

Model Based Deduction for Database Schema Reasoning

Peter Baumgartner¹, Ulrich Furbach², Margret Gross-Hardt², and Thomas Kleemann²

¹ MPI Informatik, D-66123 Saarbrücken, Germany, baumgart@mpi-sb.mpg.de

² Universität Koblenz-Landau, D-56070 Koblenz, Germany,
{uli|margret|tomkl}@uni-koblenz.de

Abstract. We aim to demonstrate that automated deduction techniques, in particular those following the model computation paradigm, are very well suited for database schema/query reasoning. Specifically, we present an approach to compute completed paths for database or XPath queries. The database schema and a query are transformed to disjunctive logic programs with default negation, using a description logic as an intermediate language. Our underlying deduction system, *KRHyper*, then detects if a query is satisfiable or not. In case of a satisfiable query, all completed paths – those that fulfill all given constraints – are returned as part of the computed models.

The purpose of computing completed paths is to reduce the workload on a query processor. Without the path completion, a usual XPath query processor would search the whole database for solutions to the query, which need not be the case when using completed paths instead.

We understand this paper as a first step, that covers a basic schema/query reasoning task by model-based deduction. Due to the underlying expressive logic formalism we expect our approach to easily adapt to more sophisticated problem settings, like type hierarchies as they evolve within the XML world.

1 Introduction

Automated theorem proving is offering numerous tools and methods to be used in other areas of computer science. An extensive overview about the state of the art and its potential for applications is given in [7]. Very often there are special purpose reasoning procedures which are used to reason for different purposes like knowledge representation [1] or logic programming [10].

The most popular methods used for practical applications are resolution-based procedures or model checking algorithms. In this paper we want to demonstrate that there is a high potential for model based procedures for database schema reasoning. Model based deduction can be based very naturally on tableau calculi [12], and in particular on the developments that started with the *SATCHMO* approach [16], which was refined later and extended in the hyper tableau calculus [6].

We start with the idea of representing a database schema as a description logic knowledge base. This idea as such is not new and has been put forward in [8, 9]. However, we found that the services usually available in description logic reasoners do not allow to express all constraints imposed by the schema in order to solve the tasks we are looking at. Indeed, the work in [8, 9] aims at different purposes, where schema reasoning tasks can be reduced to *satisfiability* of description logic knowledge bases.

We are considering the tasks of testing and optimizing certain forms of database queries as they arise in the XML world. To this end, a “pure” description logic approach was proposed before in [4]. In the present paper, the limitations of that approach are overcome by translating a schema and a given XPath like query into a disjunctive logic program (with default negation). The *KRHyper* system then detects if a query is satisfiable or not. In case of a satisfiable query, all completed paths – those that fulfill all given constraints – are returned as part of the computed models. The purpose of computing completed paths is to reduce the workload on a query processor. Without the path completion, a usual XPath query processor would search the whole database for solutions to the query, which need not be the case when using completed paths instead. The usage of a *model generation* theorem prover thus is motivated by the applications requirement to enumerate models/answers rather than querying the existence of a model.

We start with a brief review of the hyper tableau prover.

2 Theorem Proving with Hyper Tableau

Features. The Hyper Tableau Calculus is a clause normal form tableau calculus [6], which can be seen as a generalization of the SATCHMO-procedure [16]. Hyper tableau have been used in various applications (for examples see [3, 5]), where two aspects turned out to be of importance: The result of the theorem prover is a model (if the specification is satisfiable) and this model can be seen as the result of the prover’s “computation”; it can be used by the system, where the prover is embedded, for further computation steps. The second aspect is concerned with default negation. Although an entire discipline, namely knowledge representation, is emphasizing the necessity of non-monotonic constructs for knowledge representation, there are only very few sophisticated systems dealing with such constructs [17, 11].

The hyper tableau theorem prover *KRHyper* allows application tasks to be specified by using first order logic — plus possibly non-monotonic constructs — in clausal form. While *KRHyper* can be used straightforwardly to prove theorems, it also allows the following features, which are on one hand essential for knowledge based applications, but on the other hand usually not provided by first order theorem provers:

1. Queries which have the listing of predicate extensions as answer are supported.
2. Queries may also have the different extensions of predicates in alternative models as answer.
3. Large sets of uniformly structured input facts are handled efficiently.
4. Arithmetic evaluation is supported.
5. Default negation is supported.
6. The reasoning system can output proofs of derived facts.

More details about these features can be found in [21]. Also, we only note that with a simple transformation of the given rule set, *KRHyper* is sound, complete and terminating with respect to the *possible models* [18] of a stratified disjunctive logic program without function symbols (except constants).

A Small Example. Hyper tableau is a “bottom-up” method, which means that it generates instances of rule³ heads from facts that have been input or previously derived. If a hyper tableau derivation terminates without having found a proof, the derived facts form a representation of a model of the input clauses.

