bar{a}}$ agreeing with K such that $T_{\bar{a}} \models \pi'(\bar{a}, \bar{b})$. Clearly, the size of $T_{\bar{a}}$ does not depend on T . (A pumping argument similar to the one in Lemma 2.2 shows that each

part of $T_{\bar{a}}$ that gets mapped into a node of K labelled with an SCC can have size single exponential in $\|\mathcal{M}\|$, and polynomial for bounded depth.) In consequence, we can find \bar{b} and $T_{\bar{a}}$ by exhaustive search. If for some \bar{a} the search fails, there is no solution to T agreeing with K .

The last step is to merge $T_{\bar{a}}$'s into a single solution satisfying all $\pi'(\bar{a}, \bar{b})$. By the proof of Lemma 3.3, this can be done in time polynomial in the total size of the trees and the kind K . \square

As we have seen in the proof, the solution is at most polynomial in the size of the source tree and exponential in the size of the mapping. Those bounds are tight: the smallest tree conforming to the target DTD (or accepted by the target automaton) might need to be exponential, and a simple copying rule $r//_{-}(x) \rightarrow r//_{-}(x)$ makes the solution at least as large as the source tree. This however does not give matching complexity lower bounds for solution building.

To give more precise bounds let us consider solution existence, a decision version of solution building:

PROBLEM:	SOLEX
INPUT:	mapping \mathcal{M} , source tree T
QUESTION:	Is $\mathcal{M}(T)$ nonempty?

We also examine data complexity of solution existence, i.e., the complexity of the following problem:

PROBLEM:	SOLEX(\mathcal{M})
INPUT:	source tree T
QUESTION:	Is $\mathcal{M}(T)$ nonempty?

THEOREM 5.2. *SOLEX is NEXP-complete in general, and $\Sigma_3\text{P}$ -complete for mappings of bounded depth. SOLEX(\mathcal{M}) is in LOGSPACE.*

PROOF. To get the upper bounds for SOLEX proceed just like in the proof of Theorem 5.1, only instead of examining every possible kind played by \exists , choose it nondeterministically. For $\Sigma_3\text{P}$ multiple stds are eliminated during the game: \exists guesses the stds with source sides satisfied, merges them, and plays a kind for the new mapping; \forall either displays an std omitted by \exists and wins immediately, or continues with the new mapping according to the old rules.

The NEXP lower bound can be obtained by a modification of the reduction described in the proof of Theorem 3.7. The problem we reduce from is: given a nondeterministic Turing machine M and $n \in \mathbb{N}$, does M accept 0^{2^n} in at most 2^n steps. We keep the same mapping, and for the source tree we take a tree whose all c -nodes have a *zero*-child, and whose data values are distinct.

Similarly, modifying the reduction from the proof of Theorem 4.2, we get $\Sigma_3\text{P}$ -hardness for word mappings.

Let us move to SOLEX(\mathcal{M}). Since \mathcal{M} has constant size, by Lemma 3.4 we can assume that \mathcal{M} has a single std, $\mathcal{M} = \langle \mathcal{A}_s, \mathcal{A}_t, \{\pi(\bar{x}), \eta(\bar{x}) \rightarrow \pi'(\bar{x}, \bar{y}), \eta'(\bar{x}, \bar{y})\} \rangle$. The kind played by \exists is constant size and takes data values from a set A of linear size. In consequence, it can be stored in logarithmic space. Instead of guessing it, we can iterate over all possibilities. It remains to check in LOGSPACE if there is a solution of a given kind K .

By Lemma 3.3 it is enough to check if for each \bar{a} such that $\eta(\bar{a})$ holds and $T \models \pi(\bar{a})$ there is a tuple \bar{b} and a tree

T' of the kind K such that $\eta'(\bar{a}, \bar{b})$ holds and $T' \models \pi'(\bar{a}, \bar{b})$. Again, we can iterate over all \bar{a} and \bar{b} with entries from A . For fixed \bar{a}, \bar{b} , $T \models \pi(\bar{a})$ can be checked in LOGSPACE by Proposition 2.1. The remaining tests can be carried out in constant time, modulo suitable encoding of data values in K, \bar{a}, \bar{b} , which can be prepared in LOGSPACE. \square

6. MSO QUERIES

In this section we extend the language of schema mappings, by replacing patterns with formulas of monadic second-order logic (MSO). We prove that for this richer language, all the problems are still decidable. We do not establish the precise complexities.³

The general setup is similar to the one in the previous sections. A single source-to-target dependency is also an implication of the form

$$\pi(\bar{x}), \eta(\bar{x}) \rightarrow \pi'(\bar{x}, \bar{y}), \eta'(\bar{x}, \bar{y}).$$

The formulas are η and η' are boolean combinations of equality constraints on their arguments (this is slightly more general than before, since previously we had conjunctions of equalities and inequalities). The main generalization concerns π and π' . Previously, the tuples of data values selected by π and π' came from patterns. In this section, we allow a richer syntax for π and π' , where formulas of MSO are used.

To avoid confusion, we will explicitly distinguish here between queries that select tuples of nodes of a data tree (which we call *node-selecting queries*) and queries that select tuples of data values (which we call *data-selecting queries*). We use letters α, β for node-selecting queries and π for data-selecting queries. Every boolean combination η of equalities and inequalities on data values can be seen as a data-selecting query. There is also a natural transformation from node-selecting queries to data-selecting queries, defined as follows. When $\alpha(x_1, \dots, x_k)$ is a node-selecting query, then $\pi_\alpha(x_1, \dots, x_k)$ is a data-selecting query that selects a k -tuple of data values d_1, \dots, d_k in a data tree if the query α selects some k -tuple of nodes x_1, \dots, x_k in the data tree such that node x_i has data value d_i for $i = 1, \dots, k$. The idea is that in this section, we will study stds where π is not a pattern, but of the form π_α , where α is a node-selecting query of MSO that does not talk about data values.

