

# Regular n-ary Queries in Trees and Variable Independence

Emmanuel Filiot, Sophie Tison

► **To cite this version:**

Emmanuel Filiot, Sophie Tison. Regular n-ary Queries in Trees and Variable Independence. 5th IFIP International Conference on Theoretical Computer Science, Sep 2008, Milano, Italy. pp.429-443. inria-00274648v2

**HAL Id: inria-00274648**

**<https://hal.inria.fr/inria-00274648v2>**

Submitted on 21 Apr 2008

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

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

---

# Regular $n$ -ary Queries in Trees and Variable Independence

Emmanuel Filiot<sup>1</sup> and Sophie Tison<sup>1</sup>

University of Lille 1, LIFL, UMR 8022 of CNRS, INRIA Lille - Nord Europe  
Emmanuel.Filiot@lifl.fr Sophie.Tison@lifl.fr

**Summary.** Regular  $n$ -ary queries in trees are queries which are definable by an MSO formula with  $n$  free first-order variables. We investigate the variable independence problem – originally introduced for databases – in the context of trees. In particular, we show how to decide whether a regular query is equivalent to a union of cartesian products, independently of the input tree. As an intermediate step, we reduce this problem to the problem of deciding whether the number of answers to a regular query is bounded by some constant, independently of the input tree. As a (non-trivial) generalization, we introduce variable independence w.r.t. a dependence forest between blocks of variables, which we prove to be decidable.

## 1 Introduction

Querying a tree consists of selecting nodes of its domain. This task has received a special interest from the XML community as it is fundamental to information extraction and document transformation. Several formalisms have been proposed to express unary queries [13], but less work has been done on  $n$ -ary queries, *ie* the selection of  $n$ -tuples of nodes [1, 10, 15]. Nevertheless,  $n$ -ary queries are of special interest, for instance to select tuples of the form  $(name, addr, email, phone, fax)$  in an XML document representing a directory. Since the arity of the query can easily get up to 10, efficiency becomes crucial to evaluate ( $n$ -ary) queries.

On the other hand, the notion of *variable independence* has been introduced in the context of (infinite) constraint databases [3, 5, 11, 12]. Query evaluation can be improved considerably when variables are independent. In particular, complexities of many evaluation algorithms on constraint databases are related to some independence between the components of the output tuples. For instance if the query can be decomposed into a Cartesian product of queries of lower arities, then all sub-queries of the product can be evaluated independently. *Orthographic dimension* has been proposed as a measure for variable independence [11]. It corresponds to the size of the largest block of dependent variables. It is shown in [5] that this notion is well-defined, since every pair of decompositions can be intersected into a decomposition of lower maximal block size.

However, in the context of constraint databases, the structure is infinite but fixed. A natural question is whether these results carry over into the context of an infinite number of finite tree structures. The notion of variable independence is also closely related to the representation of the set of answers of an  $n$ -ary query. In particular, if the variables are independent, this allows to represent the set of answers as a Cartesian product of sets of (sub)answers of lower sizes. More generally, *aggregated answers* have been introduced in [14] as a compact representation of the set of answers. Basically, these are multipartite dags such that each part corresponds to a free variable, and the branching structure of the parts is a tree. This branching structure somehow reflects the structure of the input tree. Consider for instance the tree of Fig. 1 representing a directory. Data are omitted in this picture. Consider the ternary query  $q$  which selects all triples  $(x, y, z)$  where  $x$  is labeled **person**,  $y$  is the first name of this person, and  $z$  one of its emails. Once  $x$  is selected, then  $y$  and  $z$  are independent. The set of answers to this query can be represented as a 3-partite graph with sets of vertices  $\{V_x, V_y, V_z\}$ , where there is a directed edge from a node

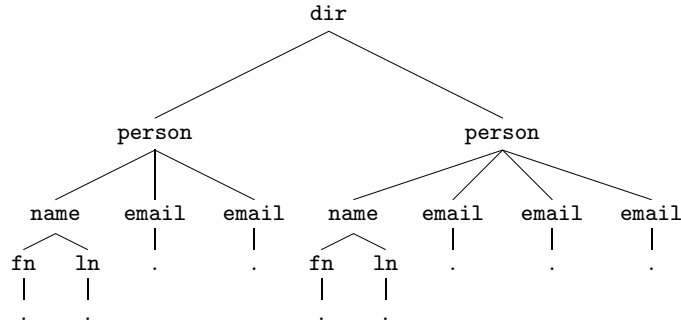


Fig. 1. A tree representing a directory

$u$  of  $V_x$  to a node  $v$  of  $V_y$ , if  $v$  is the name of  $u$ . Similarly, edge relations between  $V_y$  and  $V_z$  correspond to the  $(person, email)$  relations of the input tree. A person may have several emails, making this representation more compact than the set of all answers. Moreover, as argued in [14], this representation keeps the information on how the components of the output tuples are related in the tree. It is also particularly appropriate for post-processing tasks, such as answer searching, answer browsing, statistical computing, answer enumeration, and cascade-style querying. It raises the fundamental question of how compact this representation is. We answer this question by extending variable independence to *relative independence*. In particular, it allows more complex dependencies between variables, as emphasized by the previous example. We measure compactness in the settings of data-complexity, as we want to be independent of the query formalism.

To achieve the formal study of variable independence, we choose *Monadic Second Order Logic* to express  $n$ -ary queries, as it is often used as a yardstick logic in the context of trees [13].

We show that variable independence is decidable for MSO queries in trees, and that a decomposition is computable. We reduce this problem to testing whether a query is bounded, i.e. its number of answers is bounded by some constant independent of the input tree. We prove this problem to be decidable. We show that the notion of orthographic dimension is also well-defined in the context of trees. Finally, we introduce variable independence w.r.t. to a *dependence forest*, which introduces dependencies between blocks of variables.

Note that in the context of trees, a restricted notion of variable independence has been investigated in [15] and proved to be equivalent to non-ambiguity of tree automata.

### Acknowledgments

We are very grateful to the anonymous referees for their valuable comments, to Bruno Courcelle who pointed out some references and to Slawomir Staworko for his careful re-reading.

## 2 Preliminaries

Although XML documents are usually modeled as unranked trees [13], we consider finite binary trees only. All our results can easily be lifted to unranked trees via a binary encoding [13].

**Binary trees** We consider a finite alphabet  $\Sigma$  consisting of symbols ranged over by  $a, f, g$ . A *binary tree*  $t$  over  $\Sigma$  is inductively defined by the following grammar:  $t ::= a \mid a(t, t)$   $a \in \Sigma$ . The set of binary trees over  $\Sigma$  is denoted by  $T_\Sigma$ . The set of *nodes*  $N_t$  of a tree  $t \in T_\Sigma$  is a set of words over  $\{0, 1\}$ . We write  $\epsilon$  for the empty word and  $u.u'$  for the concatenation of  $u$  and  $u'$ . The set  $N_t$  is inductively defined by  $N_a = \epsilon$  and  $N_{a(t_0, t_1)} = \{\epsilon\} \cup \{b.u \mid b \in \{0, 1\}, u \in N_{t_b}\}$ . Nodes  $u \in N_t$  for which there is some  $b \in \{0, 1\}$  such that  $u.b \in N_t$  are called *inner-nodes*. Other nodes are called *leaves*, and the node  $\epsilon$  is called the *root* of  $t$ .

Let  $\Sigma'$  be another finite alphabet. Let  $t \in T_\Sigma$  and  $t' \in T_{\Sigma'}$  be two binary trees such that  $N_t = N_{t'}$ . The *product tree*  $t \times t'$  is the tree over  $\Sigma \times \Sigma'$  inductively defined by:  $a \times b = (a, b)$ , for all

$a \in \Sigma, b \in \Sigma'$ , and  $f(t_0, t_1) \times g(t'_0, t'_1) = (f, g)(t_0 \times t'_0, t_1 \times t'_1)$ , for all  $f \in \Sigma, g \in \Sigma', t_0, t_1 \in T_\Sigma$  and  $t'_0, t'_1 \in T_{\Sigma'}$ . More generally, we can define the product of  $n$  trees modulo associativity of  $\times$ . Trees are also viewed as structures over the signature consisting of the binary successor symbols  $S_0$  and  $S_1$  and the unary symbols  $lab_a$ , for all  $a \in \Sigma$ , interpreted by their intuitive meanings.