The following example illustrates how our hyper tableau calculus based system, *KRHyper*, proceeds to generate models. Figure 1 shows four subsequent stages of a derivation for the following input clauses⁴:

$$\begin{aligned}
 p(a) &\leftarrow & (1) \\
 q(x,y) \vee r(f(z)) \vee r(x) &\leftarrow p(x) & (2) \\
 &\leftarrow q(x,x) & (3) \\
 s(x) &\leftarrow p(x), \text{ not } r(x) & (4)
 \end{aligned}$$

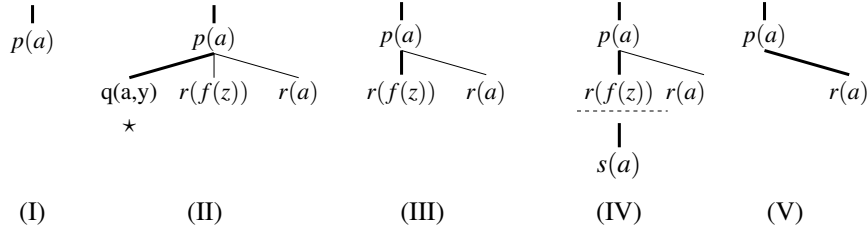


Fig. 1. Stages of a *KRHyper* derivation

KRHyper provides stratified default negation. The set of input clauses is partitioned into *strata*, according to the predicates in their heads: if a clause c_1 has a head predicate appearing in the scope of the negation operator “not” in the body of c_2 , then c_1 is in a lower stratum than c_2 . In the example, we have two strata: the lower one containing clauses (1), (2) and (3), the higher one clause (4).

As noted, a rule head may be a disjunction. In hyper tableau, disjunctions are handled by exploring the alternative branches in a systematic way. This explains the tree structure in Figure 1. Backtracking can be used to generate one model after the other.

Stage (I) shows the data structure maintained by the method, also called *hyper tableau*, after the input fact (1) has been processed. One can view the calculus as attempting to construct the representation of a model, the *active branch*, shown with bold lines in the figure. At step (I), this model fragment contradicts for example with clause (2): a model containing $p(a)$ must also contain all instances of $q(a,y)$ or of $r(f(z))$ or

³ We use an implication-style notation for clauses throughout this paper: A clause is viewed as rule $Head \leftarrow Body$, where $Head$ consists of its positive literals, combined by “ \vee ”, and $Body$ consists of its negative literals, combined by “ $;$ ” (*and*). Both the head and the body may be empty.

⁴ Here and below, the letters x, y, z denote variables, while a, b, c denote constants.

$r(a)$. The model fragment is “repaired” by derivating consequences and attaching them to the hyper tableau: The corresponding instance of clause (2) is attached to the hyper tableau. Since it has a disjunctive head, the tableau splits into three branches. The first branch is inspected and proved contradictory with clause (3) (the branch is said to be *closed*). This state is shown in (II).

Computation now tracks back and works on the second branch. With the clauses of the lower stratum, no further facts can be derived at this branch, which means that a model for the stratum has been found, as shown in step (III). Computation then proceeds with the next higher stratum: $s(a)$ can be derived by clause (4). Since no further facts can be derived, a model for the whole clause set has been found, represented by the facts on the active branch: $\{p(a), r(f(z)), s(a)\}$, as shown in (IV).

If desired, the procedure can backtrack again and continue to find another model, as shown in state (V). Another backtracking step then finally leads to the result, that there is no further model.

We conclude by noting that the *KRHyper* system implements the calculus by a combination of semi-naive rule evaluation with backtracking over alternative disjuncts and iterative deepening over a term weight bound. It extends the language of first order logic by stratified default negation and built-ins for arithmetic.

3 Flexible Database Queries for XML Data

Querying databases requires that users know the structure of a database. In fact, they have to know the *database schema* in order to formulate valid queries. In the context of complex structured data and large database schemas, knowing the complete schema is not always possible. Querying data, therefore, may be a tedious task. This section describes the application of the *KRHyper* System in order to enhance the flexibility of database queries. In particular, we focus on XML databases and address the following issues in querying XML databases:

- XML documents contain complex structured data, often rather nested. Users therefore have to navigate through these data.
- XML data usually is described by Document Type Definitions (DTDs) and more recently, XML Schema is used. Different from DTDs, XML Schema offers in some sense object oriented concepts as user defined types and aggregation as well as specialization relationships between them. During query processing and optimization this schema knowledge may be used, e.g. in order to avoid evaluation of unsatisfiable queries.
- In an XML schema so called *substitution groups* define types that can be substituted for each other, comparable to union types in other (programming) languages, though, with the difference that types which are substitutes for each other have to be related via specialization.
- Existing querying languages like XQuery [20] offer navigational expressions on the document level in order to access parts of an XML document. These languages do not cope with type expressions. Type expressions may be helpful in order to query instances of some *general* type T , resulting in instances of type T as well as of all subtypes for T .

Let us consider an example XML document representing a university with a library and researchers working in the university. A library consists of books where each book has a title, an author and an ISBN. Researchers have a name and an associated set of publications e.g. articles, monographs or some general kind of publication.

An XML schema itself also is an XML document listing the complex types together with their elements referring to other (complex) types. Furthermore by means of a so called restriction expression, it is possible to represent specialization relationships between types. Instead of the linear, XML based description of an XML schema, we use a more illustrative, graphical representation for the types and their relationship in a schema. An XML Schema is represented by a schema graph, where nodes represent the types and substitution groups of the schema and edges represent aggregation and specialization relationships between types. Starting from such a schema graph, we present an approach that allows a user to query the data, even if only parts of a database schema are known. Figure 2 shows an example schema graph.