MSO formulas. In the previous sections, we used patterns to describe π and π' in the stds. There are some limitations of patterns, however. For instance, every tuple of nodes that matches a pattern will still match it if new nodes are added to the tree. In some natural queries, this is no longer true. For example, we might be interested in node pairs x, y where y is a descendant of x and all nodes on the path from x to y have label a :

$$a(x) \wedge a(y) \wedge x \downarrow^* y \wedge \forall z (x \downarrow^* z \downarrow^* y \Rightarrow a(z)).$$

In this section, we use MSO formulas to define node-selecting queries. We use the following definition of MSO formulas: they can quantify over individual nodes, they can quantify over sets of nodes, they can use predicates to test the label of

³The questions of complexity are most interesting when the queries are given not by MSO formulas, but by automata. If the input contains MSO formulas, then any algorithm will be nonelementary, because satisfiability for MSO formulas is nonelementary.

a node, and they can use binary predicates for the ancestor and following sibling relations. A formula of this logic with k free individual variables defines a k -ary node-selecting query. (We do not use free set variables.) Note that we *do not allow* predicates for data values in MSO logic, and therefore whether or not a tuple of nodes gets selected depends only on the structure of the tree, the labels of the tree, but not the data values. This restriction is very important. (If we allow predicates for data values in MSO logic, satisfiability becomes undecidable.) We also use the name *regular queries* for node-selecting queries defined by MSO formulas.

Regular schema mappings. A *regular schema mapping* \mathcal{M} is a tuple $\langle \mathcal{S}_s, \mathcal{S}_t, \Sigma \rangle$, where \mathcal{S}_s is a source schema, \mathcal{S}_t is a target schema, and Σ is a set of stds of the form

$$\pi_\alpha(\bar{x}), \eta(\bar{x}) \longrightarrow \pi_{\alpha'}(\bar{x}, \bar{y}), \eta'(\bar{x}, \bar{y}) \quad (3)$$

where $\pi_\alpha, \pi_{\alpha'}$ are data-selecting queries obtained from MSO formulas α, α' , and η, η' are boolean combinations of equalities among \bar{x} and \bar{x}, \bar{y} respectively. (The assumption that every variable of α is used by α' is not restrictive, as α' can ignore some variables.) A pair of trees (T, T') satisfies the std (3) if for every tuple of data values \bar{a} selected by both π_α and η in T , there exists a tuple of data values \bar{b} such that $\bar{a}\bar{b}$ is selected by both $\pi_{\alpha'}$ and η' in T' . Given $T \in L(\mathcal{S}_s)$, a tree T' is a *solution* to T iff $T' \in L(\mathcal{S}_t)$ and (T, T') satisfies every std in Σ .

Once again, we would like to emphasize that the MSO formulas α, α' do not depend on the data values in the trees T and T' . The data values are only inspected by η and η' in a very restrictive way: by boolean combinations of equalities.

The goal of this section is to show that the decision problems studied in this paper are decidable also for regular schema mappings.

THEOREM 6.1. *For regular schema mappings* ABCONS *is decidable, SOLEX is decidable, and* $\text{SOLEX}(\mathcal{M})$ *is in P.*

LEMMA 6.2. *Without loss of generality, we can assume that in the regular schema mapping*

$$\mathcal{M} = \langle \mathcal{S}_s, \mathcal{S}_t, \Sigma \rangle,$$

the schemas \mathcal{S}_s and \mathcal{S}_t are trivial (select all documents) and the set Σ contains exactly one std.

Potential. By Lemma 6.2, we assume that in the regular schema mapping there are no source and target schemas and that there is only a single std

$$\pi_\alpha(\bar{x}), \eta(\bar{x}) \longrightarrow \pi_{\alpha'}(\bar{x}, \bar{y}), \eta'(\bar{x}, \bar{y}).$$

Let us write π for the conjunction of queries π_α and η , likewise let us write π' for the conjunction of queries $\pi_{\alpha'}$ and η' . These are both data-selecting queries. Let X be the set of variables in the tuple \bar{x} , and let X' be the set of variables in the tuple $\bar{x}\bar{y}$.

If X is a set of variables, we use the name X -answer set for a (finite) set of tuples of data values indexed by variables from X . Since we assume that our data values are natural numbers, an X -answer set is a finite subset of \mathbb{N}^X . For a document, we use the name *source answer set* for the set of tuples selected by π in the document (this is an X -answer set). Likewise we define the *target answer set*, using the

query π' (this is an X' -answer set). A document T' is a solution to document T if and only if the target answer set in T' , when the tuples are restricted to coordinates from X , includes the source answer set of T .

A *potential* \mathcal{P} over a set of variables X is a family of X -answer sets, i.e. a family of finite subsets of \mathbb{N}^X . If \mathcal{P} is a potential over variables X' and $X \subseteq X'$, then we write $\mathcal{P}|X$ for the potential where the tuples in the answer sets are restricted to X . If \mathcal{P} and \mathcal{Q} are potentials over the same variables, then we write $\mathcal{P} \leq \mathcal{Q}$ if every set of tuples in \mathcal{P} is a subset of some set of tuples in \mathcal{Q} .

For our fixed std, define the *source potential* (call it \mathcal{P}) as the family of its source answer sets (ranging over all possible documents), likewise define the *target potential* (call it \mathcal{Q}). The absolute consistency problem is equivalent to deciding if

$$\mathcal{P} \leq \mathcal{Q}|X.$$

The solution existence problem can also be expressed using potentials. For a given source document T , there exists a solution to T if and only if

$$\mathcal{P}_T \leq \mathcal{Q}|X$$

where \mathcal{P}_T is the potential that contains only one set, namely the source answer set of T .

Our approach to both the absolute consistency and solution existence problem is as follows. First, we develop a language for describing potentials. We show that the \leq relation is decidable for potentials presented in the language. Finally, we prove that the potentials $\mathcal{P}, \mathcal{Q}, \mathcal{P}_T$ can all be expressed in the language.

