

Data-Centric Intelligent Information Integration – From Concepts to Automation

Matthias Jarke (1,2), Manfred Jeusfeld (3), Christoph Quix (1)

(1) RWTH Aachen University, Informatik 5, Ahornstr. 55, 52074 Aachen, Germany

(2) Fraunhofer FIT, Schloss Birlinghoven, 53754 Sankt Augustin, Germany

(3) University of Tilburg, Netherlands

Email: jarke@dbis.rwth-aachen.de

Abstract. Intelligent integration of information continues to challenge database research for over 35 years. While data integration processes of all kinds are now reasonably well understood and widely used in practice, the growth and heterogeneity of data requires much higher degrees of automation to limit the need for human specialist work. This requires deeper insights in data-centric approaches of Enterprise Information Integration which focus on the semantics of information integration. Suitable recent formalizations and algorithms enable both significant improvement in schema integration, and in its automated transformation to efficient data-level integration, in a wide variety of architectural settings such as data warehouses or peer-to-peer databases. In addition to giving a short overview of developments in this field for the past twenty years, this paper focuses particularly on the challenges posed by heterogeneity in data models.

1 Introduction

Even after 35 years of database research, information integration remains one of its key challenges [SHT*77, BLN86, BDD*89, SSU90, SSU96, BBC*98, AAB*05, AAB*08]. Traditionally, there have been two different foci in this research. A focus on the *process of data integration*, for example the often cited ETL (extract- transform-load) processes in data warehousing, is nowadays a well-established industry sector called *Enterprise Application Integration (EAI)*. Among other factors, it was driven by resource constraints that often made the scheduling of integration tasks (in order to optimally exploit the buffer space available for intermediate integration results) the most critical bottleneck in information integration.

In contrast, the *data-centric approach* to information integration focuses on the semantics of the integrated information, in practitioner terms: the *data quality*. Classical techniques include the definition and analysis of formal constraints across data sources, in order to identify semantic conflicts or overlaps that could be exploited for data cleaning and integration. In industry, this strain – in the Internet propagated under the label of the Semantic Web – has caught on much later than EAI but is nowadays gaining importance under the label of *Enterprise Information Integration (EII)*. One important motivation of this is obviously the enormous costs of dealing with data quality problems. Even more important seems nowadays the advances in data mining algorithms that allow the exploitation of high-quality historical and sensor data for pattern recognition and prediction in science and industry.

For a long time, most integration tasks were solved manually, with rather limited formal or tool support. Recent practice surveys have claimed that about 40% of database-related work in industry is spent on data integration issues [Brod10], an issue so important that it caused top management attention with 68% of CEOs surveyed by IBM [Haas07].

The increasing complexity in terms of data volume, heterogeneity, and especially size and number of models, poses new challenges to the design and the development of integrated information systems. Fortune-500 enterprises now employ several thousand database systems with a few hundred relations each [Brod10]. In such settings, manual approaches for information integration are no longer feasible [Haas07, Smit07, BeHa08]. However, especially more automation in Enterprise Information Integration requires a deeper formalization of the semantic foundations e.g. in logic [Lenz02]. Even informal tasks such as the identification of shared or similar elements in different schemas must be formalized somehow, albeit with many different techniques, leading to the development of a complete new subarea of research and industry called *schema or ontology matching* [RaBe01].

In addition, while most database systems still use the relational data model, data sources and applications may include a broad range of data formats, informal media objects, and a wide variety of different modeling and metadata languages both in their operation and in their design and evolution. The heterogeneity

problem becomes even more intense in mobile multimedia applications with data as well as service integration requirements.