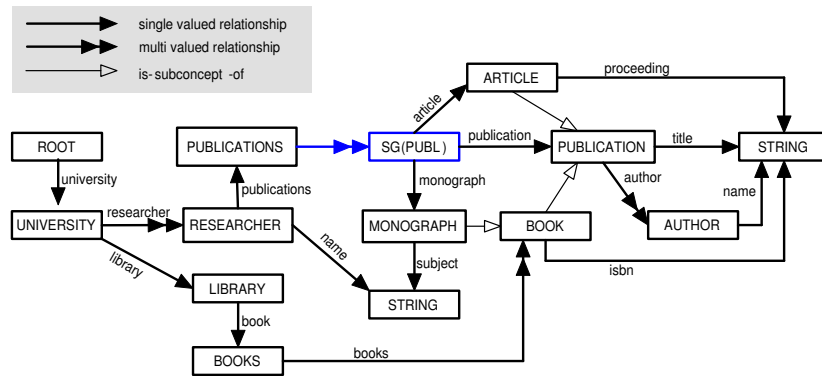


Fig. 2. Example schema graph.

This schema shows types like UNIVERSITY, PUBLICATION, AUTHOR etc. with their elements referring to other types as well as specialization relationships between types. For instance, the type UNIVERSITY has an element researcher of type RESEARCHER and an element library of type LIBRARY. Furthermore, PUBLICATION is a general type with specializations BOOK and ARTICLE. There is one substitution group SG(PUBL) contained in this schema whose general type is PUBLICATION and potential substitutes are ARTICLE and MONOGRAPH, actually specializations of PUBLICATION. In the transformation given below, we will see that substitution groups are a means to express a disjunction of disjoint concepts.

To keep the representation of schema graphs as well as the query processing simple, we only consider (XML-)elements describing complex data structures but do not cope with the term of (XML-)attributes. Nevertheless, we will use both terms in order to refer to properties of data items.

3.1 From XML Schema to Description Logics

An XML database is described by an XML Schema. The schema is represented by a graph. The nodes of the graph are the complex type identifiers; relationships between elements and their corresponding subelement types are represented by aggregation edges and "extension"-relationships, describing the generalization relationship between complex types, are represented by so called is-a edges. XML Schema supports the modelling of multiple complex types, that are extensions of the same general complex type within a substitution group, comparable to a union type in other languages.

In the following we assume a possibly infinite set \mathcal{L} of labels.

Definition 1 (XML Schema). *Let C be a set of type names and SG be a disjoint set of substitution group identifiers. An XML schema (or schema for short) is a graph $S = (C \cup SG, E_{\text{rel}} \cup E_{\text{isa}} \cup E_{\text{SG-in}} \cup E_{\text{SG-out}}, r)$ where C and SG are the nodes, E_{rel} , E_{isa} , $E_{\text{SG-in}}$, $E_{\text{SG-out}}$ are disjoint sets of edges (representing attributes/elements edges, inheritance edges, incoming edges to SG nodes, outgoing edges from SG nodes, respectively), and $r \in C$ represents the root of the schema. Every schema must satisfy the following properties:*

1. $(v, v) \notin E_{\text{rel}} \cup E_{\text{isa}} \cup E_{\text{SG-in}} \cup E_{\text{SG-out}}$, for any node v .
2. Each edge in E_{rel} and each edge in $E_{\text{SG-out}}$ is labeled with an element from \mathcal{L} . All other edges are not labeled.
3. $(v, v') \in E_{\text{SG-in}}$ if and only if $v' \in SG$ (incoming edges to nodes in SG are precisely those in $E_{\text{SG-in}}$).
4. $(v, v') \in E_{\text{SG-out}}$ if and only if $v \in SG$ (outgoing edges from nodes in SG are precisely those in $E_{\text{SG-out}}$).

The schema can be translated in a straightforward way to description logics as follows:⁵

Definition 2 (Schema to Description Logic). *Let $S = (C \cup SG, E_{\text{rel}} \cup E_{\text{isa}} \cup E_{\text{SG-in}} \cup E_{\text{SG-out}}, r)$ be a schema. The TBox for S , $T(S)$, is defined as the smallest set of inclusion statements satisfying the following properties:*

Translation of is-a links. *If $c \in C$ and $\{d \mid (c, d) \in E_{\text{isa}}\} = \{d_1, \dots, d_n\}$, for some $n \geq 1$, then $T(S)$ contains the inclusion $c \sqsubseteq d_1 \sqcup \dots \sqcup d_n$.*

Translation of elements/attributes. *If $(c, d) \in E_{\text{rel}}$ and (c, d) is labeled with l , then $T(S)$ contains the inclusion $c \sqsubseteq \exists l . d$.*

Translation of substitution groups. *If $(c, v) \in E_{\text{SG-in}}$ then $T(S)$ contains the inclusion $c \sqsubseteq v$. If $v \in SG$ and $\{d \mid (v, d) \in E_{\text{SG-out}}\} = \{d_1, \dots, d_n\}$, for some $n \geq 1$, then $T(S)$ contains the inclusion $v \sqsubseteq d_1 \sqcup \dots \sqcup d_n$ and the inclusions $d_i \sqcap d_j \sqsubseteq \perp$, for all i, j with $1 \leq i, j \leq n$ and $i \neq j$.*

This translation conforms to concept and role formations found even in basic description logic languages like \mathcal{ALC} . Although the translation does *not* result in an \mathcal{ALC} TBox ($T(S)$ might contain several inclusion statements with the same concept at the left hand side), it is easy to see that the result of the transformation can easily be brought to an \mathcal{ALC} conforming TBox (possibly cyclic).

⁵ We use standard description logic notation, see [1].

At this point we will not discuss how to employ description logic reasoners to solve the tasks we are interested in. This discussion will be postponed after our approach based on model computation has been described.