**MSO** Monadic Second Order (MSO) logic extends first-order logic with quantification over sets. We consider first-order (resp. second-order) variables ranged over by  $x, y$  (resp.  $X, Y$ ). MSO formulas consists of atomic formulas  $lab_a(x), S_0(x, y), S_1(x, y)$  or  $x \in X$ , and are closed by boolean connectives  $\wedge, \neg$  and quantification  $\exists x, \exists X$ . We let  $t, \rho \models \phi$  denotes the satisfaction relation and say that the formula  $\phi$  holds in the tree  $t$  under the variable assignment  $\rho$ . We refer the reader to [13] for more details about the semantics of MSO.

**Definition 1 (regular  $n$ -ary queries).** Let  $n \geq 0$ . An  $n$ -ary query  $q$  is a mapping from trees  $t \in T_\Sigma$  into subsets of  $N_t^n$ . It is regular (or MSO-definable) if there is an MSO-formula  $\phi(x_1, \dots, x_n)$  with  $n$  free first-order variables such that for all trees  $t$ , we have:  $q(t) = \{(\rho(x_1), \dots, \rho(x_n)) \mid t, \rho \models \phi(x_1, \dots, x_n)\}$

If  $n = 0$  (resp.  $n = 1, n = 2$ ),  $q$  is called Boolean (resp. unary, binary). We also say that  $\phi$  defines  $q$  and denote  $q$  by  $q_\phi$ .

Those queries are called regular since there is a close correspondence between MSO-definable queries and tree automata. In particular, it is well-known that every MSO-definable  $n$ -ary query  $\phi(x_1, \dots, x_n)$  on trees over  $\Sigma$  can be represented as a tree automaton  $A$  over the alphabet  $\Sigma \times \{0, 1\}^n$  [15, 19, 4]. Moreover, we can assume that  $A$  is *canonical*, i.e. for any tree  $t \in L(A)$ , and any  $i \in \{1, \dots, n\}$ , there is exactly one node of  $t$  such that the  $(i + 1)$ -th component of its label is 1 [19, 4]. The following holds: for all its nodes  $u_1, \dots, u_n$  of  $t$ ,  $t \models \phi(u_1, \dots, u_n)$  iff  $t \times \chi_{u_1} \times \dots \times \chi_{u_n} \in L(A)$ , where for all  $i \in \{1, \dots, n\}$ ,  $\chi_{u_i}$  is the tree on  $\{0, 1\}$  such that  $N_t = N_{\chi_{u_i}}$  and all nodes except  $u_i$  are labeled 0.

### 3 Boundedness Properties of Regular Queries

In this section, we prove intermediate results which are useful for Section 4, but might be of independent interest.

Let  $\phi(x_1, \dots, x_n)$  be an MSO formula whose free variables are first-order. We say that  $\phi$  is *bounded* if there is  $K \in \mathbb{N}$  such that for all  $t \in T_\Sigma$ , the number of assignments  $\rho$  of variables  $x_i$ s into  $N_t$  such that  $t, \rho \models \phi$  is bounded by  $K$ .

**Lemma 1.** *Given an MSO formula  $\phi(x_1, \dots, x_n)$  with  $n$  free first-order variables  $x_1, \dots, x_n$ , it is decidable whether  $\phi$  is bounded and, in this case, a bound is computable.*

*Proof.* As said at the end of Section 2,  $\phi$  can be represented as a canonical tree automaton  $A$  over  $\Sigma \times \{0, 1\}^n$ . Now, we can easily transform  $A$  in linear time into a bottom-up transducer  $T_A$  which takes trees  $t_1$  over  $\Sigma$  as inputs and outputs trees  $t_2$  over  $\{0, 1\}^n$  such that  $N_{t_1} = N_{t_2}$  and  $t_1 \times t_2 \in L(A)$ . Since for all trees  $t$ , we have  $|T_A(t)| = |q_\phi(t)|$ , it suffices to test whether  $|T_A(t)|$  is bounded by some constant, which is decidable in polynomial time (in the size of  $T_A$ ). This is called finite valuedness in [17, 18]. Moreover, for all fixed  $k \in \mathbb{N}$ , one can decide in non-deterministic polynomial time whether the number of images by  $T_A$  is greater than  $k$  [17, 18] (with a constant factor which is several exponentials in the size of  $k$ ). Moreover, the bound is lesser than  $2^{2^{p(|T_A|)}}$  for some polynomial  $p$  independent of  $T_A$ . Hence, based on a dichotomy algorithm, one can compute the smallest upper bound. The time complexity however is several exponentials in the size of  $T_A$ .  $\square$

Concerning time complexity, it is known that the size of the tree automaton associated with  $\phi$  might be non-elementary in the size of  $\phi$  [19], making the whole procedure possibly non-elementary. However, if the query is given by a canonical tree automaton  $A$ , testing boundedness becomes polynomial in the size of  $A$ , since testing finite valuedness of a tree transducer can be done in polynomial time [18, 17].

Lemma 1 could also be deduced from a result of [2]. This paper considers an extension of MSO on infinite trees with bounding quantifiers. In particular, for any MSO formula  $\psi(X)$ , the bounding quantified formula  $\mathbb{B}X.\psi(X)$  holds in an infinite tree  $t$ , if there is a bound  $b \in \mathbb{N}$  such that the size of any subset of the set of nodes of  $t$  that satisfies  $\psi(X)$  is bounded by  $b$ .

Two fragments are proved to be decidable: formulas of the form  $\neg\mathbb{B}X.\psi(X)$ , where  $\psi$  is in MSO, and formulas built from arbitrary MSO formulas and  $\mathbb{B}, \exists, \vee$  and  $\wedge$ .

We can easily reduce our problem to satisfiability of some formula in the first fragment. First, boundedness of an  $n$ -ary query reduces to boundedness of all its projections. Hence, we only need to consider unary queries. Now, from a formula  $\psi(x)$  in one free variable, we construct a closed formula  $\gamma$  such that  $\psi(x)$  is bounded iff  $\gamma$  is unsatisfiable. The formula  $\gamma$  is a conjunction  $\gamma = \gamma_1 \wedge \gamma_2$ . The first formula  $\gamma_1$  checks whether the model is a tree (possibly infinite) of the form  $\#(t_1, \#(t_2, \#(t_3, \dots)))$ , for some fresh symbol  $\# \notin \Sigma$  and  $t_1, t_2, t_3, \dots \in T_\Sigma$  are finite trees over  $\Sigma$ . The second formula  $\gamma_2$  has the form  $\neg\mathbb{B}X.\gamma'(X)$ , where  $\gamma'(X)$  is an MSO formula which holds in  $\#(t_1, \#(t_2, \#(t_3, \dots)))$  under some assignment  $\rho$  if there is  $i \geq 1$  such that  $\rho(X)$  corresponds to the set of nodes  $u$  of  $t_i$  such that  $t_i \models \psi(u)$ . The formula  $\gamma'$  is defined by first choosing some node  $x_0$  labeled  $\#$  and then relating  $\psi(x)$  under  $x_0$ .

However, we cannot benefit from this reduction if the query is given by a tree automaton (in term of time complexity).