Since the turn of the century, the research area of *model management* [BHP00] therefore aims at high-level methods and automated systems to support the development of metadata-intensive applications. A typical example is the definition of a *model algebra* that provides high-level abstract operators for the key model-level tasks underlying data integration:

- *Match*: the identification of correspondences between models (match, [RaBe01, ShEu05]),
- *Compose*: the (possibly multi-step) transformation between models based on specifications of their inter-relationships as a formal mapping [MHH00, ABLM10],
- *Merge*: the integration of models (schema merge, [BLN86, PaSp98]), and
- the actual *execution* of the specified data transformations [MBHR05, HHH*05].

In this paper, we present a brief review of the evolution from classical data-centric integration to the recent advances in integrated model and data management enabled by more than 20 years of research in intelligent information integration. In section 2, we summarize classical data integration efforts up to the data warehouse movement of the late 1990's. In section 3, the evolution of model management from precursors in the mid-1990s until today is reflected with a discussion of some of the most important ideas and prototypes. At the boundary between both phases, we briefly review the history of our ConceptBase system for which a key paper appeared in JIIS 1995 [JGJS95], and became one of the most-cited papers in the journal's history.

Section 4 focuses on the challenge of automatically dealing with heterogeneity in model management. A formal metamodeling framework and model management toolset must simultaneously support the model management operators and their automated data-level executability for a broad range of modeling and operational data languages under which the systems-to-be-integrated might be operated or (re-)designed. As an example, we describe our Generic Role-Based Metamodeling suite and its underlying theory. In the final section 5, we present the application of this approach to classical data integration problems such as schema matching and schema merging.

As a running example, we choose a case study in mobile traffic data integration [GQSJ12] shown in Figure 1. Data streams from mobile devices or sensor networks have to be integrated with data from classical database systems and web services. For example, consider a traffic information system which makes use of various information sources to derive accurate information of a current traffic situation. In case of an accident or a traffic jam, cars send messages (Floating Car Data FCD¹ [KDH*05]) of the event or their current state to a traffic information system which integrates, aggregates, and analyzes the received messages in real time in the context of existing database information.

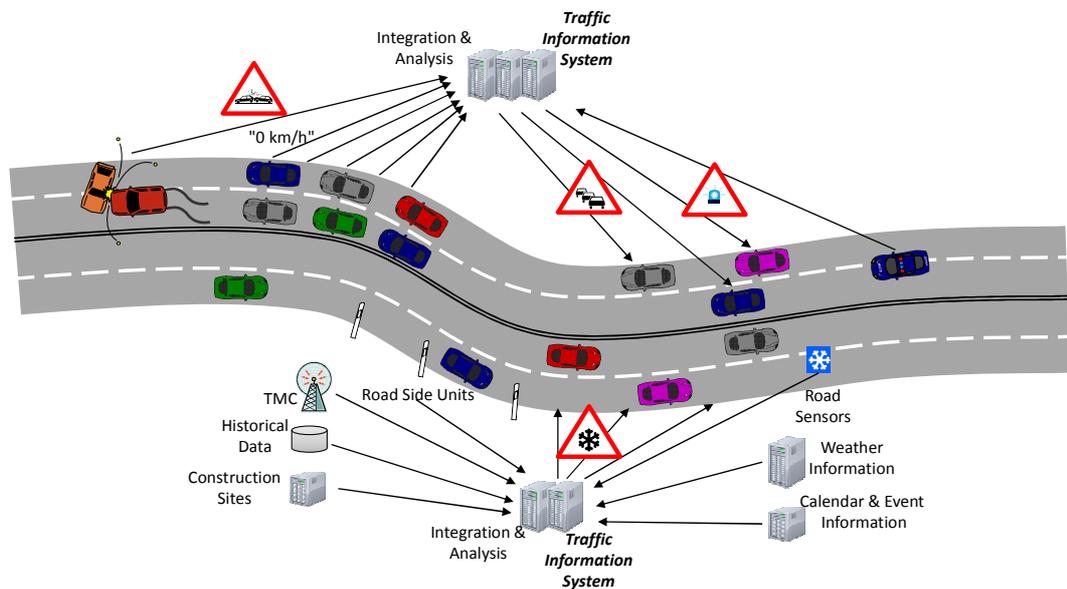


Figure 1. Information integration in a traffic information system using C2X communication

Such context data might come from heterogeneous external information sources, e.g., weather information to check the consistency of temperature data delivered by sensors in roads, road side units, or cars; traffic density information derived from aggregated C2X messages is complemented by data from TMC (Traffic Message Channel) or a database of construction sites; a baseline for the traffic state can be derived from historical information.