Potential expressions. Fix a finite set X of variables and a finite set C of constants. A *potential expression* with variables X and constants C is an expression of the form

$$e = \mu_1 | \mu_2 | \dots | \mu_n,$$

where each μ_i is a boolean combination of equalities over variables and constants. We assume that each variable and each constant is used at least once. An example of a potential expression over variables $\{x, y, z\}$ and constants $\{c\}$ is

$$x = c \wedge y = z | (x = y = z)$$

A potential expression e defines a potential \mathcal{P}_e over variables X as follows. This potential contains an answer set $A \subseteq \mathbb{N}^X$ if and only if there is some $\mu_i \in \{\mu_1, \dots, \mu_n\}$, some finite set $D \subseteq \mathbb{N}$ of data values, and some constant valuation $C \rightarrow D$ such that A contains exactly the tuples in D^X that satisfy μ_i . The potential defined by expression in the example contains exactly the answer sets of the form

$$A_{D,c} = \{(c, d, d) : d \in D\}$$

$$B_D = \{(d, d, d) : d \in D\}$$

where D is a finite set of data values and $c \in D$.

Our proof of Theorem 6.1 uses two lemmas. The first lemma says that the order \leq is decidable on potentials given by expressions. The proof is similar to the one of Theorem 3.5.

LEMMA 6.3. *The following problem is decidable (in $\Pi_4\text{P}$): Given e, f , is $\mathcal{P}_e \leq \mathcal{P}_f$?*

The second lemma says that potential expressions can be computed for the potentials that appear when solving the absolute consistency and solution existence problems.

LEMMA 6.4. *Let $\pi_\alpha(\bar{x})$ be a data-selecting query obtained from an MSO query α , and let $\eta(\bar{x})$ be a boolean combination of equalities on data values. Let \mathcal{P} be the potential of the conjunction of π_α and η . One can compute a potential expression e such that $\mathcal{P}_e \equiv \mathcal{P}$.*

The proof of Lemma 6.4 is long and included in the appendix available online [9]. Actually, our proof can be extended to show a more general result. Namely, for every potential \mathcal{P} there exists a potential expression e with $\mathcal{P}_e \equiv \mathcal{P}$ if and only if \mathcal{P} is a potential such that for any answer set $A \in \mathcal{P}$ and any permutation of data values $\rho : \mathbb{N} \rightarrow \mathbb{N}$, the potential \mathcal{P} contains also the answer set

$$\rho A = \{\rho \circ v : v \in A\}.$$

However, we do not need the more general result. On the other hand, we need to be able to compute the pattern expression (and not just know that it exists), and for this we use the assumptions of Proposition 6.4.

From the above two lemmas we get an algorithm for deciding the absolute consistency problem. First, we compute a potential expression e that represents the source potential. Then, we compute a potential expression that represents the target potential. Finally, we test the order. For the solution existence problem, the approach is the same. However, instead of computing the source potential, we only use the subset of the source potential that contains a single answer set, namely the source answer set in the given document. Any potential that contains a single answer set is easily seen to be described by a potential expression, e.g. by using constants for all data values in the answer set. Finally, once the schema mapping is fixed, the solution existence problem is in P thanks to the following lemma.

LEMMA 6.5. *Fix a potential expression e with variables X . The following problem is in P: for an X -answer set A , decide if there exists an X -answer set $B \in \mathcal{P}_e$ with $A \subseteq B$.*

7. FUTURE WORK

A summary of complexity results is shown in Fig. 3. Notably lacking are lower bounds for regular mappings. As we have mentioned, more specialized techniques are likely to bring algorithms with better complexity.

An open direction of more theoretical interest is to explore the limits of decidability. Which kinds of logic in stds guarantee decidability of absolute consistency? Or even more generally, which kinds of dependencies? We believe that our notion of potential can be helpful.

A fundamental concept in data exchange is a universal solution, i.e., a solution that can be mapped homomorphically into every other solution. Universal variants of the problems we have considered are worth exploring. On the other hand, homomorphisms of ordered trees are injective, which makes a universal solution a rare bird. Developing relaxed versions of universal solutions seems an interesting research topic.

Acknowledgments. The authors thank Leonid Libkin for inspiring discussions, and Alin Deutsch and the anonymous referees for helpful comments.

Mappings	ABCONS	SOLEX	SOLEX(\mathcal{M})
bounded-depth	$\Pi_4\text{P}$	$\Sigma_3\text{P}$	in LOGSPACE
pattern-based	$\Pi_2\text{EXP}$	NEXP	in LOGSPACE
regular	decidable	decidable	in P

Mappings	Solution building
fixed	in P
bounded-depth	in PSPACE
pattern-based	in EXPSPACE

Figure 3: Summary of complexity results

The first author was supported by the Future and Emerging Technologies (FET) programme within the Seventh Framework Programme for Research of the European Commission, under the FET-Open grant agreement FOX, number FP7-ICT-233599.

The second and the third author were supported by the *Querying and Managing Navigational Databases* project realized within the Homing Plus programme of the Foundation for Polish Science, cofinanced by the European Union from the Regional Development Fund within the Operational Programme Innovative Economy (“Grants for Innovation”).

The second author was also supported by Polish Ministry of Science and Higher Education grant no. N N201 382234.

