

## RESEARCH OUTPUTS / RÉSULTATS DE RECHERCHE

### Database Wrappers Development: Towards Automatic Generation

Thiran, Philippe; Hainaut, Jean-Luc; Houben, Geert-Jan

*Published in:*

Proceedings of the Conference on Software Maintenance and Reengineering (CSMR'05)

*Publication date:*

2005

*Document Version*

Peer reviewed version

[Link to publication](#)

*Citation for pulished version (HARVARD):*

Thiran, P, Hainaut, J-L & Houben, G-J 2005, Database Wrappers Development: Towards Automatic Generation. in *Proceedings of the Conference on Software Maintenance and Reengineering (CSMR'05): CSMR*. IEEE CS Press, pp. 207-216. <<http://csdl.computer.org/comp/proceedings/csmr/2005/2304/00/23040207abs.htm>>

#### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal ?

#### Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

# Database Wrappers Development: Towards Automatic Generation

Philippe Thiran  
Technische Universiteit Eindhoven  
The Netherlands  
ph.thiran@tue.nl

Jean-Luc Hainaut  
Université de Namur  
Belgium  
jlh@info.fundp.ac.be

Geert-Jan Houben  
Technische Universiteit Eindhoven  
The Netherlands  
g.j.houben@tue.nl

## Abstract

*Wrapping databases allows them to be reused in formerly unplanned contexts, such as web-based applications or federated systems. Indeed, a wrapper can provide external clients of an existing (legacy) database with a neutral interface and augmented capabilities. However, except in simplistic cases where the wrapper and the database schemas are similar, the wrapper must implement complex mappings: it must translate queries from the wrapper data manipulation language to the database primitives, and, conversely, translate extracted data into the external wrapper format. We have developed a generic schema mapping framework in which wrappers can be specified formally and generated automatically. This framework comprises a high-level generic data model and a set of schema transformations defined for this model. This reference model makes it possible to specify different data models in a uniform formalism. Mappings between schemas are expressed as sequences of reversible schema transformations. We show how these transformations can be used to translate data and queries between two schemas and hence to generate as much as possible of the code of the wrappers. The generation is supported by DB-MAIN, a wide-spectrum CASE tool.*

## 1. Introduction

Existing data systems contain vital information that is embedded in existing (most often legacy) databases/ flat files and application code. In many cases, data systems include the only source of years of business rules and other valuable information. Access to this information is of vi-

tal importance to new open environments like the Web and to system integration in general.

A wrapper attempts to extend the usefulness of components of the existing data systems by facilitating their integration into modern (distributed) systems. A wrapper addresses the challenge of database heterogeneity by providing a standard and common interface. This interface is made up of: (1) a *wrapper schema* of the wrapped database, expressed in a canonical data model and (2) a common query language which uses the semantics defined in the wrapper schema. Queries on the wrapper schema are also known as *wrapper queries*.

Basically, database wrapping involves two, generally different, models, namely the database model (e.g., relational or standard files) and the wrapper model (e.g., object-oriented or XML). The main function of a wrapper is the translation of queries posed on the wrapper schema to the database model, and, conversely, the translating of data from the database model to the wrapper model.

This paper focuses on the aspects of query and schema translation within wrappers. Several research projects have already investigated this issue related to query mappings. Unlike their approaches, we investigate the problem from a model-independent and schema-oriented perspective:

- *Model-independent perspective.* Current approaches for wrapping databases rely on couples of models, such as those intended to produce XML views of relational schemas ([2], [3] or [13]). In this work, we use a general formalism to reversible schema transformations [7] based on a generic high-level data model. It provides a formal and uniform description of arbitrary models and the use of schema conversions between two not necessarily equivalent models. Here we extend this work by using our schema transforma-

tions to automatically wrap queries and data between two schemas.

- *Schema-oriented perspective.* Considering the issue of mapping definition and according to [12], two main basic approaches have been used to specify them. The first and very widespread approach ([1], [2], [3], [4] or [10]) is *query-oriented* in that it provides mechanisms by which users define wrapper schema constructs as views over source schema constructs, but do not focus on the semantics of the data sources. In contrast, the second approach ([11] or [14]) is *schema-oriented* in that mappings are defined as schema transformations that are used to automate the translation of queries. A comparison of these approaches is reported in [12]. The schema-oriented approach has the further advantage of decomposing the transformation of schemas into a sequence of small steps, whereas the query-oriented approach requires to directly define constructs in one schema in terms of those in the other schema.