¹FCD is transmitted using some Car-To-Car (C2C) or Car-to-Infrastructure communication service (C2I, [SBH*10]). C2C and C2I communication is summarized under the term C2X (Car-to-X) communication.

2 Classical Data Integration

The classical procedural approach to data integration implicitly assumes a staged software architecture which was made explicit in the early 1990s as the so-called mediator architecture as depicted in Figure 2 [Wied92]. Data from several sources is integrated by a mediator which might follow a virtual or materialized integration approach.

In a virtual integration scenario first mentioned in the distributed database context by [CePe84], the mediator must reformulate the queries of the applications and integrate the data from sources on the fly. The application queries are expressed in terms of the global schema of the mediator. To retrieve the data from the source, they have to be translated into queries in terms of the local schemas of the sources. Wrappers take these reformulated queries, send them to the sources, extract the answers, and send the result back to the mediator. Wrappers may also apply some simple transformations such as translating the answers into a uniform format.

In the materialized integration approach, the data is stored in a central data repository, since the early 1990s called a data warehouse [JLVV03]. While research has focused on the fundamental and transformational aspects of this, very similar to the virtual integration scenario, industrial practice is at least equally interested in the resource-constrained scheduling of the huge bulk tasks involved in operating this architecture with the enormous data sizes of today. So-called ETL tools support the main steps within such a process which actually predated data warehousing by over ten years (EXPRESS [SHT*77]): *Extract* source data into some buffer, *Transform* them by cleaning, model unification, and merging (mediator), and *Load* them to the data warehouse.

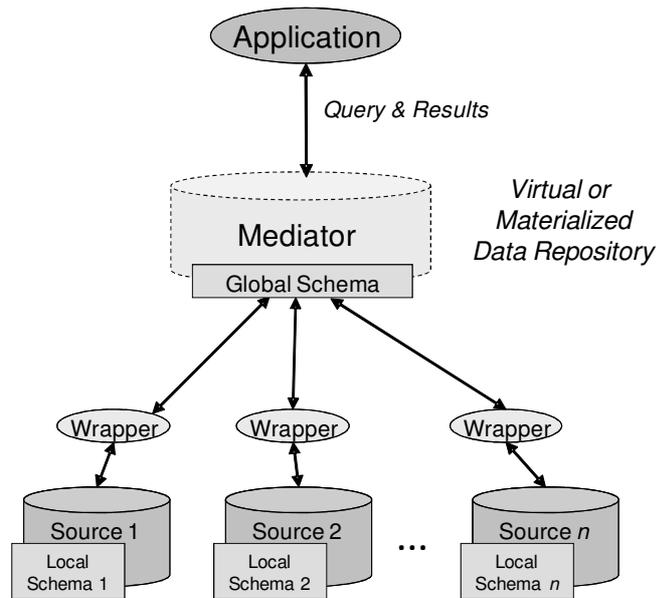


Figure 2. Architecture of an Information Integration System

On the formal side, integration methods in the 1990s began to study explicit formal mappings between the source and integrated data which could serve as a basis for “model-based” data integration where wrapper and mediator code could be generated largely automatically from the mapping specifications (e.g., Carnot [CHS91,SCJ*97], Infomaster [GKD97], Information Manifold [KLSS95], SIMS [AHK96]). While Infomaster and the Information Manifold focused on the Relational Data Model as modeling language, Carnot and SIMS used description logics [NaBr03] to express relationships on the model level as well as on the data level.