8. REFERENCES

- [1] S. Amano, L. Libkin, and F. Murlak. XML schema mappings. In *PODS 2009*, pages 33–42.
- [2] S. Amer-Yahia, S. Cho, L. Lakshmanan, D. Srivastava. Tree pattern query minimization. *VLDB J.* 11 (2002), 315–331.
- [3] M. Arenas and L. Libkin. XML data exchange: consistency and query answering. *JACM*, 55(2):7:1–72, 2008.
- [4] M. Benedikt, W. Fan, F. Geerts. XPath satisfiability in the presence of DTDs. *J. ACM* 55(2): (2008).
- [5] P. Barceló. Logical foundations of relational data exchange. *SIGMOD Record*, 38(1):49–58, 2009.
- [6] P. A. Bernstein and S. Melnik. Model management 2.0: manipulating richer mappings. In *SIGMOD*, 2007.
- [7] H. Björklund, W. Martens, T. Schwentick. Conjunctive query containment over trees. *DBPL’07*, pages 66–80.
- [8] H. Björklund, W. Martens, T. Schwentick. Optimizing conjunctive queries over trees using schema information. *MFCS’08*, pages 132–143.
- [9] M. Bojańczyk, L. Kolodziejczyk, F. Murlak. Solutions for XML Data Exchange. Available at www.mimuw.edu.pl/~fmurlak/papers/bkm10.pdf.
- [10] C. David. Complexity of data tree patterns over XML documents. In *MFCS’08*, pages 278–289.
- [11] R. Fagin, P. G. Kolaitis, R. J. Miller, and L. Popa. Date exchange: semantics and query answering. *TCS*, 336:89–124, 2005.
- [12] R. Fagin, P. G. Kolaitis, L. Popa, and W.-C. Tan. Composing schema mappings: second-order dependencies to the rescue. *ACM TODS*, 30(4):994–1055, 2005.
- [13] G. Gottlob, C. Koch, K. Schulz. Conjunctive queries over trees. *JACM* 53(2):238–272 (2006).
- [14] J. Hidders. Satisfiability of XPath expressions. In *DBPL’03*, pages 21–36.

- [15] P. G. Kolaitis. Schema mappings, data exchange, and metadata management. In *PODS 2005*, pages 61–75.
- [16] R. Miller, M. Hernandez, L. Haas, L. Yan, C. Ho, R. Fagin, and L. Popa. The Clio project: managing heterogeneity. *SIGMOD Record*, 30:78–83, 2001.
- [17] M. Murata, D. Lee, M. Mani, and K. Kawaguchi. Taxonomy of XML schema languages using formal language theory. *ACM Transactions on Internet Technology*, 5(4):1-45, 2005.
- [18] F. Neven. Automata Theory for XML Researchers. *SIGMOD Record* 31(3): 39-46 (2002).

APPENDIX

A. PROOF OF LEMMA 3.1

The witness consists of: (1) a homomorphism from $\pi(\bar{a})$ to a tree T_\forall , (2) an extended run λ, κ of \mathcal{A}_s on T_\forall , for each leaf a (short) sequence of states to which the missing subtrees should evaluate, arbitrary λ -labels in leaves, (3) a function $f: T_\forall \rightarrow K_\forall$ satisfying all three conditions for a witness for agreeing (but not necessarily a surjection), whose image has no leaves labelled with non-branching SCCs and with each node contains all its siblings. Arguing like in Lemma 2.2 we can show that T_\forall can be polynomial in $\|\pi\|$, the height of K and the branching of K . Verifying the witness can be done in time polynomial in the size of K (and the witness).

B. PROOF OF LEMMA 3.3

It suffices to prove the claim for $n = 2$, the rest follows by induction. For $i = 1, 2$ fix an extended run λ_i, κ_i and $f_i: T_i \rightarrow K$ witnessing that $T_i \in L(K)$. By the conditions on f_i , if u is not labelled with an SCC, $f_i^{-1}(u)$ is a single node. Otherwise, it is a forest. The idea of the construction of T is to replace each u in K labelled with an SCC, with a tree or a forest that combines $f_1^{-1}(u)$ and $f_2^{-1}(u)$, and satisfies the restrictions imposed by K . The witnessing function $f: T \rightarrow K$ maps this whole tree or forest to u . On old nodes f is naturally induced.

We write $S.v$ for the subtree of S rooted at v . For $v \rightarrow^* w$, $S.[v, w]$ is the forest of trees rooted at the sequence of siblings from v to w . When a forest F is substituted in S at v , $S.v$ is removed and the roots of F become children of v 's parent.

Pick u_0 labelled with an SCC X of $\mathcal{B}_{p,\sigma}$. Since K has horizontal margins of size $\|\pi\|$, there exist nodes u_i for $0 < |i| \leq \|\pi\|$ labelled with $(p_i, \sigma_i, b_i, q_i)$ such that $q_i \in X$, and $u_j \rightarrow u_{j+1}$ for $-\|\pi\| \leq j < \|\pi\|$. Hence, there is a sequence $v_{-\|\pi\|}^i \rightarrow v_{-\|\pi\|+1}^i \rightarrow \dots \rightarrow v_{-1}^i \rightarrow^* v_1^i \rightarrow v_2^i \rightarrow \dots \rightarrow v_{\|\pi\|}^i$ such that $f_i^{-1}(u_0) = T_i.[v_{-1}^i, v_1^i]$ and $f_i^{-1}(u_j) = \{v_j^i\}$ for $0 < |j| \leq \|\pi\|$. Let γ_i be the sequence of siblings from $v_{-\|\pi\|+1}^i$ to $v_{\|\pi\|}^i$. It follows that $\kappa_i(v_{-\|\pi\|}^i)\kappa_i(\gamma_i)$ is a run of $\mathcal{B}_{p,\sigma}$ on $\lambda_i(\gamma_i)$ starting in $\kappa_i(v_{-\|\pi\|}^i) \in X$ and never leaving X . Let $q_{-1}\alpha, \kappa_1(v_{\|\pi\|}^1)\alpha'$, and $\kappa_2(v_{\|\pi\|}^2)\alpha''$ be (short) runs of $\mathcal{B}_{p,\sigma}$ ending, respectively, in $\kappa_1(v_{-\|\pi\|}^1), \kappa_2(v_{-\|\pi\|}^2)$, and some $q \in X$ such that $(q, p_1, q_1) \in \delta^{p,\sigma}$ (condition 2(c) in the definition of kinds). Let $\beta, \beta', \beta'' \in Q^*$ be the corresponding words. The sequence $q_{-1}\alpha\kappa_1(\gamma_1)\alpha'\kappa_2(\gamma_2)\alpha''q_1$ is a run of on $\beta\lambda_1(\gamma_1)\beta'\lambda(\gamma_2)\beta''p_1$ starting in q_{-1} and ending in q_1 . Since all states in \mathcal{A} are reachable, there exists a forest $F = S_1, S_2, \dots, S_{|\beta|}$ such that S_j evaluates to the j th state in β for each j . Similarly, there are forests F' and F'' for β' and β'' . Concatenating them with $F_i = T_i.[v_{-\|\pi\|+1}^i, v_{\|\pi\|}^i]$ into $FF_1F'F_2F''$ we obtain the forest to substitute at u_0 .