In this paper, we extend the work in [16] and [17] on wrappers for legacy databases by exploring the query and schema mappings of such wrappers. We consider a generic schema mapping framework in which wrappers can be specified formally and generated automatically. This framework comprises a high-level generic data model and a set of schema transformations defined for this model. This reference model makes it possible to specify query and schema mappings in a uniform and unique formalism.

The paper is organized as follows. Section 2 presents the generic mechanisms of query mappings within a transformational approach. We then present in Section 3 the high-level generic data model that underpins our approach and the primitive transformations on schemas defined in terms of this model. In Section 4, we show how the schema transformations that result from our framework can be used to automatically translate queries. Section 5 deals with the development of databases wrappers based on schema transformations and its CASE support. Finally, Section 6 concludes this paper.

## 2. Transformational approach of query mapping

Query translation is the core function of a wrapper. It refers to operations that translate queries between two schemas (the database and wrapper schemas) and two languages (the database<sup>1</sup> and wrapper query languages).

Considering the issue of translating queries from one language to another one, our idea is to use an intermedi-

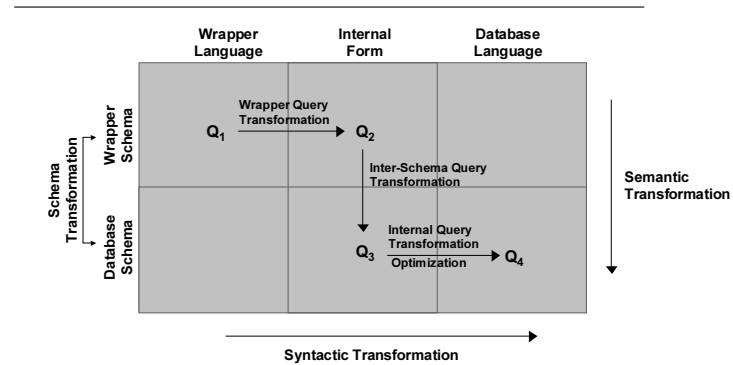


Figure 1. Language and schema mappings of a wrapper query  $Q_1$  into a DML query  $Q_4$ .

ate level, independent of all the possible operational query languages. We therefore use an internal abstract query language as the bridge for the translation rather than directly translating wrapper queries into database queries.

Considering the issue of schema mapping, our approach is to use schema transformations that provide mechanisms for formally defining schema correspondence between the database and wrapper schemas, and then, on using that equivalence to automatically perform the query mappings.

Figure 1 shows the translation process<sup>2</sup>. The wrapper query  $Q_1$  is first stripped off, creating an internal form  $Q_2$  that captures purely the semantics of the query. Next,  $Q_3$  is derived by application of the schema transformations on the constructs of  $Q_2$ . Finally,  $Q_3$  is translated into a query  $Q_4$  that complies with both the database schema and the database DML.

We can now state the three main successive steps of query translation:

- *Language mappings:* syntactic translation of the wrapper query into an internal form.
- *Inter-schema mappings:* semantics translation of the query using the schema transformation approach for defining the mappings between the database and wrapper schemas.
- *Language mappings and optimization:* syntactic translation of an internal form into a query based on the DBMS query language. Producing an efficient execution strategy depends on the syntax and expressiveness of both the wrapper (or internal) and DBMS query processing capabilities. Dealing with such issues is out of

<sup>1</sup> Or more precisely the Data Management System (DMS) language, generally called Data Manipulation Language, or DML for short.

<sup>2</sup> For simplicity, we make the hypothesis that the wrapper query language and the internal query language are equivalent. We refer to [5] for a discussion about the problem of query equivalence.

the scope of this paper however. In [16], we present some strategies for implementing query processors and optimizers in wrappers dedicated to COBOL systems.

In Section 3, we describe the formal framework of reversible transformations based on a generic data model. In Section 4, we present the internal query language based on the same data model. As we will see through these sections, reversible transformations allow internal queries to be automatically translated in either direction between two schemas.

### 3. Generic transformational framework

Query translation is a process that relies on mappings between schemas that are built within different paradigms. In the proposed approach, database and wrapper schemas are expressed in a unique wide spectrum specification model, the so-called *Generic Entity-Relationship model* (GER), from which the operational data models can be derived by specialization, that is, by selecting a subset of concepts and by defining restrictive assembly rules. As a result, it provides an ideal support for our a query translation approach based on schema transformations. Indeed, any transformation can be used whatever their underlying data model. For instance, the same schema transformation can be used in a relational schema and in an ER schema.