There are many ways to characterize the mappings between schemas. The simplest form is correspondences between individual attributes of the schemas. This type of mapping is frequently used in schema matching [RaBe01] but cannot be used directly for information integration because it does not capture more complex relationships that require restructuring and regrouping of data. Since the logic-based data warehouse research of the late 1990s, mappings are therefore often expressed as a set of query pairs. In each pair, a query q_S over a source schema S is related to a query q_G over the global schema G : $q_S \sim q_G$ [Lenz02]. The relationship “ \sim ” between the queries is a set relationship like $=$ or \subseteq , which means that the result set of the query q_S is equivalent to (a subset of, resp.) of the result set of query q_G for all valid database instances.

There are two semantic perspectives for formalizing the mappings, called *Local-As-View (LAV)* and *Global-As-View (GAV)*. In GAV mapping, any single element g of the global schema is defined as a view on the sources, i.e., $q_S \sim g$. This reflects the classical data integration perspective of the mediator approach. In the LAV approach, an element s of the source schema is defined as a view on the global schema, i.e., $s \sim q_G$. This reflects the idea of the integrated data as a partial description of a uniform “real world” about which the data sources capture perhaps incomplete, erroneous or even inconsistent observations; LAV therefore seems more relevant when we talk about important real-world problems such as semantic data quality.

At first glance, query rewriting in the case of GAV seems to be easier as a query over the global schema needs just to be unfolded, i.e., the elements of the global schema in the query are replaced with the corresponding query q_S over the sources [Lenz02]. However, in the case of constraints and incomplete sources (which is a common assumption in information integration) more complex reasoning is required to answer queries [CCGL04].

In LAV, query rewriting corresponds to the problem of answering queries using views [Hale01], which also requires reasoning; however, starting with MiniCon [PoHa01], a number of efficient algorithms have been developed. Some integration architectures which go beyond the mediator schema of Figure 2, such as Peer-to-Peer data management systems [HIM*04], require a combination of GAV and LAV mappings called GLAV.

In parallel to these developments, the growth of the Internet caused a rather different approach to intelligent information integration to emerge. The lack of central planning and the resulting irregularity of data structures in web-based systems require more flexible approaches for data management. TSIMMIS was one of the first projects that moved away from formal schemas and schema mappings as the basis for integration. Instead, it used self-describing semi-structured graph data models [GPQ*97], later replaced by the emerging (tree-oriented) XML standard. Similar to the approach taken by modern search engines, data integration in such a setting became a mix of text retrieval and graph mappings without higher-level schemata. Even though this kind of approach persists until today (as matching algorithms, see below in section 5.1), it cannot easily exploit the rich knowledge available in the schemas of the ten thousands of

databases and ontologies available today, or the documented design models which give even more semantics to these schemas. To address these challenges, model management emerged at the turn of the century.

3 The Evolution of Model Management

The creation of models and mappings in the classical data integration systems was largely a manual task. While automated support had been proposed for a few integration tasks (e.g., schema integration [BLN86] or schema matching [RaBe01]), the design, implementation, and maintenance of the integration system had to be done manually. This might have been acceptable for data management systems with a manageable schema complexity, but the systems have grown significantly in the recent years. The complexity of current “information ecosystems” [Brod10] with heterogeneity at various levels requires a methodical support for the management of models and mappings.

The importance of data models in the development of integrated information systems has been recognized by [BHP00] in their vision of *model management systems*. In such systems, models should be regarded as first-class objects, and the systems should provide operators to “work” with these models. For example, a *Match* operator should be used to compute a mapping between two models, the *Merge* operator should integrate models based on a given mapping, and the *ModelGen* operator should generate a new model by translating a given model into another modeling language. The vision was an algebra which would allow an abstract specification of complex operations on data models.

3.1 Model Management 0.1: Repositories with Formal Metamodels

The vision of model management in [BHP00] initiated new research in this area, but there has been significant research on individual topics before [Quix09b].