An *equivalence relation on  $n$ -tuples* is a  $2n$ -ary query, often denoted  $\equiv$ , such that for all trees  $t$ ,  $\equiv(t)$  is an equivalence relation on  $N_t^n$ . We let  $\equiv_t$  stands for  $\equiv(t)$ . It is regular if  $\equiv$  is regular. We say that  $\equiv$  is of bounded index if for all trees  $t$ , the number of  $\equiv_t$ -equivalence classes is bounded by some constant which does not depend on the tree. We can define a regular query which selects the minimal representatives of the equivalence classes, for some MSO-definable order on tuples. Hence, as a corollary of Lemma 1, we get:

**Corollary 1 (bounded index property).** *Let  $\equiv$  be a regular equivalence relation on  $n$ -ary tuples. It is decidable whether  $\equiv$  is of bounded index.*

*Proof.* Let  $t$  be a tree, we define a total order  $\leq_t$  on  $N_t^n$ . It suffices to start from a total order on  $N_t$  and to extend it to a lexicographic order  $\leq_t$  on  $N_t^n$ . Take for instance the lexicographic order on words over  $\{0, 1\}$  which is a total order on  $N_t$ . We can easily show that the query  $t \mapsto \leq_t$  is regular. Now, we define the  $n$ -ary query  $q_{min} : t \mapsto \{\bar{u} \mid \forall \bar{u}' \in N_t^n, \bar{u} \equiv_t \bar{u}' \implies \bar{u} \leq_t \bar{u}'\}$ . The query  $q_{min}$  is regular. Finally it suffices to verify boundedness of  $q_{min}$ , which is decidable by Lemma 1.  $\square$

### Beyond Trees

Deciding boundedness of an MSO formula can be done for classes of structures which are images of a regular set of trees by an MSO-transduction. We refer the reader to [6] for a definition of MSO-transductions. Given a set of integers  $I \subseteq \mathbb{N}$ , we say that  $I$  is *linear* if there are integers  $\alpha_0, \dots, \alpha_n \in \mathbb{N}$  such that  $I = \{\alpha_0 + \sum_{i=1}^n \alpha_i x_i \mid x_1, \dots, x_n \in \mathbb{N}\}$ . It is *semi-linear* if it is a finite union of linear sets. Let  $\sigma$  be a signature and  $C$  be a class of  $\sigma$ -structures,  $\theta$  an MSO-transduction from binary trees to  $\sigma$ -structures such that  $C$  is the image by  $\theta$  of a set of binary trees. In [6], Courcelle proves<sup>1</sup> that given an MSO formula  $\phi(X)$  over the signature  $\sigma$  with one free second-order variable  $X$ , the set  $\{\#\rho(X) \mid M, \rho \models \phi(X), M \in C\}$  is semi-linear, and we can compute the coefficients of the polynomials if the transduction  $\theta$  is known ( $\#\rho(X)$  denotes the cardinality of  $\rho(X)$ ). This is the case for instance for the class of graphs of clique-width less than  $k$ , for any fixed  $k$  [8]. To decide boundedness of an MSO formula  $\phi(x)$  (where  $x$  is first-order), it suffices to compute the above coefficients for the formula  $\Phi(X) = \forall x, x \in X \leftrightarrow \phi(x)$ . The formula  $\phi(x)$  is bounded iff the coefficients  $\alpha_0$  of the linear sets are the unique (possibly) non-null coefficients. Finally, boundedness of a formula  $\phi(x_1, \dots, x_n)$  reduces to boundedness of every projection of  $\phi$  on a single variable  $x_i$ , for all  $i \in \{1, \dots, n\}$ . Hence, boundedness is decidable, for instance, for structures of clique-width less than  $k$ , for any fixed  $k$ , or for unranked trees. However, to decide the bounded index property, we need an MSO-definable total order.

<sup>1</sup> Courcelle proves a more general result where several free variables are allowed

## 4 Variable Independence

The definition of variable independence was originally defined over a fixed structure [3, 12]. We state it over the class of binary trees. We let  $\phi$  be an MSO formula with free variables  $x_1, \dots, x_n$ , and  $P = \{B_1, \dots, B_k\}$  be a partition of  $\{x_1, \dots, x_n\}$ . We write  $\bar{x}_{B_i}$ ,  $i = 1, \dots, k$ , to denote the tuple formed by variables of  $B_i$  given in order. We say that  $\phi$  *conforms to*  $P$ , denoted  $\phi \sim P$ , if  $\phi$  is equivalent to a formula of the form  $\bigvee_{j=1}^N \phi_{j,1}(\bar{x}_{B_1}) \wedge \dots \wedge \phi_{j,k}(\bar{x}_{B_k})$ , where  $N$  is a natural and  $\phi_{j,i}$  are MSO formulas with free variables in  $B_i$ . Note that if we require  $N$  to be equal to 1, the problem becomes easy, since it suffices to test whether  $\phi$  is equivalent to the conjunction of the  $k$  projections of  $\phi$  on the variables from each block  $B_i$ .

W.l.o.g., we assume that free variables of  $\phi$  are ranked in order given by  $B_1, \dots, B_k$ . In other words, we assume  $x_1, \dots, x_n = \bar{x}_{B_1}, \dots, \bar{x}_{B_k}$  (modulo associativity). Now, for any  $i \in \{1, \dots, k\}$ , and any tuples of variables  $\bar{x}, \bar{y}$  such that  $|\bar{x}| = |\bar{y}| = |B_i|$ , we let  $\psi_\phi^i(\bar{x}, \bar{y})$  be the formula defined by:

$$\begin{aligned} \psi_\phi^i(\bar{x}, \bar{y}) &= \forall \bar{x}_{B_1} \dots \forall \bar{x}_{B_{i-1}} \forall \bar{x}_{B_{i+1}} \dots \forall \bar{x}_{B_k} \phi(\bar{x}_{B_1}, \dots, \bar{x}_{B_{i-1}}, \bar{x}, \bar{x}_{B_{i+1}}, \dots, \bar{x}_{B_k}) \\ &\quad \leftrightarrow \\ &\quad \phi(\bar{x}_{B_1}, \dots, \bar{x}_{B_{i-1}}, \bar{y}, \bar{x}_{B_{i+1}}, \dots, \bar{x}_{B_k}) \end{aligned}$$

For any  $i \in \{1, \dots, k\}$ , and any tree  $t \in T_\Sigma$ , we let  $R_i^t$  be the binary relation on  $N_t^{|B_i|}$  defined by  $R_i^t = \{(\bar{u}, \bar{v}) \mid t \models \psi_\phi^i(\bar{u}, \bar{v})\}$ . Intuitively,  $\bar{u}$  and  $\bar{v}$  are equivalent if one can substitute  $\bar{u}$  with  $\bar{v}$ , in any tuple selected by  $\phi$  whose  $i$ -th block is  $\bar{u}$ , and conversely.

**Lemma 2.** *Let  $i \in \{1, \dots, k\}$  and  $t \in T_\Sigma$ . The following are true:*

1.  $R_i^t$  is an equivalence relation on  $N_t^{|B_i|}$ ;
2. if  $\phi$  is equivalent to some formula of the form  $\bigvee_{j=1}^N \phi_{j,1}(\bar{x}_{B_1}) \wedge \dots \wedge \phi_{j,k}(\bar{x}_{B_k})$ , then the number of  $R_i^t$ -equivalence classes is bounded by  $2^N$ .

*Proof.* We only prove the second point, as the first is easy. Given some natural  $i \in \{1, \dots, k\}$ , some tree  $t$  and some node tuples  $\bar{u}, \bar{v}$  of length  $|B_i|$ , we let  $\bar{u} \equiv_i^t \bar{v}$  if there is some set  $F \subseteq \{1, \dots, N\}$  (possibly empty) such that for all  $j \in F$ , we have  $t \models \phi_{j,i}(\bar{u}) \wedge \phi_{j,i}(\bar{v})$ , and for all  $j \in \{1, \dots, N\} \setminus F$ , we have  $t \models \neg \phi_{j,i}(\bar{u}) \wedge \neg \phi_{j,i}(\bar{v})$ . We can easily prove that  $\equiv_i^t$  is an equivalence relation on  $N_t^{|B_i|}$  which has at most  $2^N$  equivalence classes. We now prove that  $\equiv_i^t$  is a refinement of  $R_i^t$ , which will be sufficient to conclude. Let  $\bar{u}, \bar{v}$  be two node tuples of length  $|B_i|$  such that  $\bar{u} \equiv_i^t \bar{v}$ . Let  $\bar{w}_1, \dots, \bar{w}_{i-1}, \bar{w}_{i+1}, \dots, \bar{w}_k$  be node tuples. We have  $t \models \phi(\bar{w}_1, \dots, \bar{w}_{i-1}, \bar{u}, \bar{w}_{i+1}, \dots, \bar{w}_k)$  iff  $t \models \bigvee_{j=1}^N \phi_{j,1}(\bar{w}_1) \wedge \dots \wedge \phi_{j,i}(\bar{u}) \wedge \dots \wedge \phi_{j,k}(\bar{w}_k)$  iff (by definition of  $\equiv_i^t$ )  $t \models \bigvee_{j=1}^N \phi_{j,1}(\bar{w}_1) \wedge \dots \wedge \phi_{j,i}(\bar{v}) \wedge \dots \wedge \phi_{j,k}(\bar{w}_k)$  iff  $t \models \phi(\bar{w}_1, \dots, \bar{w}_{i-1}, \bar{v}, \bar{w}_{i+1}, \dots, \bar{w}_k)$ . Hence we get  $(\bar{u}, \bar{v}) \in R_i^t$ .  $\square$

**Lemma 3.** *If for all  $i \in \{1, \dots, k\}$ , there is some  $m_i \in \mathbb{N}$  such that for all trees  $t \in T_\Sigma$ , the number of equivalence classes of  $R_i^t$  is bounded by  $m_i$ , then  $\phi \sim P$ .*

*Proof.* Let  $i \in \{1, \dots, n\}$ . We define a successor relation between equivalence classes, then we introduce formulas  $cl_i^l(\bar{x}), l = 1, \dots, m_i$ , to define the  $l$ -th equivalence class of  $R_i^t$ , in any tree  $t \in T_\Sigma$ .