This section gives a short overview of the model and of the transformation techniques. More details of this approach can be found in [14].

#### 3.1. Generic Entity-Relationship model

For the need of this paper, the GER can be perceived as an enriched variant of the standard entity-relationship model. It includes the concepts of *entity type*, *attribute*, *value domain* and *relationship type*. Attributes can be atomic or compound, mandatory or optional, single-valued or multivalued. The roles of a relationship type can be labelled; it has a cardinality constraint (a pair of integers stating the range of the number of relationships in which any entity can appear). An attribute has a cardinality constraint too, that states how many values can be associated with each parent instance (default is 1-1 and does not appear in graphical schemas). In general, several properties hold, and must be declared, among the components of an entity type: uniqueness, referential and existence constraints are just some of them. Due the wide variety of such properties, the GER includes the generic concept of property *group*, or group for short. A group is any subset of components (attributes and/or roles) of an entity type on which one or several properties are defined. The label(s) of the group

specifies its properties (*id* for identifier, *ref* for referential, *excl* for exclusive, and so on). For example, a group of attributes of entity type  $E$  can be declared identifier and referential. This group models such relational pattern as a primary key that simultaneously is a foreign key.

This generic data model can be specialized into any operational model. A specialized model is built by selecting generic constructs and structural constraints, and by renaming constructs to make them comply with the concept taxonomy of the specialized model. As an illustration, the relational model, considered as an operational database model, can be precisely defined as follows (standard ER, UML class diagrams, IMS, Cobol, OO or XML DTD and Schema<sup>3</sup> can be defined in the same way):

- *Selecting constructs.* We select the following constructs: entity types, domains, attributes, identifiers and reference attributes.
- *Structural constraints.* An entity type has at least one attribute. The valid attribute cardinalities are  $[0-1]$  and  $[1-1]$ . An attribute must be atomic.
- *Renaming constructs.* An entity type is called a table, an attribute is called a column, an identifier, a key and a group of reference attributes, a foreign key.

#### 3.2. Mapping definition

A transformation consists in deriving a target schema  $S'$  from a source schema  $S$  by replacing construct  $C$  (possibly empty) in  $S$  with a new construct  $C'$  (possibly empty).

More formally, considering instance  $c$  of  $C$  and instance  $c'$  of  $C'$ , a transformation  $\Sigma$  can be completely defined by a pair of mappings  $\langle T, \tau \rangle$  such that  $C' = T(C)$  and  $c' = \tau(c)$ .  $T$  is the structural mapping, that explains how to replace construct  $C$  with construct  $C'$  while  $\tau$ , the instance mapping, states how to compute instance  $c'$  of  $C'$  from any instance  $c$  of  $C$ .

**3.2.1. Inverse transformation.** Each transformation  $\Sigma_1 \equiv \langle T_1, \tau_1 \rangle$  can be given an inverse transformation  $\Sigma_2 \equiv \langle T_2, \tau_2 \rangle$ , usually denoted  $\Sigma^{-1}$ , such that, for any structure  $C$ ,  $T_2(T_1(C)) = C$ .

So far,  $\Sigma_2$  being the inverse of  $\Sigma_1$  does not imply that  $\Sigma_1$  is the inverse of  $\Sigma_2$ . Moreover,  $\Sigma_2$  is not necessarily reversible. These properties can be guaranteed only for a special variety of transformations<sup>4</sup>, called symmetrically reversible.  $\Sigma_1$  is said to be a symmetrically reversible trans-

3 In [15], we show how XML structures can be represented in terms of the GER.

4 In [7], a proof system has been developed to evaluate the reversibility of a transformation.

formation, or more simply semantics-preserving, if it is reversible and if its inverse is reversible too.

From now on, unless mentioned otherwise, we will work on the structural part of transformations, so that we will denote a transformation through its T part.

**3.2.2. Some typical transformations.** We propose in Figure 2 the most common transformational operators. In particular, these transformations are sufficient to carry out the transformation of most ER schemas into the relational schemas [6], and conversely. Experience suggests that a collection of about thirty of such techniques can cope with most database engineering processes, at all abstraction levels and according to all current modelling paradigms.