Let us now take u_0 labelled with a branching SCC X .

Since K has vertical margins of size $\|\pi\|$, there exists a path $u_{-\|\pi\|}, u_{-\|\pi\|+1}, \dots, u_0$ such that u_j is labelled with $(p_j, \sigma_j, b_j, q_j)$ for some $p_j \in X$ for all $-\|\pi\| < j < 0$. Consequently, there is a path $v_{-\|\pi\|}^i, v_{-\|\pi\|+1}^i, \dots, v_{-1}^i$ in T_i such that $f_i^{-1}(u_j) = \{v_j^i\}$ and $f_i^{-1}(u_0)$ is the set of descendants of v_{-1}^i (excluding v_{-1}^i). As X is branching, there is a tree S and a run assigning p_{-1} to the root of S , and $p_{-\|\pi\|}$ to two incomparable nodes w_1, w_2 . Substituting $T_i.v_{-\|\pi\|}^i$ at w_i in S for $i = 1, 2$ we get the tree to substitute at u_0 .

The last case is u_0 labelled with a non-branching X . There exists a path $u_{-\|\pi\|}, u_{-\|\pi\|+1}, \dots, u_{\|\pi\|}$ in K (visiting u_0) such that u_j 's label is $(p_j, \sigma_j, b_j, q_j)$ with $p_j \in X$ for $0 < |j| \leq \|\pi\|$, as well as paths $v_{-\|\pi\|}^i, v_{-\|\pi\|+1}^i, \dots, v_{-1}^i$ and $v_0^i, v_1^i, v_2^i, \dots, v_{\|\pi\|}^i$ in T_i such that v_0^i is a descendant of v_{-1}^i , $f_i^{-1}(u_j) = \{v_j^i\}$ for $0 < |j| \leq \|\pi\|$, and $f_i^{-1}(u_0)$ is the set of those descendants of v_{-1}^i that are not descendants of v_0^i (excluding v_{-1}^i , including v_0^i). There exist trees S_1, S_2, S_3 , corresponding runs $\lambda^1, \lambda^2, \lambda^3$, nodes w_1, w_2, w_3 , and a state $p \in X$ such that $\lambda^1(\varepsilon) = p_{-1}, \lambda^1(w_1) = p_{-\|\pi\|}, \lambda^2(\varepsilon) = p_{\|\pi\|}, \lambda^2(w_2) = p_{-\|\pi\|}, \lambda^3(\varepsilon) = p_{\|\pi\|}, \lambda^3(w_3) = p$ and the children of u_0 are labelled with elements of $\mathcal{C}_{p,\sigma} \cup Q \times \Gamma \times \mathbb{N} \times Q^{p,\sigma}$, where $\sigma = \ell_{S_3}(w_3)$ (condition 2 of the definition of kinds). Let $T'_i = T_i.v_{-\|\pi\|}^i$. Informally speaking, we want to insert a “path” of trees $S_1T'_1S_2T'_2S_3$ between u_0 's parent and u_0 's children. First we perform a sequence of substitutions: $[S_3]$ is S_3 with the descendants of w_3 replaced by the descendants of u_0 , $[T'_2S_3]$ is T'_2 with $[S_3]$ substituted at $v_{\|\pi\|}^2$, $[S_2T'_2S_3]$ is S_2 with $[T'_2S_3]$ substituted at w_2 , $[T'_1S_2T'_2S_3]$ is T'_1 with $[S_2T'_2S_3]$ substituted at $v_{\|\pi\|}^1$, $[S_1T'_1S_2T'_2S_3]$ is S_1 with $[T'_1S_2T'_2S_3]$ substituted at w_1 ; and then substitute $[S_1T'_1S_2T'_2S_3]$ at u_0 .

From the construction it follows that T agrees with K with the witnessing function f described in the beginning. Let us see that $T \models \pi(\bar{a})$ whenever $T_i \models \pi(\bar{a})$. Take a witnessing homomorphism $h_i: \pi \rightarrow T_i$. For each node u_0 labelled with an SCC we define a set V_{u_0} where h_i has to be modified. Outside of these sets, h is defined by $\{h(x)\} = f^{-1}(f_i(h_i(x)))$. If $f_i^{-1}(u_0)$ and $h_i(\pi)$ are disjoint, V_{u_0} is empty. When they are not disjoint, there are three cases.

Assume first that u_0 is labelled with an SCC of $\mathcal{B}_{p,\sigma}$. We keep the symbols from the construction above. Since the image of h_i has cardinality at most $\|\pi\|$, and one of its elements is in $f_i^{-1}(u_0)$, by the pigeonhole principle there exist $-\|\pi\| \leq j' \leq -1$ and $1 \leq j'' \leq \|\pi\|$ such that $v_{j'}$ and $v_{j''}$ are not in the image of h_i . Let $V_{\text{below}}, V_{\text{above}}$ be a partition of π 's vertices into those mapped to $Z_{u_0} = T_i.[v_{j'+1}, v_{j''-1}]$ and the remaining ones. Take vertices $x \in V_{\text{below}}, y \in V_{\text{above}}$. If $(y, x) \in E_c$, then $h_i(y) = v$ (the parent of $v_{j'}$ and $v_{j''}$). The case $(x, y) \in E_c$ is not possible at all. Similarly, $(x, y) \notin E_n, (y, x) \notin E_n, (x, y) \notin E_d$, and edges $E_d(y, x), E_f(x, y), E_f(y, x)$ can only be witnessed by paths visiting $v, v_{j''}$, and $v_{j'}$ respectively. Hence, we can move the image of V_{below} to the copy of Z_{u_0} in T , and we still find witnessing paths for all those edges. Let $V_{u_0} = V_{\text{below}}$.