As already seen in the proof of Corollary 1, there is an MSO-definable total order  $\leq$  on node tuples. We now define a successor relation  $S_\phi^i$  between the minimal representatives (for  $\leq$ ) of the equivalence relation defined by  $\psi_\phi^i$ . Now, let the formula  $min_\phi^i(\bar{x})$  holds if  $\bar{x}$  is the minimal representative of some equivalence class. It can be defined by  $\forall \bar{y}, \psi_\phi^i(\bar{x}, \bar{y}) \rightarrow \bar{x} \leq \bar{y}$ . The relation  $S_\phi^i$  is now defined by the following MSO formula:

$$S_\phi^i(\bar{x}, \bar{y}) = \bar{x} < \bar{y} \wedge min_\phi^i(\bar{x}) \wedge min_\phi^i(\bar{y}) \wedge \neg(\exists \bar{z}, \bar{x} < \bar{z} < \bar{y} \wedge min_\phi^i(\bar{z}))$$

We let  $s_0(\bar{x})$  stand for  $\neg \exists \bar{z}, S_\phi^i(\bar{z}, \bar{x})$  and  $s_l(\bar{x}), l \in \mathbb{N}$  stands for  $\exists \bar{y}, s_{l-1}(\bar{y}) \wedge S_\phi^i(\bar{y}, \bar{x})$ . We now define  $cl_i^l(\bar{x})$  by  $\exists \bar{y}, s_l(\bar{y}) \wedge \psi_\phi^i(\bar{x}, \bar{y})$ , for all  $1 \leq l \leq m_i$ . Intuitively  $s_l(\bar{x})$  holds in  $t$  under some

assignment  $\rho$  if  $\rho(\bar{x})$  is the minimal representative of the  $l$ -th equivalence class of  $R_i^l$ , while  $cl_i^l(\bar{x})$  holds in  $t$  under  $\rho$  if  $\rho(\bar{x})$  belongs to the  $l$ -th equivalence class of  $R_i^l$ .

Finally, we let  $L$  be the set of tuples of naturals  $\bar{l} = (l_1, \dots, l_k)$  such that  $1 \leq l_i \leq m_i$ ,  $i = 1, \dots, k$ , and we denote by  $\beta^{\bar{l}}(\bar{x}_1)$  the formula  $\exists \bar{x}_2 \dots \bar{x}_k, \phi(\bar{x}_1, \bar{x}_2, \dots, \bar{x}_k) \wedge \bigwedge_{j=1}^k cl_j^{l_j}(\bar{x}_j)$ . We let  $\phi^{\bar{l}}$  be the formula  $\beta^{\bar{l}}(\bar{x}_1) \wedge cl_2^{l_2}(\bar{x}_2) \dots \wedge cl_k^{l_k}(\bar{x}_k)$ . We now prove that  $\phi$  is equivalent to  $\bigvee_{\bar{l} \in L} \phi^{\bar{l}}$ .

Let  $t \in T_\Sigma$  and  $\bar{u}_1, \dots, \bar{u}_k$  node tuples of  $t$  such that  $t \models \phi(\bar{u}_1, \dots, \bar{u}_k)$ . For all  $i \in \{1, \dots, k\}$ , we necessarily have  $t \models \psi_\phi^i(\bar{u}_i, \bar{u}_i)$ . Hence, there is some natural  $l_i \in \{1, \dots, m_i\}$  such that  $t \models cl_i^{l_i}(\bar{u}_i)$ ,  $i = 1, \dots, k$ . Let  $\bar{l} = (l_1, \dots, l_k)$ . It is easy to see that  $t \models \phi^{\bar{l}}(\bar{u}_1, \bar{u}_2, \dots, \bar{u}_k)$ .

Conversely, suppose there is some tuple  $\bar{l} = (l_1, \dots, l_k) \in L$  such that  $t \models \phi^{\bar{l}}(\bar{u}_1, \dots, \bar{u}_k)$ . In particular, we have  $t \models \beta^{\bar{l}}(\bar{u}_1)$ , hence there are some node tuples  $\bar{u}'_2, \dots, \bar{u}'_k$ , such that  $t \models \phi(\bar{u}_1, \bar{u}'_2, \dots, \bar{u}'_k)$ , and for each  $i \in \{2, \dots, k\}$ ,  $t \models cl_i^{l_i}(\bar{u}'_i)$ . We now prove by induction that for all  $p \in \{2, \dots, k\}$ , we have  $t \models \phi(\bar{u}_1, \dots, \bar{u}_p, \bar{u}'_{p+1}, \dots, \bar{u}'_k)$ . It is true for  $p = 1$  by hypothesis. Suppose that it is true at rank  $p > 1$ . Since  $t \models cl_{p+1}^{l_{p+1}}(\bar{u}_{p+1})$ , we also have  $t \models \psi_\phi^{p+1}(\bar{u}_{p+1}, \bar{u}'_{p+1})$ . By induction hypothesis, we have  $t \models \phi(\bar{u}_1, \dots, \bar{u}_p, \bar{u}'_{p+1}, \dots, \bar{u}'_k)$ , and by definition of  $\psi_\phi^{p+1}$ , we easily get  $t \models \phi(\bar{u}_1, \dots, \bar{u}_p, \bar{u}_{p+1}, \bar{u}'_{p+2}, \dots, \bar{u}'_k)$ .  $\square$

As a consequence of Lemma 2 and 3, and Corollary 1, we get the main result:

**Theorem 1.** *Given an MSO formula  $\phi$  with free variables  $x_1, \dots, x_n$  and a partition  $P$  of  $\{x_1, \dots, x_n\}$ , it is decidable whether  $\phi \sim P$  holds or not. If it holds, a decomposition of  $\phi$  is computable.*

*Orthographic Dimension* The notion of orthographic dimension has been introduced in [11] to measure the degree of independence between variables, over a fixed database. We define it for any tree structure. Given a formula  $\phi(x_1, \dots, x_n)$ , the *orthographic dimension*  $d_\phi$  of  $\phi$  is defined by  $d_\phi = \min_{P: \phi \sim P} \max_{B \in P} |B|$ .

Theorem 1 gives us a naive algorithm to compute  $d_\phi$ : for each partition  $P$  of  $\{1, \dots, n\}$ , test whether  $\phi \sim P$  and compute  $\max_{B \in P} |B|$ . But as we next show, we can restrict the tests to 2-partitions.

Given two partitions  $P, P'$  of  $\{1, \dots, n\}$ , we write  $P \sqcap P'$  for the refinement of  $P$  and  $P'$ . Formally, we have  $P \sqcap P' = \{B \cap B' \mid B \in P, B' \in P'\} - \emptyset$ . Theorem 1 of [5] states well-definedness of the notion of orthographic dimension, for first-order logic on any vocabulary, over a fixed structure. In particular, it means that there is a unique partition whose largest block is equal to the orthographic dimension, such that the formula conforms to it. However, the proof given in [5] also works if the structure is not fixed. Moreover, it also works as soon as the logic contains first-order quantifiers, negations, and Boolean connectives. Hence, it also proves the following:

**Theorem 2.** *Let  $\phi$  be an MSO formula in  $n$  free first-order variables  $x_1, \dots, x_n$ , and  $P, P'$  be two partitions of  $\{1, \dots, n\}$ . If  $\phi \sim P$  and  $\phi \sim P'$ , then  $\phi \sim P \sqcap P'$ .*

Now, let  $\mathcal{P} = \{P \mid P \text{ is a 2-partition of } \{1, \dots, n\} \text{ and } \phi \sim P\}$ . From Theorem 2, we can deduce that the orthographic dimension of  $\phi$  is the size of the largest block of  $\bigcap_{P \in \mathcal{P}} P$ . Moreover, by Theorem 1, we can compute a decomposition which corresponds to the orthographic dimension.