**3.2.3. Structural analysis of schema transformations.** A transformation is known to replace construct C with construct C' in schema S, to yield new schema S'. The effect of a transformation T in schema S can be specified as follows. We define a schema S as a set of constructs. Therefore, set-theoretic relations and operators apply on schemas. Let us consider the structural functions C<sub>-</sub>, C<sub>+</sub> and C<sub>0</sub>:

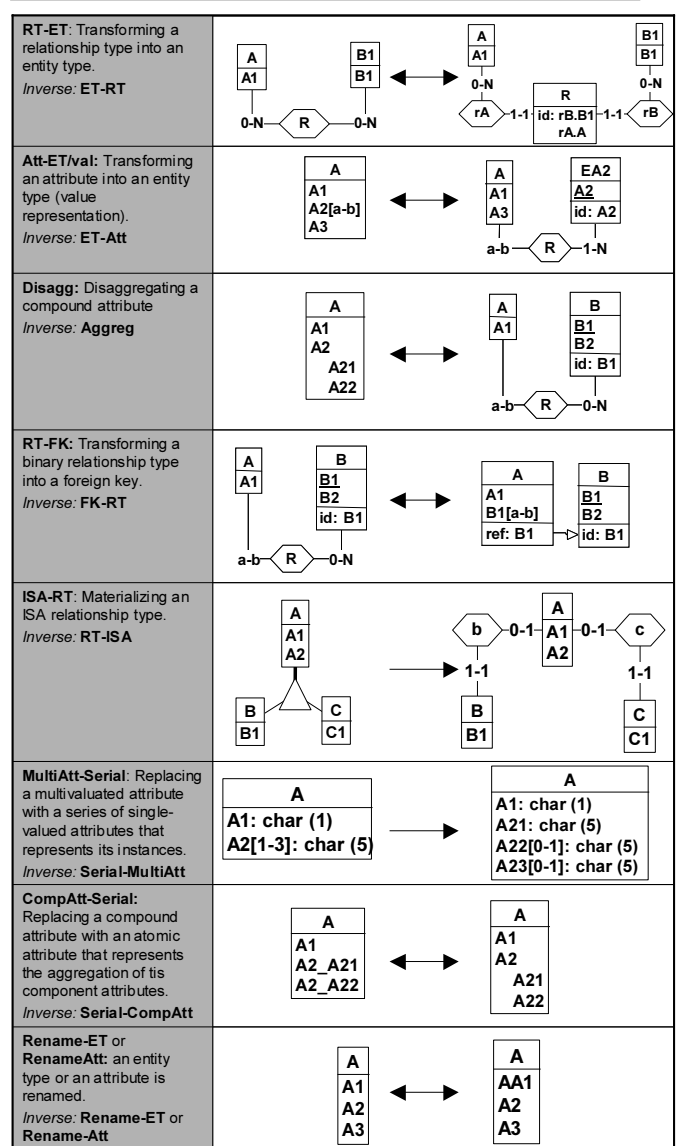
- C<sub>-</sub> returns the constructs of S that have disappeared in S';
- C<sub>+</sub> returns the new constructs that appear in S';
- C<sub>0</sub> returns the constructs of S that are concerned by T, but that are preserved by transformation (the catalytic constructs of T).

**3.2.4. Transformation sequence.** A transformation sequence is a list of n primitive transformations: S<sub>1</sub>-to-S<sub>2</sub> = (T<sub>1</sub> T<sub>2</sub> ... T<sub>n</sub>). For instance, the application of S<sub>1</sub>-to-S<sub>2</sub> = (T<sub>1</sub> T<sub>2</sub>) on a schema S<sub>1</sub> consists of the application of T<sub>2</sub> on the schema that results from the application of T<sub>1</sub>, so that we obtain S<sub>2</sub>.