3.2 From Description Logics to Model Computation

The following translation is the standard relational translation from description logics to predicate logic. For our purpose, it is enough to work in a restricted setting, where all inclusions in a TBox are of a particular form, which is obtained as the result of the transformation in Definition 2.

Definition 3 (Description Logic to Rules – Basic Version). *Let S be a schema. The rules for S , $R(S)$, are defined as the smallest set of rules satisfying the following properties:*

1. *if $c \sqsubseteq c_1 \sqcup \dots \sqcup c_n \in T(S)$ then $R(S)$ contains the rule $c_1(x) \vee \dots \vee c_n(x) \leftarrow c(x)$*
2. *if $c \sqsubseteq \exists l.d \in T(S)$ then $R(S)$ contains the rules $l(x, f_{c,l,d}(x)) \leftarrow c(x)$ and $d(f_{c,l,d}(x)) \leftarrow c(x)$. ($f_{c,l,d}$ is a unary function symbol whose name contains c , l and d , as indicated.)*
3. *if $c \sqcap d \sqsubseteq \perp \in T(S)$ then $R(S)$ contains the rule $\text{false} \leftarrow c(x), d(x)$.*

Using this transformation, simple graph reachability problems can be reduced easily to model computation problems. Speaking in terms of the schema graph, to compute a path, say, from a node c to a node d in a schema S , it suffices to add to $R(S)$ the fact $c(a) \leftarrow$ (for some constant a) and the rules $\text{found} \leftarrow d(x)$ and $\text{false} \leftarrow \text{not found}$, where found is a predicate symbol not occurring in $R(S)$. Each model of the thus obtained program corresponds to exactly one path from c to d in S . However, this approach works only in a satisfactory way if the schema does not contain any circle.

Example 1 (Cycle). Consider a TBox consisting of the two inclusions $c \sqsubseteq \exists l.d$ and $d \sqsubseteq \exists k.c$. It can be obtained by translating a suitable schema containing a circle. Its translation to rules gives the following program:

$$\begin{array}{ll} l(x, f_{c,l,d}(x)) \leftarrow c(x) & d(f_{c,l,d}(x)) \leftarrow c(x) \\ k(x, f_{d,k,c}(x)) \leftarrow d(x) & c(f_{d,k,c}(x)) \leftarrow d(x) \end{array}$$

Now, any Herbrand model as computed by bottom-up procedures will be infinite and contains $c(a)$, $c(f_{d,k,c}(f_{c,l,d}(a)))$, $c(f_{d,k,c}(f_{c,l,d}(f_{d,k,c}(f_{c,l,d}(a)))))$ and so on. Therefore, *KRHyper* and related procedures will not terminate.

3.3 Blocking by Transformation

In this section we give an improved transformation in order to guarantee termination of the model computation. This will be achieved by a “loop check” similar to the blocking technique known from the description logic literature [14, e.g.]. The idea is to re-use an individual already known to belong to a certain concept instead of adding a new individual to it in order to satisfy an existentially quantified role constraint. In the example, the individual a can be re-used instead of $f_{d,k,c}(f_{c,l,d}(a))$ in order to put $f_{c,l,d}(a)$ into the k -relation to some individual belonging to c . This re-use technique will be described now. It will guarantee the termination of our reasoning algorithm.

Definition 4 (Description Logic to Rules – Improved Version). Let S be a schema. The rules for S , $R(S)$, are defined as the smallest set of rules satisfying the following properties:

1. if $c \sqsubseteq c_1 \sqcup \dots \sqcup c_n \in T(S)$ then $R(S)$ contains the rule $c_1(x) \vee \dots \vee c_n(x) \leftarrow c(x)$
2. $R(S)$ contains the fact $\text{equal}(x, x) \leftarrow$.
3. if $c \sqsubseteq \exists l.d \in T(S)$ then $R(S)$ contains the following rules:

$$\text{new}_{c,l,d}(x) \vee \text{old}_{c,l,d}(x) \leftarrow c(x) \quad (1)$$

$$\text{false} \leftarrow \text{new}_{c,l,d}(x), \text{old}_{c,l,d}(x) \quad (2)$$

$$l(x, f_{c,l,d}(x)) \leftarrow \text{new}_{c,l,d}(x) \quad (3)$$

$$d(f_{c,l,d}(x)) \leftarrow \text{new}_{c,l,d}(x) \quad (4)$$

$$l(x, z) \leftarrow \text{old}_{c,l,d}(x), c(y), l(y, z), d(z) \quad (5)$$

$$\text{false} \leftarrow \text{old}_{c,l,d}(x), \text{not some}_{c,l,d} \quad (6)$$

$$\text{some}_{c,l,d} \leftarrow c(x), l(x, y), d(y) \quad (7)$$

$$\text{false} \leftarrow \text{new}_{c,l,d}(x), \text{new}_{c,l,d}(y), \text{not equal}(x, y) \quad (8)$$

4. if $c \sqcap d \sqsubseteq \perp \in T(S)$ then $R(S)$ contains the rule $\text{false} \leftarrow c(x), d(x)$.