#### Relation to the Answer Set Representation

If  $\phi \sim P$ , then for any tree  $t \in T_\Sigma$ , we can represent  $q_\phi(t)$  by an aggregated answer of size  $O(n|t|^{d_\phi})$ , computable in time  $O(n|t|^{d_\phi})$  ( $\phi$  is assumed to be fixed). Indeed,  $\phi$  is equivalent to a (computable) formula of the form:  $\bigvee_{i=1}^N \phi_{i,1}(x_{B_1}) \wedge \dots \wedge \phi_{i,k}(x_{B_k})$  for some natural  $N$ . For every  $i \in \{1, \dots, N\}$ , and  $j \in \{1, \dots, k\}$ , we compute an automaton  $A_{i,j}$  over  $\Sigma \times \{0, 1\}^{|B_j|}$  such that we can compute in time  $O(|t|^{d_\phi} |A_{i,j}|)$  the set  $q_{\phi_{i,j}}(t)$  [7]. The answer representation can be identified to the collection of  $k$ -tuples  $(q_{\phi_{i,1}}(t), \dots, q_{\phi_{i,k}}(t))$ .

## 5 Relative Variable Independence

In this section, we generalize variable independence w.r.t. a partition to variable independence w.r.t. a dependence forest on free variables. Consider for instance the tree of Fig. 1, and let  $\phi(x, y, z)$  be an MSO formula, where  $x$  denotes a person,  $y$  its first name and  $z$  one of its emails. Once the interpretation of  $x$  is fixed, then  $y$  and  $z$  are independent. We call this *relative independence*. Let  $T$  be the tree  $x(y, z)$ , we say that  $\phi$  conforms to  $T$ , denoted  $\phi \sim T$ . We next show that relative independence is decidable as a consequence of Theorem 1.

As a slight generalization, we allow dependence forests to specify dependences between sets of variables instead of single variables. For example, if a formula  $\phi(x, y, z, w)$  conforms to some dependence forest  $\{x, y\}(\{z\}, \{w\})$ , it means that once  $x$  and  $y$  are selected, then  $z$  and  $w$  are independent.

Formally, let  $V = \{x_1, \dots, x_n\}$  be a finite set of variables. A *dependence forest*  $F$  over  $V$  is a forest whose nodes are labeled by subsets of  $2^V$  and such that the set of labels occurring in  $F$  form a partition of  $V$ . If  $F$  has only one root, then it is called a *dependence tree*. We often denote by  $\{T_1, \dots, T_k\}$  a dependence forest consisting of the dependence trees  $T_1, \dots, T_k$ , and by  $V'(F)$  a dependence tree consisting of a dependence forest  $F$  rooted by a set  $V' \subseteq V$ . Although we take a graph point of view, we use the same notations as for binary trees to denote the set of nodes of some forest  $F$  over  $V$  and its labeling function, respectively by  $N_F$  and  $lab_F : N_F \rightarrow 2^V$ . Finally, we confuse the set and tuple notations for variables, so that we sometimes write  $V = \bar{x}$ , for some tuple of variables  $\bar{x}$ .

Let  $\phi(x_1, \dots, x_n)$  be an MSO formula in  $n$  free first-order variables  $x_1, \dots, x_n$ . Let  $\mu$  be a mapping from  $N_F$  into MSO formulas. We say that  $\mu$  is *admissible* for  $F$  if for all nodes  $u, v \in N_F$ , if  $u$  is the parent of  $v$ , then  $\mu(v)$  is an MSO formula  $\psi$  whose free variables are  $lab_F(u) \cup lab_F(v)$ ; if  $u$  is the root of  $F$ , then we require that  $\mu(u)$  is an MSO formula whose free variables are  $lab_F(u)$ . If  $\mu$  is admissible, we naturally extend it to an MSO formula  $\mu(F) = \bigwedge_{u \in N_F} \mu(u)$  in free variables  $x_1, \dots, x_n$ .

**Definition 2.** We say that  $\phi$  conforms to  $F$ , denoted  $\phi \sim F$ , if there is a finite sequence  $\mu_1, \dots, \mu_N$  of admissible mappings for  $F$  such that  $\phi$  is equivalent to  $\bigvee_{i=1}^N \mu_i(F)$ .

We prove decidability of relative independence. We start by a base lemma (Lemma 4), for forests of the form  $\bar{x}(\{\bar{y}, \bar{z}\})$ . Then we give a recursive algorithm for the general case, that uses Lemma 4 and Theorem 1.

**Lemma 4.** Given an MSO formula  $\phi(\bar{x}, \bar{y}, \bar{z})$  with free variables  $\bar{x}, \bar{y}, \bar{z}$ , it is decidable whether  $\phi$  is equivalent to some disjunction of the form  $\bigvee_{i=1}^n \alpha_i(\bar{x}, \bar{y}) \wedge \beta_i(\bar{x}, \bar{z})$ , for some natural  $n$ , and MSO formulas  $\alpha_i, \beta_i$  for  $i = 1, \dots, n$ . Moreover if it holds, a disjunction is computable.

*Proof.* Intuitively, we fix the interpretation of  $\bar{x}$  by extending the alphabet with Boolean tuples. This gives a formula  $\phi_{\bar{x}}(\bar{y}, \bar{z})$ , and we test whether  $\phi_{\bar{x}}(\bar{y}, \bar{z}) \sim \{\bar{y}, \bar{z}\}$ .

More formally, we first transform  $\phi(\bar{x}, \bar{y}, \bar{z})$  into  $\phi_{\bar{x}}(\bar{y}, \bar{z})$ , interpreted on trees over the alphabet  $\Sigma \times \{0, 1\}^m$ , where  $m = |\bar{x}|$ , such that the following property holds ( $P_1$ ): for all trees  $t$ , all nodes  $u_1, \dots, u_m \in N_t$ , and all node tuples  $\bar{v}, \bar{w}$ , we have  $t \models \phi(u_1, \dots, u_m, \bar{v}, \bar{w})$  iff  $t \times \chi_{u_1} \times \dots \times \chi_{u_m} \models \phi_{\bar{x}}(\bar{v}, \bar{w})$ , where the trees  $\chi_{u_i}$  are defined at the end of Section 2.

This can be done by repeating exhaustively the following transformation rule on  $\phi$ : replace each atom of the form  $P(\bar{x}_1, x, \bar{x}_2)$ , where  $x$  is the  $i$ -th component of  $\bar{x}$  and is a free occurrence in  $\phi$ , by  $\exists x, \bigvee_{(f, \bar{b}) \in \Sigma \times B_i} lab_{(f, \bar{b})}(x) \wedge P(\bar{x}_1, x, \bar{x}_2)$ , where  $B_i \subseteq \{0, 1\}^m$  is the set of Boolean tuples whose  $i$ -th component is 1 for  $i = 1, \dots, m$ . Hence,  $\phi(\bar{x}, \bar{y}, \bar{z})$  rewrites to some formula  $\phi'_{\bar{x}}(\bar{y}, \bar{z})$ . We define  $\phi_{\bar{x}}$  by  $\phi'_{\bar{x}} \wedge \phi_{can}$ , where  $\phi_{can}$  is a sentence which ensures that all models  $t \in T_{\Sigma \times \{0, 1\}^m}$  of  $\phi_{\bar{x}}$  are canonical (ie all nodes except one have their  $i$ -th component set to 0  $i = 1, \dots, m$ ). We call  $R_1$  the transformation from  $\phi$  to  $\phi_{\bar{x}}$ .