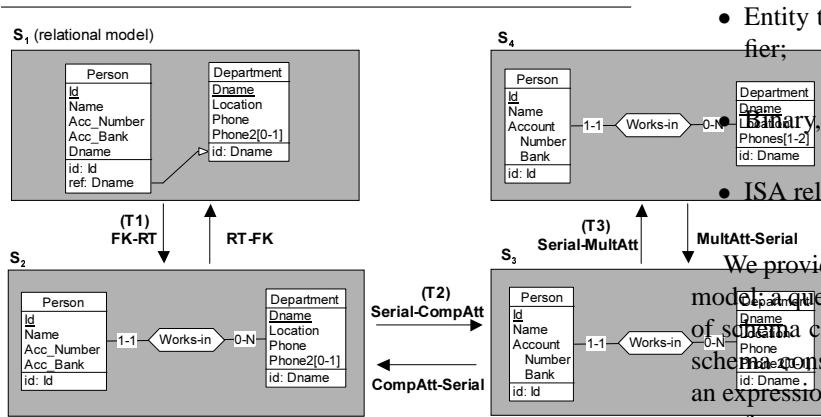
As for schema transformation, a transformation can be inverted. The inverse sequence S<sub>2</sub>-to-S<sub>1</sub> can be derived from the sequence S<sub>1</sub>-to-S<sub>2</sub> and can be defined as follows: if S<sub>1</sub>-to-S<sub>2</sub> = (T<sub>1</sub> T<sub>2</sub> ... T<sub>n</sub>) then S<sub>2</sub>-to-S<sub>1</sub> = (T<sub>n</sub><sup>-1</sup> ... T<sub>2</sub><sup>-1</sup> T<sub>1</sub><sup>-1</sup>) where T<sub>i</sub><sup>-1</sup> is the inverse of T<sub>i</sub>; and hence S<sub>1</sub> = S<sub>2</sub>-to-S<sub>1</sub>(S<sub>2</sub>). In other words, S<sub>2</sub>-to-S<sub>1</sub> is obtained by replacing each origin schema transformation by its inverse and by reversing the operation order.

The concepts of sequence and its inverse are used for defining the mappings between two schemas. The transformational approach then consists in defining a (reversible) transformation sequence which, applied to the source schema, produces the target schema.

As an illustration, Figure 3 shows a sequence of three transformations often used in database engineering process. The first one (FK-RT) replaces a foreign key into a relationship type, the second one (Serial-CompAtt) aggregates two attributes and the third one (Serial-MultiAtt)



**Figure 2. Major generic schema transformations with their inverse. Entity type and attribute names as well as cardinalities a,b,c, d must be replaced with actual values.**



**Figure 3. Sequence of schema transformations: a foreign key transformation followed by an aggregation transformation and a transformation of serial attributes into a multi-valued one.**

transform a serie of single-value attributes into a multivalued attribute.

**3.2.5. Model translation.** A model translation is a particular case of schema conversions [9]. It consists in translating a schema expressed in a data model  $M_s$  into a schema expressed in another data model  $M_t$  where  $M_s$  and  $M_t$  are two different submodels (i.e., subsets) of GER.

## 4. Schema and query mapping

In this section, we show how a schema transformation sequence can be used to automatically translate queries between a pair of schemas. More precisely, for a schema transformation sequence between two schemas  $S_1$  and  $S_2$ , we show how this sequence can be used to automatically translate queries posed on  $S_2$  to queries posed on  $S_1$ .

### 4.1. Model and query language

For simplicity and clarity, we consider a binary model defined as a sub-model of the generic data model described above. This model is compliant with standard files, SQL2 and ER models. It is expressive and generic enough to describe all the main structures and constraints that are explicitly offered by these data models:

- Atomic or compound attributes; single-valued or multivalued attributes;
- Reference, identifier and access groups;

- Entity types with at least one attribute and one identifier;
- Binary, non cyclic relationship types, without attribute;
- ISA relations.

We provide a simple query language based on this binary model. A query (named *Query* here below) is a conjunction of schema constructs. A query answer is a set instances of schema constructs. Any query *Query* over a schema  $S$  is an expression whose variables are constructs of  $S$ . The syntax of a query is:

**Query** ::= Construct | Predicate | [and, Query, Query {, Query}] | [or, Query, Query {, Query}] | [not, Query]  
**Predicate** ::= [eq, Atom, Atom] | [less, Atom, Atom]

Construct identifies a schema construct being added or deleted by a transformation. In other words, this is one of the constructs that take part in the definition of a schema transformation signature. Construct includes variable(s) used to instantiate instances of the construct and it takes one of the forms presented in Table 1. The underscore character is an anonymous variable. Atom represents a variable declared in a schema construct. When eq refers two variables of the same query, we can simplify the query and omit this predicate, e.g. we need only write [att, Person, Id, EP, 4] instead of [and, [att, Person, Id, EP, ID], [eq, [ID, 4]]]. In table 2, we illustrate the Construct representations of three constructs of Figure 3.