Some comments are due. The difference to the previous version is the translation of inclusions of the form $c \sqsubseteq \exists l.d$. In order to explain it, suppose that the concept c is populated with some individual, say, a . That is, when constructing a model, $c(a)$ already holds true. Now, the program above distinguishes two complementary cases to satisfy the constraint $\exists l.d$ for a : either a new l -connection is made between a and some new individual in d , or an existing (“old”) l -connection between some individual from c and from d is re-used. That exactly one of these cases applies is guaranteed by the rules (1) and (2). The rules (3) and (4) are responsible to establish a new connection, while the rule (5) is responsible to re-use an existing connection. To achieve the desired effect, some more constraints are needed: as said, re-using an existing connection is realized by applying the rule (5). It establishes the connection $l(x, z)$, where x stands for the object the connection is to be established from (a in the example), and z stands for the re-used object from d . However, there is no guarantee per se that the rule’s subgoals $c(y)$, $l(y, z)$ and $d(z)$ are satisfied. This, however, is achieved by the rules (6) and (7): whenever the program chooses to re-use an old connection, i.e. to build a model containing this choice, by (6) and (7) this can succeed only if the mentioned subgoals are satisfiable. Finally, the rule (8) acts as a “loop check”: with it, it is impossible that between individuals belonging to the concepts c and d more than one new l -connection is made. Only new l -connections cause insertion of more complex atoms⁶ and thus are the only source for non-termination. With the rule (8) there is a finite bound on the complexity then for a given program.

The program above is intended to be run by a bottom-up model computation procedure like *KRHyper* (Section 2). Together with some more rules and facts obtained by further transformation steps this yields an algorithm that is similar to usual tableau

⁶ Complexity being measured as the tree depth of the atoms.

algorithms developed for description logics. On the one side, our translation and the reasoning tasks to be solved do not quite match those in existing algorithms. This is because of the use of default negation to filter out nonintended models (see Section 3.5). Another difficulty we encountered with existing systems is their inability to actually output the computed models. From our application point of view this is problematic, as the answer to the tasks to be solved *is* the model (see again Section 3.5).

On the other side, it suffices for our purpose to work with TBoxes that are quite simple and do not involve constructs that are notoriously difficult to handle (like the combination of inverse roles, transitive roles and number restrictions). This allows us to use the above rather simple “loop checking” technique, which is inspired by the blocking technique developed for an ABox/TBox reasoner in [14].

3.4 Query Language

Existing query languages use path queries that navigate along the structures of the XML data. For instance, in order to access the name of all researchers of a university in *XPath* [19] one may use the XPath expression `/university/researchers/researcher/name`. Path queries usually allow some form of “abbreviation”. For instance, with `//researcher/name` one addresses all descendants of the “root” that are *researcher*-elements and navigate to their names. However, because path queries work directly on the XML data and not on the schema, it is not possible to query those elements from a data source that belongs to the same type or concept. In particular, in order to ask e.g. for all kinds of publications, one would have to construct the union of path queries navigating to publication, book, article and monograph, explicitly.

This problem has been addressed in [13] where concepts or type expressions, respectively, have been added to the query language. Querying instances of types or concepts is well known in object oriented databases. Furthermore, path expressions allow to navigate through the nested structure of the data. We assume a syntax similar to that applied in object oriented databases [15]. We aim at a query language that combines schema expressions as used in object oriented query languages with a flexible navigation mechanism as given by “abbreviated” path expressions as e.g. provided by XPath.

Let A denote a set of attribute names.

Definition 5 (Query Syntax). A path term is an expression of the form $c[x] op_1 a_1[x_1] \dots op_m a_m[x_m]$, where $m \geq 0$, c is a type name, a substitution group identifier (cf. Def. 1) or the symbol \top , $op_i \in \{., !\}$, $a_i \in A$, and x, x_i , for $i = 1, \dots, m$ are variables.

A conjunctive query expression is a conjunction of path terms $p_1 \wedge \dots \wedge p_n$, where $n \geq 1$. A disjunctive query expression is a disjunction of conjunctive query expressions $e_1 \vee \dots \vee e_n$, where $n \geq 1$.

By simply a *query expression* we mean a disjunctive query expression, which includes the case of a conjunctive query expression as a disjunction with one element.

A path term is an expression that starts in a concept and navigates through a schema by a sequence of attributes. Variables are used to “hold” the spots during this navigation. At the instance level, a path term describes a set of paths in a data source. The result of a path term basically is a relationship, where all variables occurring in a path term are bound to elements in the XML document. Actually, there are two possibilities to

traverse a schema. First, by explicit navigation that specifies a path step by step. We use the “!” operator for this kind of navigation. Second, a path term may specify only some attributes occurring on one or several paths in a schema; then the “.” operator is used.

For instance, $\text{BOOK}[b]!\text{title}[t]$ describes all instances b of type BOOK, with a title t as subelement. The result, basically, is a binary relation containing values for b and t . Compared to an XPath-expression, the “!” operator matches “/” and the “.” operator can be compared to “//”. Different from an XPath expression, in our query syntax variables can be specified and type expressions are possible. If a user does not want to specify an explicit type, the most general type \top can be taken. As will be shown below (see Section 3.5), using type expressions in a query allows a user to query for different elements described by the same general type in one simple expression. For instance, $\text{BOOK}[b]!\text{title}[t]$ retrieves all BOOK elements with their titles as well as all MONO-GRAPH elements, because of the underlying specialization/subsumption relationship.

Furthermore, using type expressions in a query provide a query processor with a possibility to validate purely by means of the schema if a query is satisfiable. Consider another query: $\text{PUBLICATION}[p].\text{isbn}[i] \wedge \top[p].\text{proceedings}[x]$. This query asks for all instances p of PUBLICATION with their isbn and their proceedings. Actually, there are no such instances. This fact can be established automatically, as will be shown below. This means it is useless to try this query on any concrete database satisfying the schema, as the result will be empty anyway.