Then it suffices to test whether  $\phi_{\bar{x}}(\bar{y}, \bar{z}) \sim \{\bar{y}, \bar{z}\}$ , which is decidable by Theorem 1. If it holds, then  $\phi_{\bar{x}}(\bar{y}, \bar{z})$  is equivalent to some formula of the form  $\psi_{\bar{x}}(\bar{y}, \bar{z}) = \bigvee_{i=1}^n \alpha_{i, \bar{x}}(\bar{y}) \wedge \beta_{i, \bar{x}}(\bar{z})$ , for some MSO formulas  $\alpha_{i, \bar{x}}, \beta_{i, \bar{x}}$ . We next consider the following transformation rule  $R_2$ : replace each



**Algorithm 1** Testing Relative Independence

---

```

procedure D( $\phi, F$ )
2:   case  $F$  is a leaf or is of the form  $\overline{x}(\overline{y})$ :
      return  $\phi$ 
4:
      case  $F$  is of the form  $\{\overline{x}_1, \dots, \overline{x}_k\}$ :
6:   test whether  $\phi \sim \{\overline{x}_1, \dots, \overline{x}_k\}$  as in the proof of Theorem 1 and return a decomposition. Otherwise
      breaks.

8:   case  $F$  is of the form  $\overline{x}(\overline{y}, \overline{z})$ :
      test whether  $\phi \sim \overline{x}(\overline{y}, \overline{z})$  as in the proof of Lemma 4 and return a decomposition. Otherwise breaks.
10:
      case  $F$  is of the form  $\overline{x}(\overline{y}(F'), F'')$ :
12:    $\overline{z}', \overline{z}'' \leftarrow$  sets of variables occurring in  $F', F''$ 
       $\bigvee_i \alpha_i(\overline{x}, \overline{y}, \overline{z}') \wedge \beta_i(\overline{x}, \overline{z}'') \leftarrow$  D( $\phi, \overline{x}(\{\overline{y} \cup \overline{z}', \overline{z}''\})$ )
14:   return  $\bigvee_i$  D( $\alpha_i, \overline{y}(F', \overline{x})$ )  $\wedge$  D( $\beta_i, \overline{x}(F'')$ )

16:   case  $F$  is of the form  $\{T_1, \dots, T_k\}$ :
      for  $i \in \{1, \dots, k\}$  do
18:      $\overline{x}_i \leftarrow$  set of variables occurring in  $T_i$ 
       $\bigvee_{i=1}^n \bigwedge_{j=1}^k \alpha_i^j(\overline{x}_j) \leftarrow$  D( $\phi, \{\overline{x}_1, \dots, \overline{x}_k\}$ )
20:   return  $\bigvee_{i=1}^n$  D( $\alpha_i^1, T_1$ )  $\wedge \dots \wedge$  D( $\alpha_i^k, T_k$ )

```

---

atom of the form  $lab_{f, \overline{b}}(x)$  by  $lab_f(x) \wedge \bigwedge_{b_i=1} x_i = x$ , where  $b_i$  denotes the  $i$ -th component of  $\overline{b}$ ,  $i = 1, \dots, m$ . Suppose that  $\overline{x} = x_1, \dots, x_m$ . Applying exhaustively this transformation rule on  $\psi_{\overline{x}}$  leads to a formula  $\psi(\overline{x}, \overline{y}, \overline{z})$  of the form  $\bigvee_{i=1}^n \alpha_i(\overline{x}, \overline{y}) \wedge \beta_i(\overline{x}, \overline{z})$  interpreted on trees over  $\Sigma$ . We have the following property ( $P_2$ ): for all trees  $t$ , all nodes  $u_1, \dots, u_m \in N_t$ , and all node tuples  $\overline{v}, \overline{w}$ , we have  $t \models \psi(u_1, \dots, u_m, \overline{v}, \overline{w})$  iff  $t \times \chi_{u_1} \times \dots \times \chi_{u_m} \models \psi_{\overline{x}}(\overline{v}, \overline{w})$ .

Finally, we prove this algorithm to be correct. Suppose that it returns a decomposition. By combining properties  $P_1$  and  $P_2$ , we can prove correctness of this decomposition. Conversely, suppose that  $\phi$  is equivalent to some formula of the form  $\bigvee_{i=1}^n \alpha_i(\overline{x}, \overline{y}) \wedge \beta_i(\overline{x}, \overline{z})$ . It is easy to see that  $\phi_{\overline{x}}$  is equivalent (here we use canonicity of its models) to  $\bigvee_{i=1}^n \alpha_{i, \overline{x}}(\overline{y}) \wedge \beta_{i, \overline{x}}(\overline{z})$ , where  $\phi_{\overline{x}}$ ,  $\alpha_{i, \overline{x}}(\overline{y})$  and  $\beta_{i, \overline{x}}(\overline{z})$  are obtained by applying the rewrite rule  $R_1$  on respectively  $\phi$ ,  $\alpha_i(\overline{x}, \overline{y})$  and  $\beta_i(\overline{x}, \overline{z})$ . Hence,  $\phi_{\overline{x}} \sim \{\overline{y}, \overline{z}\}$ , and the proof follows since the algorithm of the proof of Theorem 1 is sound.  $\square$

Now, we extend the result of Lemma 4 to full independence forests:

**Theorem 3.** *Given a formula  $\phi$  in free variables  $V = \{x_1, \dots, x_n\}$  and a dependence forest  $F$  over  $V$ , it is decidable whether  $\phi \sim F$  holds or not.*

*Proof.* Consider Algorithm 1. The inputs are a formula  $\phi$  with free variables  $V$  and a dependence forest  $F$  over  $V$ . The symbols  $T_1, \dots, T_k$  denote dependence trees while  $F', F''$  denote (possibly empty) dependence forests.

First note that the algorithm terminates. Indeed, the number of nodes of the forest strictly decreases at each recursive call except for the 4th case when  $F''$  is empty, but in this case the height of the forest strictly decreases.

Now, we can prove (inductively and by using Theorem 1 and Lemma 4 for the basic cases) that if this algorithm returns a formula, then it is a decomposition of the input formula  $\phi$  w.r.t. the input forest  $F$ . It suffices to push up the disjunctive connectives to get a sequence of admissible mappings for  $F$ .

Conversely, let  $\phi$  (resp.  $F$ ) be an input formula (resp. an input dependence forest), such that we have  $\phi \sim F$ . We prove by induction that the algorithm outputs a decomposition. We use the fact that the decomposition is not arbitrary, but has a particular form, derived from the algorithms given in the proofs of Theorem 1 and Lemma 4. The first case is obvious, and the two next

cases have already been proved. First remark that we have the following property (\*): let  $F$  be a dependence forest,  $\gamma_1, \gamma_2$  two formulas,  $\alpha$  a sentence and  $\beta(\bar{x})$  a formula such that  $\bar{x}$  is a label of  $F$ . If  $\gamma_1 \sim F$  and  $\gamma_2 \sim F$ , we have  $\gamma_1 \vee \gamma_2 \sim F$ ,  $\gamma_1 \wedge \gamma_2 \sim F$ ,  $\neg\gamma_1 \sim F$ ,  $\alpha \wedge \gamma_1 \sim F$ , and  $\beta(\bar{x}) \wedge \gamma_1 \sim F$ .