Next, suppose that u_0 is labelled with a branching SCC. Again, there exists $-\|\pi\| \leq j \leq -1$ such that v_j is not in the image of h_i . Let $V_{\text{below}}, V_{\text{above}}$ be a partition of π 's vertices into those mapped to Z_{u_0} , consisting of the descendants of v_j , and the remaining ones. This time the only possible edge between $x \in V_{\text{below}}$ and $y \in V_{\text{above}}$ is $E_d(y, x)$, and the witnessing path must visit v_j . Hence, the image of V_{below} can be moved to the copy of Z_{u_0} in T , without violating the existence of witnessing paths. Let $V_{u_0} = V_{\text{below}}$.

Finally, let u_0 be labelled with a non-branching SCC. There exist $-\|\pi\| \leq j' \leq -1$ and $1 \leq j'' < \|\pi\|$ such that $v_{j'}$ and $v_{j''}$ are not in the image of h_i . Let $V_{\text{below}}, V_{\text{between}}, V_{\text{above}}$ be a partition of π 's vertices into those mapped to $T_i.v_{j''}$, $Z_{u_0} = T_i.v_{j'} \setminus T_i.v_{j''}$, and elsewhere. The only possible edges between $x \in V_{\text{below}}, y \in V_{\text{between}}$, and $z \in V_{\text{above}}$ are $E_d(z, y), E_d(y, x)$, and $E_d(z, x)$ with the witnessing paths visiting $v_{j'}, v_{j''}$ and both of these, respectively. Consequently, we can move the image of V_{between} to the copy of Z_{u_0} in T , preserving the existence of witnessing paths. Let $V_{u_0} = V_{\text{between}}$.

As the margins of different nodes are disjoint, Z_u and $Z_{u'}$ are either disjoint, or one is contained in the other. The same holds for their inverse images, V_u and $V_{u'}$. On each maximal V_u define h by moving $h_i(V_u)$ to T as described above. Correctness of h follows from the construction.

C. PROOF OF LEMMA 6.3

Recall that an *equality type* \sim on a set of constants and/or variables X is simply an equivalence relation on X . For $s, t \in X$, \sim agrees with an equality $s = t$ if $s \sim t$. Similarly, \sim agrees with $s \neq t$ if $s \not\sim t$. Finally, \sim agrees with a (partial) valuation v on X if $v(s) = v(t) \Leftrightarrow s \sim t$ for all $s, t \in \text{Dom } v$.

Without loss of generality we may assume that $\text{Var } e \subseteq \text{Var } f$. Let $M_e = \text{Var } e$, $N_e = \text{Const } e$, and similarly for f .

We claim that the problem is solved by the following algorithm.

1. Universally guess μ in e and $v: \text{Const } e \rightarrow \{1, 2, \dots, N_e\}$.
2. Existentially guess ν in f and $w: \text{Const } f \rightarrow \text{Rg } v \cup \{\perp_i | i = 1, 2, \dots, N_f\}$.
3. Universally guess a disjunct p in μ and an equality type \sim_1 on $\text{Var } e \cup \text{Const } e$.
4. Existentially guess a disjunct q in ν and an equality type \sim_2 on $\text{Var } f \cup \text{Const } f$.
5. Return YES iff \sim_1 violates p or v , or the following conditions hold:

- (a) \sim_2 agrees with q and w ;
- (b) $\sim_2 \upharpoonright_{\text{Var } e} = \sim_1 \upharpoonright_{\text{Var } e}$;
- (c) $\forall x \in \text{Var } e \forall d \in \text{Const } f \ x \sim_2 d \implies \implies \exists c \in \text{Const } e \ x \sim_1 c \wedge v(c) = w(d)$.

Assume the algorithm says YES for $e = \mu_1 | \mu_2 | \dots | \mu_n$ and $f = \nu_1 | \nu_2 | \dots | \nu_n$. We need to show that for every D, v, i there exist D', v', i' such that $\text{val}_{D,v} \mu_i \sqsubseteq \text{val}_{D',v'} \nu_{i'}$.

Let ν and w be the existential player's response to $\mu = \mu_i$ and v . Let w' be obtained from w by replacing each \perp_i with a fresh value a_{\perp_i} , not used in D . Let D' be the set D extended with all those fresh values, plus a fresh value a_y for each variable y used only in f . It suffices to show that $\text{val}_{D,v} p \sqsubseteq \text{val}_{D',w'} \nu$ for every p , disjunct of α .

Pick a tuple \bar{a} from $\text{val}_{D,v} p$. Let \sim_1 be the equality type induced by \bar{a} and the valuation v . Let q and \sim_2 be the response to p and \sim_1 . By definition \sim_1 does not violate p nor v . Hence, the conditions (a)-(c) hold. We need to find values in D' for the new variables (i.e., the ones used only in f) so as to satisfy \sim_2 .

By condition (a) all constants within a single \sim_2 -class have the same value. By condition (b) all old variables within a single \sim_2 -class have the same value. If some \sim_2 -class contains a constant and an old variable, then by (c) they have

the same value too. Hence, all constants and old variables within a single class have the same value. We can safely propagate this value to all new variables within this \sim_2 -class. If a \sim_2 -class contains only new variables, choose one of them, say y , and set all variables from this class to a_y . Let \bar{b} be the obtained valuation of variables in f . By construction, \bar{b} respects the equality type \sim_2 . By (a), and the definition of w' , \bar{b} also respects q and w' . Hence, $\bar{b} \in \text{val}_{D',w'} q \sqsubseteq \text{val}_{D',w'} \nu$.