Construct	Syntax	Semantics
Entity type	[ent, Name, Et]	represents an entity type called Name, and Et can be instantiated with instances of Name
Attribute	[att, OwnerName, AttName, Owner, Att]	represents an attribute AttName of a construct OwnerName. The type of OwnerName can be either an entity type or a compound attribute. OwnerName contains the name(s) of the parent(s) of the attribute. Att can be instantiated with a value of the attribute associated with the instance Owner of OwnerName
Relationship type	[rel, ET1Name, RTName, ET2Name, ET1, ET2]	represents a relationship RTName between entities ET1Name and ET2Name. ET1 and ET2 can be instantiate with entity instances involved in the relationship

**Table 1. Syntax and semantics of the main constructs of the generic data model.**

Construct	Query Language Representation
Entity type Person	[ent, Person, EP]
Attribute Id of Person	[att, Person, Id, EP, ID]
Relationship type Works-in	[rel, Person, Works-in, Department, EP, ED]

**Table 2. Examples of construct representation (EP, ED and ID represent variables).**

## 4.2. Schema transformation and query substitution

Let us assume that a schema  $S_1$  is transformed into a schema  $S_2$  and the queries posed on  $S_1$  have to be translated to queries posed on  $S_2$ . Consider first the case where  $S_1$  is transformed into  $S_2$  by a single primitive transformation  $T$ . The only cases we need to consider in order to translate a query  $Q_1$  posed on  $S_1$  to an equivalent query  $Q_2$  on  $S_2$  are to apply renamings and to substitute occurrences of constructs of  $C_-(T)$  (Table 3). For transformation sequences, the substitutions are successively applied in order to obtain the final query  $Q_2$ .

Transformation	Signature	Substitution
RenameET	(name') ← RenameET(name)	$Q_2 = [\text{name}'/\text{name}] Q_1$
RenameAtt	(name')← RenameAtt(ET, name)	$Q_2 = [\text{ET}, \text{name}'/\text{ET}, \text{name}] Q_1$
Other	(S2)← T(S1)	$Q_2 = [C_-(T)/\text{query}] Q_1$
Other (inverse)	(S2)← T <sup>-1</sup> (S1)	$Q_1 = [C_-(T^{-1})/\text{query}] Q_2$

**Table 3. Schema transformation and query substitution.**

## 4.3. Illustration

We illustrate these notions by giving examples of query and update translation between two schemas. We consider the primitive schema transformations  $T_1$ ,  $T_2$  and  $T_3$  and their inverse between the pair of schemas  $S_1$  and  $S_4$  illustrated in Figure 3.  $T_1$ ,  $T_2$ ,  $T_3$  and their inverse are defined in Figure 4 below by means of: (1) their name (2) their structural function  $C_-$  expressed in the schema form; and (3) the queries  $query$  that state how the extents of each constructs of  $C_-$  can be recovered from the extents of the remaining schema constructs  $C'$ .

For any query on  $S_4$ , the table of Figure 4 can be used to translate constructs of  $S_4$  (the wrapper schema) into ones on  $S_1$  (the database schema), resulting in query on  $S_1$ .

T1: FK-RT and its inverse T1 <sup>-1</sup> : RT-FK		Direct transformations
C <sub>-</sub> (T1) direct	Query	
[att, Person, Dname, EP, X]	[and, [rel, Person, Works-in, Department, EP, ED], [att, Department, Dname, EP, X]]	1
C <sub>-</sub> (T1 <sup>-1</sup> ) inverse	Query	3
[rel, Person, works-in, Department, EP, ED]	[and, [att, Person, Dname, EP, X], [att, Department, Dname, ED, X]]	
T2: Serial-CompAtt and its inverse T2 <sup>-1</sup> : CompAtt-Serial		2
C <sub>-</sub> (T2) direct	Query	
[att, Person, Acc_Number, EP, X]	[and, [att, Person, Account, EP, AC], [att, Account, Number, AC, X]]	2
[att, Person, Acc_Bank, EP, X]	[and, [att, Person, Account, EP, AC], [att, Account, Number, AC, X]]	
C <sub>-</sub> (T2 <sup>-1</sup> ) inverse	Query	2
[and, [att, Person, Account, EP, AC], [att, Account, Number, AC, X]]	[att, Person, Acc_Number, EP, X]	
[and, [att, Person, Account, EP, AC], [att, Account, Number, AC, X]]	[att, Person, Acc_Bank, EP, X]	
T3: Serial-MultiAtt and its inverse T3 <sup>-1</sup> : MultiAtt-Serial		3
C <sub>-</sub> (T3) direct	Query	
[att, Department, Phone, ED, X]	[att, Department, Phones, ED, X]	3
[att, Department, Phone1, ED, X]	[att, Department, Phones, ED, X]	
C <sub>-</sub> (T3 <sup>-1</sup> ) inverse	Query	1
[att, Department, Phones, ED, X]	[or, [att, Department, Phone, ED, X], [att, Department, Phone1, ED, X]]	
		Inverse transformations

**Figure 4. Example of schema transformations and the queries that state how the extents of each construct of  $C_-$  can be recovered from the extents of the remaining schema constructs  $C'$  (T).**

Translation of a query from  $S_4$  to  $S_1$  (Figure 5). "Find the persons that are reachable via phone number 040-303030" is translated into equivalent queries in  $S_4$  and in  $S_1$  by applying the substitutions 1 and 3.