Early approaches to schema integration and matching are summarized in [BLN86]. They were mainly based on abstract, conceptual modeling languages (e.g., variants of the EER model), or directly operated on the relational schemas with intra-relational and inter-relational dependencies [CaVi83, BiCo86]. In both cases, the mapping languages were rather weak as only one-to-one correspondences could be expressed.

With the increasing size and complexity of database systems in the late 1980s, the necessity for formal methods for the management of complex data models became evident. Business IT researchers like [Dolk88] first stated the requirement for a theory for models similar to the relational database theory. Such a theory should include formal definitions of models and operations on models and could be used as a basis for the implementation of a model management system. This work was based on a draft of the *Information Resource Dictionary System (IRDS)* standard [IRDS90] which was accepted in 1990. IRDS clarified the terminology of modeling systems as a four-level hierarchy. A decade later, the Unified Modeling Language (UML) community adopted the same approach with slightly different terminology in its MOF standard (Meta Object Facility [MOF05]). The metamodel hierarchy in Figure 3 shows a MOF-based repository hierarchy for our running example: at the lowest level reside data instances which are described by a model (or schema) on the next higher level. The model is expressed in some modeling language (or metamodel) which is located at the third level. The highest level contains a metametamodel which can be used to define metamodels.

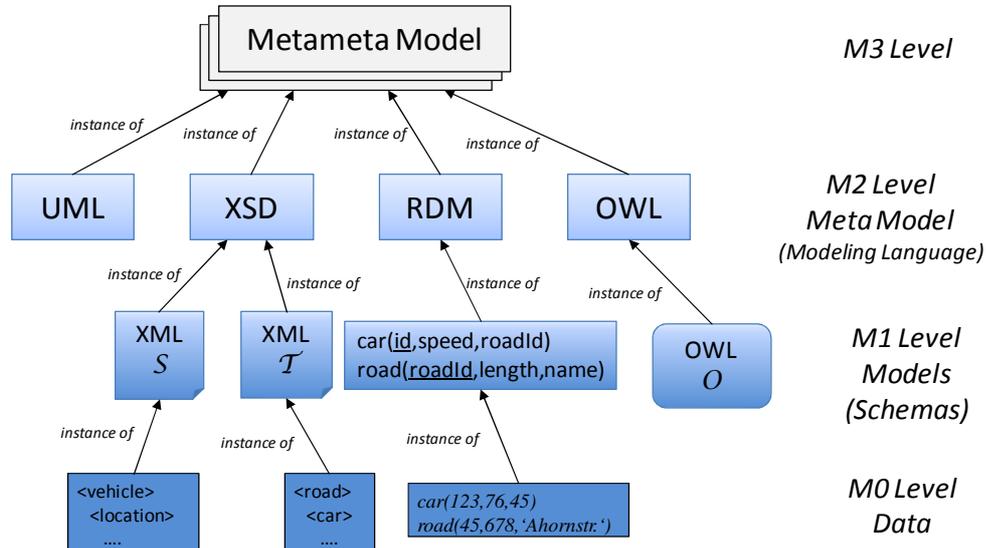


Figure 3. Metamodeling hierarchy according to the Meta Object Facility [MOF05]

The IRDS framework introduced the concept of so-called metadata repositories for purposes such as development process traceability, information integration, and model transformation [Quix2009a]. For example, solutions for forward and reverse engineering between ER models and relational databases were developed using generic metametamodels [AtTo96, JeJo95], i.e., a uniform representation at

the M3 level). Generic metamodels at the M2 level also addressed several aspects of database schema integration [SpPa94, PoBe03]. However, these early solutions mainly concentrated on the transformations at the model level and did not pay detailed attention to automated transformation on the data level.