Suppose that  $F$  is of the form  $\bar{x}(\bar{y}(F'), F'')$ . Since  $\phi \sim F$ , in particular,  $\phi \sim \bar{x}(\bar{y} \cup \bar{z}', \bar{z}'')$ , where  $\bar{z}', \bar{z}''$  are defined as in Algorithm 1. Now, we inspect the proof of Lemma 4. Let  $\phi_{\bar{x}}$  be the result of applying the rewriting rule  $R_1$  of this proof. It is clear, by hypothesis, that we have  $\phi_{\bar{x}} \sim \{\bar{y} \cup \bar{z}', \bar{z}''\}$ . Hence, algorithm of the proof of Theorem 1 outputs a decomposition of the form  $\bigvee_{\bar{l} \in L} \beta_1^{\bar{l}}(\bar{y}, \bar{z}') \wedge cl_2^{l_2}(\bar{z}'')$ , exactly as defined in the proof of Theorem 1. We let  $\psi^{-1}$  be the result of applying exhaustively the rewrite rule  $R_2$  of the proof of Lemma 4 on  $\psi$ , for all formulas  $\psi$ . Hence,  $D(\phi, \bar{x}(\bar{y} \cup \bar{z}', \bar{z}''))$  returns the formula  $\bigvee_{\bar{l} \in L} (\beta_1^{\bar{l}})^{-1}(\bar{x}, \bar{y}, \bar{z}') \wedge (cl_2^{l_2})^{-1}(\bar{x}, \bar{z}'')$ . It remains to prove that formulas  $(\beta_1^{\bar{l}})^{-1}$  and  $(cl_2^{l_2})^{-1}$  satisfy  $(\beta_1^{\bar{l}})^{-1} \sim \bar{y}(\bar{x}, F')$  and  $(cl_2^{l_2})^{-1} \sim \bar{x}(F'')$ . We only prove it for formulas  $(cl_2^{l_2})^{-1}$ , as the proof for formulas  $(\beta_1^{\bar{l}})^{-1}$  is analogous. So let us fix some natural  $l_2$ . By going back to the definition of formula  $cl_2^{l_2}$ , we can prove that  $(cl_2^{l_2})^{-1}$  is equivalent to a formula of the form  $\Gamma = \exists \bar{z}_0, \gamma(\bar{x}, \bar{z}_0) \wedge \forall \bar{u}, \phi(\bar{x}, \bar{u}, \bar{z}') \leftrightarrow \phi(\bar{x}, \bar{u}, \bar{z}_0)$ , for some  $\gamma$ . Now, since  $\phi \sim F$ , it is easy to see that  $\phi$  is equivalent to a formula of the form  $\Psi = \bigvee_{i=1}^n \epsilon_i^1(\bar{x}, \bar{y}, \bar{z}') \wedge \epsilon_i^2(\bar{x}, \bar{z}'')$ , such that  $\epsilon_i^1 \sim \bar{y}(\bar{x}, F')$  and  $\epsilon_i^2 \sim \bar{x}(F'')$  for  $i = 1, \dots, n$ . We replace in  $\Gamma$  the formula  $\phi$  by  $\Psi$ , and, after a series of rewritings (by pushing up disjunctions and pushing down quantifiers), we can prove that  $(cl_2^{l_2})^{-1}$  is equivalent to a formula of the form  $\bigvee_{i=1}^n \gamma_i(\bar{x}) \epsilon_i^2(\bar{x}, \bar{z}'') \vee \bigvee_{P \subseteq \{1, \dots, n\}} \gamma'_P(\bar{x}) \wedge \bigwedge_{i \in P} \neg \epsilon_i^2(\bar{x}, \bar{z}'')$ , for formulas  $\gamma_i, \gamma'_P$  depending only on  $i$  and  $P$ . The conclusion follows by property (\*) and the fact that every  $\epsilon_i^2$  satisfies  $\epsilon_i^2 \sim \bar{x}(F'')$ .

Suppose now that  $F$  is of the form  $\{T_1, \dots, T_k\}$ , let  $\bar{x}_i$  be the variables occurring in  $T_i$  for  $i = 1, \dots, k$ , and let  $P = \{\bar{x}_1, \dots, \bar{x}_k\}$ . Since  $\phi \sim F$ , in particular,  $\phi \sim P$ . Hence  $D(\phi, P)$  outputs a decomposition of the form  $\bigvee_{\bar{l} \in L} \beta_1^{\bar{l}}(\bar{x}_1) \wedge cl_2^{l_2}(\bar{x}_2) \wedge \dots \wedge cl_k^{l_k}(\bar{x}_k)$ , exactly as defined in the proof of Theorem 1. We have to prove that for all  $\bar{l} = (l_1, \dots, l_k) \in L$ , we have  $\beta_1^{\bar{l}} \sim T_1$ , and  $cl_i^{l_i} \sim T_i$  for  $i = 2, \dots, k$ . This is sufficient, since by induction hypothesis,  $D$  will output a decomposition of every  $\beta_1^{\bar{l}}$  and  $cl_i^{l_i}$ . We only prove it for formulas  $cl_i^{l_i}(\bar{x}_i)$ , as it is similar for formulas  $\beta_1^{\bar{l}}$ . Let us fix some  $\bar{l}$  and  $i$ . We come back to the definition of  $cl_i^{l_i}$ , and we can easily show that it is of the form  $\exists \bar{y}, \gamma(\bar{y}) \wedge \psi_\phi^i(\bar{x}_i, \bar{y})$ , for some formula  $\gamma$  which selects the minimal representatives of the  $l_i$ -th equivalence class of the relation defined by  $\psi_\phi^i$ , where  $\psi_\phi^i$  has been defined in Section 4. Now, since  $\phi \sim F$ ,  $\phi$  is equivalent to a formula of the form  $\Psi = \bigvee_{p=1}^n \bigwedge_{j=1}^k \epsilon_p^j(\bar{x}_j)$  such that every  $\epsilon_p^j$  satisfy  $\epsilon_p^j \sim T_j$ . Next, in  $\exists \bar{y}, \gamma(\bar{y}) \wedge \psi_\phi^i(\bar{x}_i, \bar{y})$ , we replace  $\phi$  by  $\Psi$  ( $\psi_\phi^i$  can be viewed as the result of applying the function  $\psi^i$  on  $\phi$ , and we just replace  $\psi_\phi^i(\bar{x}_i, \bar{y})$  by  $\psi_\Psi^i(\bar{x}_i, \bar{y})$ ). We get a formula equivalent to  $cl_i^{l_i}$  which, after a series of rewritings preserving equivalence (by moving up disjunctions and pushing down quantifiers), rewrites to a formula (equivalent to  $cl_i^{l_i}$ ) of the form  $\bigvee_{p=1}^n \phi_p \wedge \epsilon_p^i(\bar{x}_i) \vee \bigvee_{Q \subseteq \{1, \dots, n\}} \psi_Q \wedge \bigwedge_{p \in Q} \neg \epsilon_p^i(\bar{x}_i)$ , for some closed formulas  $\phi_p, \psi_Q$  depending on  $p$  and  $Q$ . The conclusion follows by using property (\*) and the fact that every  $\epsilon_p^i$  satisfies  $\epsilon_p^i \sim T_i$ .  $\square$

Similarly as the case of variable independence, if  $\phi(x_1, \dots, x_n)$  conforms to  $F$ , then for any tree  $t \in T_\Sigma$ ,  $q_\phi(t)$  can be represented by an aggregated answer of size  $O(n|t|^b)$  ( $\phi$  is assumed to be fixed, and necessarily  $|F| \leq n$ ). The parameter  $b$  denotes the maximal sum of the size of a label of  $F$  plus the size of the label of its father if it exists.

Note that variable independence w.r.t. a dependence forest subsumes variable independence w.r.t. a partition, since a partition can be viewed as a dependence forest consisting of a set of leaves. Moreover, as stated by the next theorem, there is an MSO formula  $\phi$  such that there is no dependence forest  $F$  such that: (i) labels of  $F$  are singletons (ii)  $\phi \sim F$ . Nevertheless, we know that on trees, every MSO formula is equivalent to an existentially quantified Boolean combination of MSO formulas in two free variables [16, 10].

**Theorem 4.** *There is an MSO formula  $\phi$  such that there is **no** dependence forest  $F$  whose labels are singletons and such that  $\phi \sim F$ .*

*Proof.* Let  $x \preceq y$  be an MSO formula which holds in a tree if  $y$  is a descendant of  $x$ . It is well-known that it can easily be defined as the reflexive and transitive closure of  $S_1 \vee S_2$ , this closure being definable in MSO.

Now, let  $\phi(x, y, z)$  be an MSO formula defined by:

$$\begin{aligned} \phi(x, y, z) = \exists \alpha \quad & \alpha \preceq x \wedge \alpha \preceq y \wedge \alpha \preceq z \\ & \wedge \forall \alpha' \quad \alpha' \preceq x \wedge \alpha' \preceq y \implies \alpha' \preceq \alpha \\ & \wedge \forall \alpha' \quad \alpha' \preceq x \wedge \alpha' \preceq z \implies \alpha' \preceq \alpha \end{aligned}$$

For all trees  $t$  and all nodes  $u, v, w$  of  $t$ , we have  $t \models \phi(u, v, w)$  iff the least common ancestor of  $u$  and  $v$  is equal to the least common ancestor of  $u$  and  $w$ .