We conclude this section by informally describing the semantics of queries: we say a query expression $q = e_1 \vee \dots \vee e_n$ (cf. Definition 5) is *satisfiable in a schema S* (cf. Definition 1) iff there is a conjunctive query expression e_i , ($1 \leq i \leq n$), and a substitution σ mapping each variable in e_i to a type name or a substitution group identifier, such that (i) S satisfies each path term in $e_i\sigma$ according to the just indicated semantics of the operators “!” and “.”, and (ii) no type name or substitution group identifier is visited more than twice on the traversal of S as specified by $e_i\sigma$. Any such substitution is called a *solution for q* .

The rationale behind condition (ii) is to allow queries admitting solutions representing circles in the schema graph, but no circle should be followed more than once.

3.5 Translating Queries to Rules

The translation of a path term p of the form $c[x] \text{op}_1 a_1[x_1] \dots \text{op}_m a_m[x_m]$ (cf. Definition 5) is the following list of atoms $tr(p)$:

$$c(x), tr_{sel}(x, \text{op}_1 a_1[x_1]), tr_{sel}(x_1, \text{op}_2 a_2[x_2]), \dots, tr_{sel}(x_{m-1}, \text{op}_m a_m[x_m]) ,$$

where

$$tr_{sel}(x, \text{op } a[y]) = \begin{cases} a(x, y) & \text{if } \text{op} = ! \\ \text{role_filler_ref_trans}(x, z), a(z, y) & \text{if } \text{op} = . , \\ & \text{where } z \text{ is a fresh variable} \end{cases}$$

The purpose of $tr(p)$ is to translate the path term p into a sequence of subgoals that correspond to traversing the schema as prescribed by p and thereby assigning values to

the variables mentioned in p . Notice the translation distinguishes between the operators “.” and “!”. The former stands for the presence of an attribute immediately at the current point of the traversal and thus translates into a corresponding role filler subgoal. The latter is similar, but it allows to follow an arbitrary number of attributes first, by means of the `role_filler_ref_trans` relation.

(1) The translation of a conjunctive query expression $p_1 \wedge \dots \wedge p_n$, where $n \geq 1$ and each p_i (for $i = 1, \dots, n$) is of the form just mentioned and written as $p_i = c^i[x^i]op_1^i a_1^i[x_1^i] \dots op_{m_i}^i a_{m_i}^i[x_{m_i}^i]$ consists of the following rules:

$$\text{solution_path}(x^1, (x_{m_1}^1, \dots, x_{m_n}^n)) \leftarrow \text{role_filler_ref_trans}(\text{init}, x^1), c^1(x^1), \\ \text{tr}(p_1), \dots, \text{tr}(p_n)$$

Observe that the selector x^1 mentioned in the first path term p_1 , is treated in a special way. Its type c^1 is treated as the “start type”, and a path from the root to it is computed in the first argument of the `solution_path` predicate.

The translation of a disjunctive conjunctive query expression $e_1 \vee \dots \vee e_n$ is the union of the translation of each e_i , for $i = 1, \dots, n$.

(2) The following rules constrain the admissible models to those that contain at least one “solution path”:

$$\text{false} \leftarrow \text{not some_solution_path} \\ \text{some_solution_path} \leftarrow \text{solution_path}(x, y)$$

(3) Let E_{rel} be a set of attribute names as mentioned in Definition 1. The set E_{rel} is reified by the set of rules

$$T(E_{\text{rel}}) = \{\text{role_filler}(x, y) \leftarrow l(x, y) \mid l \text{ is the label of some attribute in } E_{\text{rel}}\}$$

The reflexive-transitive closure of $T(E_{\text{rel}})$ is obtained as follows:

$$\text{role_filler_ref_trans}(x, x) \leftarrow \\ \text{role_filler_ref_trans}(x, z) \leftarrow \text{role_filler}(x, y), \text{role_filler_ref_trans}(y, z)$$

In the following definition all transformations introduced so far are collected and combined. It states the complete transformation applied to a schema and a query.

Definition 6 (Transformation of a Schema and a Query). *Let $S = (C \cup SG, E_{\text{rel}} \cup E_{\text{isa}} \cup E_{\text{SG-in}} \cup E_{\text{SG-out}}, r)$ be a schema as in Definition 1 and q a query expression. The transformation of S and q consists of the union of $R(S)$ of Definition 4, the result of the transformation step (1) applied to q , the rules from (2), the rules from (3) (both $T(E_{\text{rel}})$ and the rules for the reflexive-transitive closure of $T(E_{\text{rel}})$), and the facts*

$$\top(x) \leftarrow \text{root}(\text{init}) \leftarrow$$

As said earlier, the purpose of our model based approach is to detect if a XPath like query is satisfiable or not. In case of a satisfiable query, a fully completed path – one that fulfills all given constraints – is returned as part of the model. Furthermore, every such fully completed path will be computed. The following theorem states this result more formally and summarizes important properties of our transformation.

Theorem 1 (Soundness and Solution Completeness). *Let S be a schema as in Definition 1 and q a query expression. Let T be a tableau derived by $KRHyper$ applied to the transformation of S and q . Then the following holds:*

1. *Every open branch in T contains an (at least one) atom of the form $\text{solution_path}(t, (t_1, \dots, t_n))$ and it denotes a solution for q , by means of the substitution $\{x/t, x_1/t_1, \dots, x_n/t_n\}$, where x_1, \dots, x_n are the variables occurring in some conjunctive query expression of q , in this order.*
2. *For every solution $\{x/t, x_1/t_1, \dots, x_n/t_n\}$ for q there is an open branch in T containing the atom $\text{solution_path}(t, (t_1, \dots, t_n))$.*

From the theorem it follows easily that T does not contain an open branch (i.e. T is closed) iff q is not satisfiable in S , as expected.