In parallel to the initial IRDS standardization, a logic-based approach to dealing with an unbounded number of meta levels was investigated in the Telos project jointly conducted between the University of Toronto and several European projects in the late 1980s [MBJK90]. An important feature of Telos is the strong and highly efficient formalization in Datalog with stratified negation [Jeu92] from which excerpts are briefly reviewed in section 4.3. Based on this formalization, we developed the deductive metadatabase system ConceptBase whose complete description was published in JIS 1995 [JGJS95] and whose present version is still widely used in several thousand installations worldwide [JJN*09]. Based on the Datalog formalization, ConceptBase was the first metadata repository to also offer effective query optimization, integrity constraint evaluation, and incremental view maintenance at the data level [StJa00], as well as viewpoint resolution [NiJa99] and requirements traceability [RaJa01] at the meta level simultaneously, thus providing an early example of fully automated model-based code generation.

With the confluence of structured data, text and multimedia capabilities, the World Wide Web, and mobile communications, the range of data models has grown well beyond what could be covered by these early approaches. In [HaK110], an overview of metadata interoperability in heterogeneous media repositories is given. Although the survey focusses on media repositories, it also addresses interoperability between “structural” modeling languages, such as UML, XML Schema, and OWL. The authors classify approaches according to the MOF hierarchy and argue that effective interoperability between systems can only be achieved if data transformations at the instance level are also addressed. Furthermore, they distinguish between standardization and mapping approaches. The former propose metadata standards to enable interoperability, whereas the latter build relationships between different metamodels. Mapping approaches are more complex, but are advantageous in open environments such as the Web, as in these cases, no central authority can enforce a standard [HaK110].

3.2 Model Management 1.0: Algebra of Model Operators

The increasing complexity of information systems [Brod10] requires techniques for automating the tasks of creating models and mappings. The original vision of model management aimed at providing support for these tasks [BHP00], even though a complete automation was expected to be hard to achieve. The creation of models and mappings was considered a design activity which requires a deep understanding of the semantics of the modeled systems. Such tasks are AI-complete, i.e., it requires human intelligence to solve these problems [BMPQ04]. Another important motivation for the definition of a model management algebra was the observation that many applications that deal with models require a significant amount of code for loading and navigating models in graph-like structures. A model management system based on a formal algebra should simplify the development of model-oriented applications in the same way as data management systems based on relational algebra simplified the development of data-oriented applications [BHP00].

First model management systems such as Rondo [MRB03a, MRB03b] and COMA [DoRa02] applied simple, abstract model representations in which a model is represented as a directed, labeled graph. Other approaches for model management also focused on models, e.g., there is a huge research area on schema matching [RaBe01, ShEu05]. Although a graph representation is often sufficient for basic schema matching tasks, semantic details of the models (such as constraints) cannot be easily represented. Mappings are often just represented as a set of pairwise correspondences between nodes in the graphs. MISM (Model-Independent Schema Management) uses a richer representation for schemas [ABB*09]. Schemas are described in a generic way using a multi-level dictionary. However, the system uses a set-theoretic approach for some model management operators (e.g., *Merge*), i.e. again only correspondences ('equivalence views' in the terminology of MISM) are used as the mapping formalism.

In the original vision of model management [BHP00], mappings had a weak representation and were seen as a special type of a model which might include expressions to describe the semantics of a mapping in more detail (e.g., by using a SQL query). However, to automate operations on models and mappings, mappings have to be represented in a separate formalism which is more expressive than just simple correspondences.

3.3 Model Management 2.0: Mappings as First-Class Citizens

There have been several attempts aiming at combining a rich modeling language with powerful mapping languages. For example, in the European DWQ project (Foundations of Data Warehouse Quality [JLVV03]), a semantically rich metamodel [JJQV99] was combined with an information integration approach based on description logics [CGL*01]. Similarly, the Italian MOMIS system used an object-oriented modeling language to support the integrated querying of heterogeneous information sources [BCVB01].