## 5. Wrapper development

Since the mapping between wrapper and database schemas is formally defined, we can expect them to be a sound basis to build the wrapper in a systematic way. Indeed, while the structural mapping  $T$  of a transformation defines a rewriting rule that can be used to transform the input query, its instance mapping  $\tau$  states how the instance of the target construct can be derived from that of the source construct. Therefore, these mappings can be used to define the query translation logic and the data transformation rules of the wrapper that implements this transformation. This analysis is still valid for transformation

Query on S4	Query on S1
<pre>[and, [att, Person, Id, EP, ID], [rel, Person, Works-in, Department, EP, ED], [att, Department, Phone, ED, '040-303030']]</pre>	<pre>[and, [att, Person, Id, EP, ID], [att, Person, Dname, EP, X], [att, Department, Dname, ED, X], [or, [att, Department, Phone, ED, '040-303030'], [att, Department, Phone1, ED, '040-303030']] ]</pre>

**Figure 5. Example of a query translation from S4 to S1.**

sequences, so that complete wrappers can be formally specified by such sequences.

## 5.1. Baselines

Practically, our approach of wrapper development can be summarized as follow:

- *Semi-hardcoded wrapper.* The wrappers are developed as program components dedicated to a specific database model and to a specific database. It comprises two parts, namely a model layer, in which the aspects specific to a given data model (e.g., RDB or standard files) are coped with, and a database layer that is dedicated to the specific database schema. While the model layer is common to all the databases built in this model, the wrapper/database schemas mapping is hardcoded rather than interpreted from mapping tables as is the case in other approaches.
- *Database reverse engineering.* Since most legacy databases have no associated documentation, the latter must be rebuilt through database reverse engineering techniques. These techniques yield all the necessary information to specify and develop the wrappers (i.e., the schema and mapping definitions). To accomplish this, we rely on a proven approach, namely the DB-MAIN DBRE methodology [6]. This approach has been already integrated in the wrapper development methodology presented in [16] and [17] and will not be discussed further in this paper. Its key feature is threefold. First all the schemas, whatever their modeling language, are expressed in the GER. Secondly, it uses the same transformational approach than that of this paper. Thirdly, it produces a wrapper schema that captures not only the structures explicitly expressed in the DDL code but also implicit structures implemented in the program codes (for instance, a foreign key can be explicitly declared in a modern relational schemas but they are

implicit (hidden) in old relational schemas – Oracle V5).

- *Operational CASE support.* the wrapper generation is supported by a specific module of the general-purpose CASE tool DB-MAIN.

## 5.2. Wrapper generation support

The wrapper generation is supported by the DB-MAIN tool, a general-purpose database engineering CASE and meta-CASE environment that offers sophisticated database application engineering toolsets. DB-MAIN includes advanced processors such as DDL parsers, transformation toolboxes, reverse engineering processors and schema analysis tools. In particular, DB-MAIN offers a rich set of transformational operators (including semantics-preserving ones) that offers a rich set of transformational operators (including semantics-preserving ones) that allow developers to define mappings in a systematic and formalized, though intuitive way. Another interesting feature of DB-MAIN is the meta-CASE layer, which allows method engineers to customize the tool and to add new concepts, functions, models and even new methods. In particular, DB-MAIN offers a complete development language, *Voyager 2*, through which new functions and processors can be developed and seamlessly integrated into the tool. Further details on DB-MAIN can be found in [8]. In the limited scope of this paper, we describe the two main *Voyager 2* programs dedicated to the wrapper code generation.

*History analyzer.* DB-MAIN automatically generates and maintains a history log of all the transformations that are applied when the developer carries out any engineering process such as wrapper schema definition. This history is completely formalized in such a way that it can be replayed, analyzed and transformed. An history basically is a procedural description of inter-schema mappings. The history analyzer parses history logs and transforms them into non-procedural annotations that define the inter-schema object mappings.