In the statement of the theorem it is implicitly assumed that the tableau construction by $KRHyper$ terminates. That this is indeed the case was argued for in Section 3.3. The other properties stated can be shown with a careful analysis of the properties of the transformation of S and q . The idea is show that traversing the schema according to a given solution is simulated by the model construction, and vice versa. An important detail concerns solutions representing circles in the schema graph. As explained at the end of Section 3.4, following circles *once* is admissible, which matches exactly with what the blocking technique in Section 3.3 achieves.

3.6 Example

The query expression $\text{BOOK}[b]!\text{title}[t]$ from above translates into the following program:

$$\begin{aligned}
 \top(x) &\leftarrow & (1) \\
 \text{root}(\text{init}) &\leftarrow & (2) \\
 \text{solution_path}(b, t) &\leftarrow \text{role_filler_ref_trans}(\text{init}, b), & (3) \\
 &\quad \text{BOOK}(b), \text{role_filler_ref_trans}(b, z), \text{title}(z, t) \\
 \text{false} &\leftarrow \text{not some_solution_path} & (4) \\
 \text{some_solution_path} &\leftarrow \text{solution_path}(x, y) & (5) \\
 \text{role_filler_ref_trans}(x, x) &\leftarrow & (6) \\
 \text{role_filler_ref_trans}(x, z) &\leftarrow \text{role_filler}(x, y), & (7) \\
 &\quad \text{role_filler_ref_trans}(y, z) \\
 \text{role_filler}(x, y) &\leftarrow \text{university}(x, y) & (8) \\
 &\quad \vdots & (9) \\
 &\quad \vdots \\
 \text{role_filler}(x, y) &\leftarrow \text{isbn}(x, y) & (10)
 \end{aligned}$$

Notice the rule (3) is the translation of the given (single conjunct) query expression according to the scheme (1) above. The rules starting from (6) stem from the translation of the schema graph in Figure 2 according to the scheme (3).

Now, suppose that this rule set is combined with the translation of the schema graph in Figure 2 according to Definitions 2 and 3 (or Definition 4 instead). The unique model contains

$$\begin{aligned} \text{solution_path}(&f_{\text{BOOKS,books,BOOK}}(f_{\text{LIBRARY,book,BOOKS}}(\\ &f_{\text{UNIVERSITY,library,LIBRARY}}(f_{\text{ROOT,university,UNIVERSITY}}(\text{init}))))), \\ &f_{\text{PUBLICATION,title,STRING}}(f_{\text{BOOKS,books,BOOK}}(f_{\text{LIBRARY,book,BOOKS}}(\\ &f_{\text{UNIVERSITY,library,LIBRARY}}(f_{\text{ROOT,university,UNIVERSITY}}(\text{init})))))) \end{aligned}$$

Observe that the path from the ROOT concept to the BOOK concept is coded in the names of the Skolem function symbols in the first argument of `solution_path`. This path is extended towards a path to the title attribute in the second argument; this extension encodes, in terms of the schema graph, moving from the BOOK type to its supertype PUBLICATION and then moving to the title attribute. We note that the query expression `BOOK[b].title[t]` would have given the same result.

As a second example consider `PUBLICATION[p].isbn[i] \wedge \top [p].proceedings[x]` from above. Its translation according to scheme (1) is

$$\begin{aligned} \text{solution_path}(p, (i, x)) \leftarrow &\text{role_filler_ref_trans}(\text{init}, p), \\ &\text{PUBLICATION}(p), \text{isbn}(p, i), \\ &\top(p), \text{proceedings}(p, x) \end{aligned} \quad (3)$$

The rest of the transformation is the same as in the previous examples and is omitted. This time, the `solution_path` relation is empty, as the body of rule (3) cannot be satisfied. But then, with rules (4) and (5) this rule set is unsatisfiable. This is the expected result, because, as mentioned above, the query expression is unsatisfiable (in terms of the schema graph).⁷

4 Conclusion

In this paper we aimed to demonstrate that automated deduction techniques, in particular those following the model computation paradigm, are very well suited for database schema/query reasoning. We started by showing how to represent an XML schema graph in basic description logic. This representation then is transformed into the predicate logic language of the first order theorem prover *KRHyper*. We also proposed a flexible query language containing constructs for path specification similar as in the XPath query language, and we showed how to transform the query language into *KRHyper*'s input language. The models computed by *KRHyper* then encode fully specified XML path queries. The purpose of computing fully specified XML path queries is to reduce the workload on a query processor, because, in general, fully specified queries involve much less search on the database than XPath queries. To our knowledge an approach comparable to ours, based on model based deduction, has not been considered before.

An obvious question concerns the correctness of our approach, i.e. termination, soundness and completeness of the combination of *KRHyper* and the transformed

⁷ The runtime of *KRHyper* on both examples is negligible.

database schema and query. While termination has been argued for in Section 3.3, it is hardly possible to establish the other two properties: there is no formal semantics of XML path queries in the literature. However,

We found description logic to be helpful as an intermediate language to represent XML schema graphs. Moreover, our transformation includes a blocking technique comparable to those used in description logic reasoners [14, e.g.]. However, we went beyond the “classical” description paradigm in two aspects. The first aspect concerns the use of *models* to represent the solutions of the given problem of computing paths (as explained). As a second aspect, we found *default negation* very helpful to formulate constraints on acceptable models. From this point of view, we believe having shown there are applications for description logics where model computation and default negation is an issue. We perceive our approach to knowledge representation also as an original contribution of our work. A similar direction was proposed only recently in [2], where related ideas as presented here are exploited in the context of natural language processing. In a wider context, we speculate that *Semantic Web* applications would profit from knowledge representation languages as discussed here.