The Clio project between IBM and the University of Toronto [HMH01, HHH*05, FHH*09] introduced a strong mapping language based on tuple-generating dependencies (tgds) [BeVa84, AHV95]. The well-defined, formal basis and the ability to easily translate the mappings into executable code (e.g., queries in SQL or XQuery) proved a significant advantage and caused a re-thinking of the whole definition of model management, dubbed Model Management 2.0 [BeMe07].

In model management 2.0, it has been realized that the *representation of mappings* is at least equally important as the representation of models [BeMe07]. Mappings are involved in all model management operations, and mappings are at the core of any integration approach. Furthermore, model management is not only a design time issue. The runtime system has also to be taken into account, because mappings have to be executed eventually to perform data transformation tasks. Thus, it is not sufficient to hide the semantics of a mapping in a string expression in some arbitrary language. Model management systems must be able to understand the semantics of a mapping in order to enable mapping operations (e.g., composition, inversion) and produce mappings as output (e.g., in match and merge operations). While Clio mostly explored this issue in the context of (nested) relational data models, the following section will investigate the extension to the management of heterogeneous data models.

4 Towards Heterogeneous Model Management

The heterogeneity of data management systems and modeling languages used in the web, but also in enterprise information systems, requires an integration approach, which is able to cope with the different modeling formalisms in a single

uniform framework. Moreover, this has to be done at the same time at the model level as well as the data level.

A rich modeling language is especially required for the definition of mappings between data models. A mapping states how the data of one model is related to the data of another model. It is important to note that mappings relate data and not only models. Because of this, mappings need to be very expressive in order to be able to represent rich data transformations. *Executability* of a mapping language means that it must be possible to apply a mapping such that it enables automatic code generation that executes the data transformations specified in the mapping.

Summarizing, a model management system needs to address at least the three lower levels of Figure 3, i.e., the data level for expressive data translations, the model level for operations on models such as *Match* and *Merge*, and the metamodel level to enable a generic representation of heterogeneous data models.

Defining a mapping language between all individual pairs of different modeling languages would be a daunting to impossible task, as for each pair, syntax and semantics of two individual formalisms have to be interlinked. A generic modeling language simplifies this task by providing a uniform basis for the definition of mappings. Model management operations have to be based on formal languages with rich semantics. These formalisms are necessary in order to support the development of model management systems which have to produce in the end mappings and models in existing, formal languages (e.g., SQL, XML Schema, XQuery). Moreover, a formal basis is required to prove characteristics of model management operations, e.g., that a model transformation is correct and complete, or a merged schema is minimal, but also preserves all information of the input schemas.

Last not least there is a trend that the restrictions of relational and even XML database systems are too tight for new applications that require high availability and scalability for the web [Voge07, Ston10]. Thus, new data management systems with another set of modeling languages are being introduced (e.g., NoSQL database systems [Catt10]), which again require transformations and mappings of existing data models. Consequently, the management of heterogeneous data models will be a running challenge also for future information integration projects, and model management systems must be extensible and flexible to support also future modeling languages.

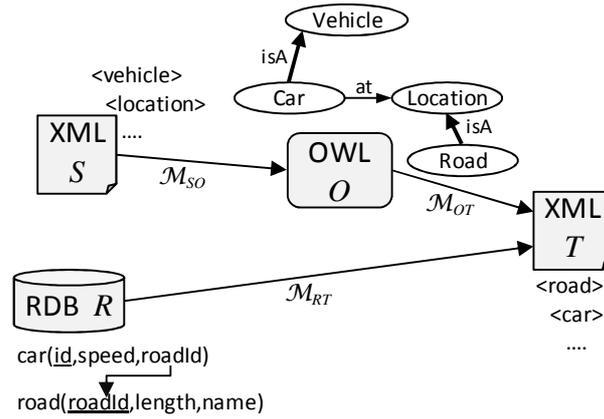


Figure 4. A simple, heterogeneous data integration scenario

To illustrate the concepts and algorithms in the following subsections, Figure 4 extracts a simplified heterogeneous schema integration setting from the scenario in Figure 1. An XML web service S provides information about vehicles and their location. A relational data stream R gives Floating car information about cars, their speed and the roads on which the cars are traveling. These two sources should be integrated into an XML database T , which has information about roads and for each road a list of cars which are driving on this road. The XML documents use a different vocabulary, therefore, we might use an ontology O as the semantic bridge between the XML documents. The relational schema R is mapped directly to the target XML schema T .