*Wrapper encoders.* The wrappers are automatically generated from the mapping annotations. Two wrapper interfaces are provided, namely SQL-based through a variant of JDBC, and object-based. At the current time, wrapper encoders for COBOL files and relational data structures are available. b

## 6. Conclusions

Data wrapping is one of the most powerful techniques to bridge existing (most often legacy) databases with modern architectures. Its main goal is model conversion, that addresses two inverse streams, namely query translation and data transformation. In this paper, we have focused on

the query translation process in database wrappers. This process relies on a special kind of inter-schema mapping, namely sequences of transformations defined on a unique and generic data model. By replacing the schemas constructs names in the wrapper query with their database equivalent, we produce a database query that can be executed on the actual data. This systematic approach can be automated, in such a way that wrappers can be generated based on the schema mappings. A specific plug-in has been developed for the DB-MAIN CASE tool. Considering two schemas and their mappings, expressed by a sequence of transformations, it generates wrappers for COBOL files and relational databases. The approach and the tool have been applied, among others, for building federated databases mixing both legacy and modern technologies.

## References

- [1] S. Bergamaschi, S. Castano, D. Beneventano, and M. Vinci. Retrieving and integrating data for multiple sources: the momis approach. *Data and Knowledge Engineering*, 36, 2001.
- [2] M. J. Carey, D. Florescu, Z. G. Ives, Y. Lu, J. Shanmugasundaram, E. J. Shekita, and S. N. Subramanian. XPERANTO: Publishing object-relational data as XML. In *WebDB (Informal Proceedings)*, pages 105–110, 2000.
- [3] M. Fernandez, W. Tan, and D. Suci. Silkroute: Trading between relations and xml, 2000.
- [4] H. Garcia-Molina, Y. Papakonstantinou, D. Quass, A. Rajaraman, Y. Sagiv, J. D. Ullman, V. Vassalos, and J. Widom. The TSIMMIS approach to mediation: Data models and languages. *Journal of Intelligent Information Systems*, 8(2):117–132, 1997.
- [5] G. Graefe. Query evaluation techniques for large databases. *ACM Computing Surveys*, 25(2):73–170, 1993.
- [6] J.-L. Hainaut. Introduction to database reverse engineering. Technical report, University of Namur, 2002.
- [7] J.-L. Hainaut. *Transformation of Knowledge, Information and Data: Theory and Applications*, chapter Transformation-based Database Engineering. IDEA Group, 2004.
- [8] J.-M. Hick, V. Englebert, J. Henrard, D. Roland, and J.-L. Hainaut. The db-main database engineering case tool (version 6.5) - functions overview. Db-main technical manual, Institut d’informatique, University of Namur, 2002.
- [9] D. Lee, M. Mani, F. Chiu, and W. W. Chu. NeT and CoT: Translating relational schemas to XML schemas using semantic constraints. In *ACM International Conference on Information and Knowledge Management*, 2002.
- [10] I. Manolescu, D. Florescu, and D. K. Kossmann. Answering {XML} queries over heterogeneous data sources. In *VLDB*, pages 241–250, 2001.
- [11] P. McBrien and A. Poulouvasilis. Automatic migration and wrapping of database applications - a schema transformation approach. In *International Conference on Conceptual Modeling / the Entity Relationship Approach*, pages 96–113, 1999.
- [12] P. McBrien and A. Poulouvasilis. Schema evolution in heterogeneous database architectures. In *CAiSE’02*, 2002.
- [13] J. Shanmugasundaram, J. Kiernan, E. J. Shekita, C. Fan, and J. Funderburk. Querying XML views of relational data. In *Proceedings of the 27th VLDB Conference*, 2001.
- [14] P. Thiran. *Legacy Database Federation - A Combined Reverse-Forward Approach*. Phd thesis, University of Namur, October 2003.
- [15] P. Thiran, F. Estievenart, J.-L. Hainaut, and G.-J. Houben. Exporting databases in xml a conceptual and generic approach. *Proc. of CAiSE Workshops (WISM04)*, 2004.
- [16] P. Thiran and J.-L. Hainaut. Wrapper development for legacy data reuse. In *WCRE Proceedings*, 2001.
- [17] P. Thiran, G.-J. Houben, J.-L. Hainaut, and D. Benslimane. Updating legacy databases through wrappers: Data consistency management. In I. C. Press, editor, *WCRE*, 2004.