We understand this paper as a first step, that covers a basic schema/query reasoning task by model-based deduction. Due to the underlying expressive logic formalism we expect our approach to easily adapt to more sophisticated problem settings, like type hierarchies as they evolve within the XML world. One big advantage of such a declarative approach over, say, explicitly programmed algorithms is the possibility to easily add further constraints. We intend to explore this potential in future work.

References

1. F. Baader, D. Calvanese, D. McGuinness, D. Nardi, and P. Patel-Schneider, editors. *Description Logic Handbook*. Cambridge University Press, 2002.
2. P. Baumgartner and A. Burchardt. Logic Programming Infrastructure for Inferences on FrameNet. In *The European Conference on Logics in Artificial Intelligence*, LNAI. Springer Verlag, Berlin, Heidelberg, New-York, 2004. To appear.
3. P. Baumgartner, P. Fröhlich, U. Furbach, and W. Nejdl. Semantically Guided Theorem Proving for Diagnosis Applications. In M. E. Pollack, editor, *15th Int. Joint Conf. on Artificial Intelligence (IJCAI 97)*, pages 460–465, Nagoya, 1997. Morgan Kaufmann.
4. P. Baumgartner, U. Furbach, M. Gross-Hardt, and T. Kleemann. Optimizing the Evaluation of XPath using Description Logics. In *Proc. INAP2004, 15th Int. Conf. on Applications of Declarative Programming and Knowledge Management*, Potsdam, 2004.
5. P. Baumgartner, U. Furbach, M. Gross-Hardt, and A. Sinner. ‘**Living Book**’ :– ‘**Deduction**’, ‘**Slicing**’, ‘**Interaction**’. – system description. In F. Baader, editor, *CADE-19 – The 19th Int. Conf. on Automated Deduction*, LNAI. Springer, 2003.
6. P. Baumgartner, U. Furbach, and I. Niemelä. Hyper Tableaux. In *Proc. JELIA 96*, number 1126 in LNAI. European Workshop on Logic in AI, Springer, 1996.
7. W. Bibel and P. H. Schmitt, editors. *Automated Deduction. A basis for applications*. Kluwer Academic Publishers, 1998.
8. D. Calvanese, G. D. Giacomo, and M. Lenzerini. Answering queries using views in description logics. In *Proc. DL’99, Description Logic Workshop*, 1999.
9. D. Calvanese, G. D. Giacomo, and M. Lenzerini. Representing and reasoning on XML documents: a description logic approach. *J. of Logic and Computation*, pages 295–318, 1999.

10. J. Dix, U. Furbach, and I. Niemelä. Nonmonotonic Reasoning: Towards Efficient Calculi and Implementations. In A. Voronkov and A. Robinson, editors, *Handbook of Automated Reasoning*, pages 1121–1234. Elsevier-Science-Press, 2001.
11. T. Eiter, W. Faber, N. Leone, and G. Pfeifer. Declarative problem-solving using the DLV system. In *Logic-based artificial intelligence*, pages 79–103. Kluwer Academic Publishers, 2000.
12. U. Furbach. Automated deduction. In W. Bibel and P. Schmitt, editors, *Automated Deduction. A Basis for Applications*, volume I: Foundations. Calculi and Refinements. Kluwer Academic Publishers, 1998.
13. M. Gross-Hardt. Querying concepts — an approach to retrieve xml data by means of their data types. In *17. WLP - Workshop Logische Programmierung*, Technical Report. Technische Universität Dresden, 2002.
14. V. Haarslev and R. Möller. Expressive ABox Reasoning with Number Restrictions, Role Hierarchies, and Transitively Closed Roles. In *KR2000: Principles of Knowledge Representation and Reasoning*, pages 273–284. Morgan Kaufmann, 2000.
15. M. Kifer, W. Kim, and Y. Sagiv. Querying object-oriented databases. In *SIGMOD*, 1992.
16. R. Manthey and F. Bry. SATCHMO: a theorem prover implemented in Prolog. In E. Lusk and R. Overbeek, editors, *Proc. of the 9th Conf. on Automated Deduction, Argonne, Illinois, May 1988*, volume 310 of *LNCS*, pages 415–434. Springer, 1988.
17. I. Niemelä and P. Simons. Efficient implementation of the well-founded and stable model semantics. In *Proc. of the Joint Int. Conf. and Symposium on Logic Programming*, Bonn, Germany, 1996. The MITPress.
18. C. Sakama. Possible Model Semantics for Disjunctive Databases. In W. Kim, J.-M. Nicholas, and S. Nishio, editors, *Proc. First Int. Conf. on Deductive and Object-Oriented Databases (DOOD-89)*, pages 337–351. Elsevier Science Publishers B.V. (North-Holland) Amsterdam, 1990.
19. W3C. XPath specification. <http://www.w3.org/TR/xpath>, 1999.
20. W3C. XQuery 1.0: An XML query language. <http://www.w3.org/TR/xquery/>, 2001.
21. C. Wernhard. System Description: KRHyper. Fachberichte Informatik 14–2003, Universität Koblenz-Landau, Institut für Informatik, 2003.