Next, assume that the algorithm says NO. Let μ, v be the "universal" guess witnessing this answer. Let $A \subseteq \mathbb{N} \setminus \text{Rg } v$ with $|A| \geq M_e + N_f$. Define $D = \text{Rg } v \cup A$. We will show that for every D', v' , and ν in f , it holds that $\text{val}_{D,v} \mu \not\sqsubseteq \text{val}_{D',v'} \nu$.

Fix ν, v', D' . Let

$$w(d_i) = \begin{cases} v'(d_i) & \text{if } v'(d_i) \in \text{Rg } v, \\ \perp_{s(v'(d_i))} & \text{if } v'(d_i) \notin \text{Rg } v, \end{cases}$$

where $s: \text{Rg } v' \rightarrow \{1, 2, \dots, N_f\}$ is a fixed injection. Let p, \sim_1 be the "universal" reply to ν, w . We claim that $\text{val}_{D,v} p \not\sqsubseteq \text{val}_{D',v'} \nu$.

Since the answer is NO, \sim_1 must agree with p and v . Let $B = \text{Rg } v' \setminus \text{Rg } v$. Since $|B| \leq N_f$, there is a tuple $\bar{a} \in (D \setminus B)^{M_e}$ that respects \sim_1 and v . Suppose there is a tuple $\bar{b} \in \text{val}_{D',v'} \nu$ that extends \bar{a} . Then there is a disjunct q in ν such that $\bar{b} \in \text{val}_{D',v'} q$. Let \sim_2 be the equality type induced by \bar{b} and v' . By definition of w , \sim_2 agrees with w . Since \sim_2 also agrees with q , it satisfies (a). Furthermore, it clearly satisfies condition (b). To verify (c), suppose that some \sim_2 -class contains an old variable x and a constant d . This means that $v'(d) \in D$ and in consequence $v'(d) \in \text{Rg } v$. Hence, $v'(d) = v(c)$ for some $c \in \text{Const } e$. It follows that $x \sim_1 c$ and $w(d) = v'(d) = v(c)$. This shows that q, \sim_2 is a good answer to p, \sim_1 , which gives a contradiction.

D. PROOF OF LEMMA 6.4

The purpose of this section is to prove Lemma 6.4. The key element is a normal form for MSO queries, as stated in Theorem D.1.

Embellishment. Consider the following query α : "select any tuple (x, y, z) where y is the root and $x = z$ is any node with label a ". Although this query is superficially a ternary query, there is redundancy in the output it produces. First of all, the coordinate y is always the same node, namely the root. Second, the coordinates x and z give the same node. In a certain sense, this query is really just an embellishment on the unary query β , which selects nodes with label a . In this case, the embellishment means that any output x of β is expanded into a triple (x, root, x) . This intuitive notion is defined below.

A *regular constant* is a unary query, which produces at most one node in each input tree. A regular constant can be interpreted as a partial function from trees to nodes, which maps a tree T to the unique node that satisfies the query. Typical regular constants include "the root", or "the leftmost leaf with label a ". A *tuple of regular constants* is defined in a similar way.

Let Δ be a set of regular constants. We use a function

$$\pi: \{1, \dots, k\} \rightarrow \{1, \dots, n\} \cup \Delta.$$

to make a k -ary query $\pi\alpha$ out of an n -ary query α . This function π is called an *embellishment*. Each tuple selected by

the query $\pi\alpha$ will be obtained from a single tuple selected by α . The idea is that π indicates which coordinates of a tuple in $\pi\alpha$ are taken from the tuple in α , and which coordinates are replaced by constants. Formally speaking, if α is an n -ary query, then $\pi\alpha$ is a k -ary query, which selects a tuple (x_1, \dots, x_k) in a tree T whenever there is a tuple (y_1, \dots, y_n) selected by α in T such that for any coordinate $i \in \{1, \dots, k\}$ we have:

- if $\pi(i) \in \{1, \dots, n\}$, then $x_i = y_{\pi(i)}$;
- if $\pi(i) \in \Delta$ then y_i is the interpretation of the constant $\pi(i)$ in the tree T .

Observe that if $\pi\alpha$ selects some tuple in T , then *all* the constants from Δ used by π must be defined in T .

The *domain* of a query is the set of trees where it produces at least one tuple. The domain of a regular query is a regular tree language. The observation on constants above can be rephrased as follows: the domain of $\pi\alpha$ is included in the intersection of the domains of all queries in Δ .

Basis. Let L be a regular tree language. We say that a set Φ of regular queries is a *basis* for a query α over L if over trees from L , the query α can be decomposed as a union of embellishments of queries from Φ and some tuples of regular constants. More formally, there are finitely many tuples of regular constants $\gamma_1, \dots, \gamma_k$, queries $\beta_1, \dots, \beta_n \in \Phi$ and embellishments π_1, \dots, π_n that use constants defined on L , such that for any $T \in L$, the set of tuples $\alpha(T)$ is the same as

$$\gamma_1(T) \cup \dots \cup \gamma_k(T) \cup (\pi_1\beta_1)(T) \cup \dots \cup (\pi_n\beta_n)(T).$$

The set Φ is the basis, while the expression above, including the representations of γ_i and π_i is called the *basis embellishment*.

We say that a set of queries Φ is *unboundedly disjoint in a set of trees* if for every threshold n , there is a tree in the set such that each of the queries in Φ selects at least n tuples with disjoint carriers.

Below is the main result of this section, which says that after splitting into a finite number of cases, every query has a basis that is unboundedly disjoint.

THEOREM D.1. *Let L be a regular tree language and α a regular query. One can compute a partition of L into regular languages L_1, \dots, L_n and regular queries $\alpha_1, \dots, \alpha_n$ such that for each $i \in \{1, \dots, n\}$, α_i is equivalent to α on L_i , and α_i has an unboundedly disjoint basis on L_i , such that each two MSO constants in the basis embellishment are always different or always equal. The bases can also be computed.*

From Theorem D.1 to Lemma 6.4. Before proving Theorem D.1, we first show how it implies Lemma 6.4.