We now prove by applying algorithm 1 that there is no forest  $F$  whose labels are singletons such that  $\phi \sim F$ , by proving it for each forest over  $\{x, y, z\}$ . Let  $n \geq 0$  and  $t_n$  be the tree over the alphabet  $\{a\}$  inductively defined by  $t_0 = a$  and  $t_n = a(t_{n-1}, a)$ . For all  $n$ , we denote by  $v_0, v_1, \dots, v_n$  the nodes  $1^n, 1^{n-1}, \dots, \epsilon$  respectively, and, if  $n > 0$ , by  $w_1, \dots, w_n$  the nodes  $1^{n-1}.2, 1^{n-2}.2, \dots, 2$  respectively.

1.  $F = \{x, y, z\}$  or  $F = \{x, y(z)\}$  or  $F = \{x, z(y)\}$ . We apply algorithm 1. The formula  $\psi_\phi^1(x, x')$  is defined by  $\forall y, z \phi(x, y, z) \leftrightarrow \phi(x', y, z)$ . We prove that for all  $n$ , and all  $i, j \leq n$ ,  $v_i \neq v_j$  implies  $t_n \not\models \psi_\phi^1(v_i, v_j)$ . Indeed, if  $v_i \neq v_j$  such that  $\text{desc}(v_i, v_j)$ , then we have  $t \models \phi(v_i, v_i, v_j)$  but  $t \not\models \phi(v_j, v_i, v_j)$ . Hence, the number of classes of the equivalence relation defined by  $\psi_\phi^1$  is at least  $n$ , which is unbounded. So Algorithm 1 breaks;
2.  $F = \{y, x(z)\}$  or  $F = \{y, z(x)\}$ . The formula  $\psi_\phi^2(y, y')$  is defined by  $\forall x, y, z \phi(x, y, z) \leftrightarrow \phi(x, y', z)$ . We prove that for all  $n$ , and all  $i, j \leq n$ ,  $v_i \neq v_j$  implies  $t_n \not\models \psi_\phi^2(v_i, v_j)$ . Indeed, if  $v_i \neq v_j$  such that  $\text{desc}(v_i, v_j)$ , then we have  $t \models \phi(v_j, v_i, v_i)$  but  $t \not\models \phi(v_j, v_j, v_i)$ . Hence, the number of classes of the equivalence relation defined by  $\psi_\phi^2$  is at least  $n$ , which is unbounded. So Algorithm 1 breaks;
3.  $F = \{z, y(x)\}$  or  $F = \{z, x(y)\}$ . Those cases are symmetric to the previous ones.
4.  $F = x(y, z)$ . We let  $\psi(y, y') = \forall z, \phi_x(y, z) \leftrightarrow \phi_x(y', z)$  where  $\phi_x$  has been defined in the proof of Lemma 4. By definition of Algorithm 1, if the equivalence relation defined by  $\psi$  has an unbounded index, then the algorithm fails. This is what we next prove. Let  $n \leq 0$ . We fix  $x$  by a Boolean in the tree  $t_n$ : we let  $t'_n$  be the tree over  $\{a\} \times \{0, 1\}$  such that  $N_{t_n} = N_{t'_n}$  and all nodes are labeled  $(a, 0)$  except  $v_0$  which is labeled  $(a, 1)$ . It is easy to see that for all  $i \geq 1$ , we have  $t'_n \models \phi_x(v_i, w_i)$  and for all  $j > i$ , we have  $t'_n \not\models \phi_x(v_j, w_i)$ . Hence there are at least  $n$  equivalence classes for the relation defined by  $\psi$ . So Algorithm 1 breaks.
5.  $F = y(x, z)$ . Similarly to the previous case, we fix a variable. Let  $\psi(z, z') = \forall x, \phi_y(x, z) \leftrightarrow \phi_y(x, z')$ . Let  $t'_n$  be the tree defined in the previous case. Hence  $y$  is fixed to denote the node  $v_0$ . We can prove that for all  $i > 0$ , we have  $t'_n \models \phi_y(v_i, w_i)$  but for all  $j > i$ , we have  $t'_n \not\models \phi_y(v_i, w_j)$ . Hence there are at least  $n$  equivalence classes for the relation defined by  $\psi$ . So Algorithm 1 breaks.
6.  $F = z(x, y)$ . This case is symmetric to the previous one. □

### Further Extensions

First note that all the results presented in the paper also hold for FO-queries, as we do not use second order variables in decompositions (but in this case we need to add in the tree structure a total order on the nodes).

We would like to investigate independence problems for more general classes of structures  $C$ . Indeed, we can give two sufficient conditions for relative independence to be decidable on  $C$ : (i) boundedness of an MSO formula with first-order variables is decidable on  $C$ , (ii) there is a computable MSO-definable total order on the elements of the structures of  $C$ . The first point has already been detailed in Section 3, while the second point is studied in [9]. This is the case for instance for unranked tree structures, over the signature consisting of the first-child and next-sibling predicates, and predicates to test the labels.

Finally, we would like to extend independence w.r.t. a dependence forest to independence w.r.t. a dependence graph. The techniques presented here do not seem to be easily extendable to graphs, even for a clique of size 3 for instance. In particular, we cannot use an inductive proof based on Lemma 4 anymore.

## References

1. A. Berlea. On-the-fly tuple selection for XQuery. In *International Workshop on XQuery Implementation, Experience and Perspectives*, June 2007.
2. M. Bojanczyk. A bounding quantifier. In *14th Annual Conference of the EACSL on Computer Science Logic*, 2004.
3. J. Chomicki, D. Goldin, G. Kuper, and D. Toman. Variable independence in constraint databases. *IEEE Transactions on Knowledge and Data Engineering*, 15(6):1422–1436, 2003.
4. H. Comon, M. Dauchet, R. Gilleron, C. Löding, F. Jacquemard, D. Lugiez, S. Tison, and M. Tommasi. Tree automata techniques and applications. Available on: <http://www.grappa.univ-lille3.fr/tata>, 2007.
5. S. Cosmadakis, G. Kuper, and L. Libkin. On the orthographic dimension of definable sets. *Inf. Process. Lett.*, 79(3):141–145, 2001.
6. B. Courcelle. Structural properties of context-free sets of graphs generated by vertex replacement. *Inf. Comput.*, 116(2):275–293, 1995.
7. B. Courcelle. Linear delay enumeration and monadic second-order logic. 2007. To appear in *Discrete Applied Mathematics*.
8. B. Courcelle and S. Oum. Vertex-minors, monadic second-order logic, and a conjecture by seese. *J. Comb. Theory Ser. B*, 97(1):91–126, 2007.
9. Bruno Courcelle. The monadic second-order logic of graphs x: linear orderings. *Theor. Comput. Sci.*, 160(1-2):87–143, 1996.
10. E. Filiot, J. Niehren, J.-M. Talbot, and S. Tison. Polynomial time fragments of xpath with variables. In *ACM Symposium on Principles of Database Systems*, 2007.
11. S. Grumbach, P. Rigaux, and L. Segoufin. On the orthographic dimension of constraint databases. In *7th International Conference on Database Theory*, pages 199–216. Springer-Verlag, 1999.
12. L. Libkin. Variable independence for first-order definable constraints. *TOCL*, 4(4):431–451, 2003.
13. L. Libkin. Logics over unranked trees: an overview. *Logical Methods in Computer Science*, 3(2):1–31, 2006.
14. Holger Meuss, Klaus U. Schulz, and François Bry. Towards aggregated answers for semistructured data. In *ICDT '01: Proceedings of the 8th International Conference on Database Theory*, pages 346–360, 2001.
15. J. Niehren, L. Planque, J.-M. Talbot, and S. Tison. N-ary queries by tree automata. In *10th International Symposium on Database Programming Languages*, volume 3774, pages 217–231, 2005.
16. T. Schwentick. On diving in trees. In *25th International Symposium on Mathematical Foundations of Computer Science*, pages 660–669, 2000.
17. H. Seidl. Ambiguity, valuedness and costs. 1992. Habilitation Thesis.
18. H. Seidl. Equivalence of finite-valued tree transducers is decidable. *Math. Syst. Theory*, 27(4):285–346, 1994.
19. J. W. Thatcher and J. B. Wright. Generalized finite automata with an application to a decision problem of second-order logic. 2:57–82, 1968.