Let π_α , α and η be as in the assumptions of Lemma 6.4. Let X be the variables used by these queries. Apply Theorem D.1 to α , yielding a partition of all documents into regular languages L_1, L_2, \dots, L_n , and queries $\alpha_1, \alpha_2, \dots, \alpha_n$ such that for each $i \in \{1, 2, \dots, n\}$, α_i is equivalent to α on L_i and α_i has an unboundedly disjoint basis on L_i .

Let $\beta = \alpha_i$ and $M = L_i$ for some i . Let

$$\gamma_1 \cup \gamma_2 \cup \dots \cup \gamma_k \cup \pi_1\beta_1 \cup \pi_2\beta_2 \cup \dots \cup \pi_m\beta_m$$

be an unboundedly disjoint basis for β given by Theorem D.1. Every two MSO constants in the presentation above (entries

of $\gamma_1, \gamma_2, \dots, \gamma_k, \pi_1, \pi_2, \dots, \pi_m$) either always return two different nodes, or are the same MSO constant. To each of MSO constants γ assign a constant $c(\gamma)$.

Define $\mu(\beta)$ as $p_1 \vee p_2 \vee \dots \vee p_k \vee q_1 \vee q_2 \vee \dots \vee q_m$, where each p_i is $\eta \wedge \bigwedge_{j=1}^l x_j = c(\gamma_i(j))$, and each q_i is the conjunction of η , equalities $x_j = c(\pi_i(j))$ for j such that $\pi_i(j)$ is an MSO constant, and equalities $x_j = x_{j'}$ for j, j' such that $\pi_i(j) = \pi_i(j')$.

Let $e = \mu(\alpha_1) | \mu(\alpha_2) | \dots | \mu(\alpha_n)$. By the definition of unboundedly disjoint basis, $\mathcal{P}_e \equiv \mathcal{P}$.

Proof of Theorem D.1. The rest of this section is devoted to proving Theorem D.1. The algorithm works as follows. It inputs the query α represented by a basis and a basis embellishment. Initially, the basis is just α and the embellishment is trivial (i.e. it does not modify the query). The algorithm will work inductively, splitting the language L into smaller cases and progressively simplifying the bases involved. The induction is by the position of the basis in the well-founded order defined below.

A set Φ of regular node-selecting queries is defined to be smaller than a set Ψ , written $\Phi < \Psi$, if Φ can be reached from Ψ by applying a finite number of operations of the form: “replace a single query by a finite set of queries with strictly smaller arities”. This ordering is well-founded.⁴

The induction base is when the basis is empty, and the statement of the theorem vacuously holds.

We now do the induction step. Consider any domain V , and any n -ary relation $R \subseteq V^n$, not necessarily the output of a regular query. A *blockade* for R is a set of $X \subseteq V$ such that every tuple from R intersects X (i.e. contains at least one element of X , on some coordinate).

FACT D.2. *Let $R \subseteq V^n$ and $k \in \mathbb{N}$. Either R contains at least k disjoint tuples, or there is a blockade of size kn .*

LEMMA D.3. *Let L be a regular tree language, and Φ a finite set of regular queries. One of the following cases holds, and it is decidable which one:*

- *Unbounded.* The queries in Φ are unboundedly disjoint on L ; or
- *Blockade.* There is a blockade size $K \in \mathbb{N}$, computable from L and Φ , such that in each tree $T \in L$, there is some query $\beta \in \Phi$ such that the set of tuples $\beta(T)$ has a blockade of size at most K .

PROOF. Let K be a number large enough with respect to the representations of L and Φ . Using pumping one can show that if each query in Φ selects $i \geq K$ disjoint tuples in some tree from L , then there is some other tree in L where each query from Φ will select at least $i + 1$ disjoint tuples. The rest follows from Fact D.2. \square

We now proceed to prove Theorem D.1. Let Φ be a basis for α over L . Apply Lemma D.3 to L and Φ . If the case “unbounded” holds, then we are done. Otherwise, case “blockade” holds. Let K be the constant given by the lemma.

For a query $\beta \in \Phi$, define L_β to be the set of trees $T \in L$ such that the set of tuples $\beta(T)$ has a blockade of size at most K . Since having a blockade of a fixed size is MSO

⁴Although well-founded, this ordering has famously bad complexity properties. Probably the proof can be improved in this respect.

definable, the language L_β is a regular tree language. The following lemma shows that if blockades exist, then they can be represented by regular constants.

LEMMA D.4. *Let β be a regular query and K a constant. There is a set Δ of K regular constants such that for any tree T , if there is some blockade of $\beta(T)$ of size at most K , then the interpretation of Δ is such a blockade.*

PROOF. Observe that “ X is a blockade of size at most K for the nodes selected by α ” is a regular property $\varphi(X)$ of a set X of nodes in a tree T . We can use a selector (these are finite trees), to write a formula $\hat{\varphi}(X)$ which selects in any tree exactly one set, call it X_T , that satisfies φ . The i -th constant in Δ is interpreted in a tree T as the i -th node in X_T , in the lexicographic order, or the last node of X_T if X_T has less than i nodes. \square

For each $\beta \in \Phi$, apply Lemma ?? with K being the blockade size given by Lemma D.3, yielding a set Δ_β of regular constants. Let β' be the query obtained from β by restricting output tuples to those that contain a constant from Δ_β on at least one coordinate. Let Φ_β be the set obtained from Φ by replacing β with β' . This new set is smaller than Φ with respect to our induction order. By construction, β and β' have the same output on trees from L_β . It follows that over trees from L_β , the set Φ_β is a basis for α . Splitting L_β again if needed we can assume that every two MSO constants used by the new basis embellishment are different on each tree in L_β , or equal on each tree in L_β . The result follows by applying the induction assumption.