4.1 Requirements for a Generic Model Management System

A generic representation of models is a prerequisite for building a model management system. Without a generic representation, model operations would have to be implemented for each modeling language that should be supported by the system. Especially for the task of model transformation, a generic representation of models is advantageous as the necessary transformations have just to be implemented for the generic representation. Such a generic representation is called a *generic metamodel*. A generic metamodel should be able to represent models originally represented in different metamodels (or modeling languages) in a generic way without losing much detailed information about the semantics of the model.

First implementations of model management systems used rather simple graph representations of models, e.g., Rondo [MRB03b]. Although the graph-based

approach might allow an efficient implementation of operations which do not rely on a detailed representation of the models (such as schema matching), it makes it more difficult to implement more complex operations (such as model transformation or schema integration). Schema integration methods often use rather abstract metamodels extending traditional conceptual modeling formalisms, such as the ERC+ model in [SpPa94]. They usually focus on the conflicts at the semantic level, because such modeling languages are good at representing the semantics of a model. Schema integration approaches using a more concrete metamodel (e.g., the relational model with constraints and dependencies [CaVi83, BiCo86]) assume that these semantic conflicts have already been resolved and provide solutions for integrating schemas without conflicting constraints.

It is apparent, that a richer modeling language is required for model management operations that explicitly deal with the semantics of model elements in a heterogeneous setting. Early examples of such metamodels have been used, e.g., for model transformation [AtTo96, ACB06, MBM07].

In [AtTo96], the authors describe a metamodel consisting of “superclasses” of the modeling constructs in the native metamodels. The transition between this internal representation and a native metamodel is described as a set of patterns. This induces the concept of a *supermodel* which is the union of patterns defined for any supported native metamodel. This model representation has been expressed in a relational model dictionary [ACB05], and was used for the generic *ModelGen* implementation *MIDST* [ACB06].

In [JeJo95], a ConceptBase metametamodel to enable the translation of models between different modeling languages is used. A model element in a concrete modeling language is mapped to the metametamodel, in which models can be rearranged, and then translated into the desired target modeling language. Another metamodel following the approach of generalizing metaclasses is Vanilla [PoBe03] which has been used to implement model merging.

4.2 Role-Based Metamodeling with *GeRoMe*

Our metamodel *GeRoMe* provides a generic, yet detailed representation of data models originally represented in different languages [KQCJ07]. In its role-based modeling approach [BaDa77, RiSc91, WCL97], an object is regarded as playing roles in collaborations with other objects. This allows describing the properties of

model elements as accurately as possible while using only metaclasses and roles from a relatively small set. This strongly reduces a well-known problem in metamodeling, as follows.

A classical approach for modeling a generic metamodel could define a hierarchy of metaclasses that represent an abstraction of concrete modeling elements in existing modeling languages (e.g. [JeJo95]). One could then map the concrete modeling elements to exactly one of these metaclasses. However, modeling elements in different modeling languages have often similar or overlapping semantics, but rarely a truly equivalent semantics. When two model elements from different metamodels are mapped to the same metaclass in the generic metamodel, this implies that the elements have the same semantics in the view of the generic metamodel. Thus, differences and details of model elements are lost due to the abstraction in the generic metamodel. A solution to this problem could be to model the detailed semantics of model elements by intersection classes, i.e., classes which inherit features from several base classes representing basic modeling features (e.g., aggregation, inheritance, association). However, many intersection classes would be necessary to represent all the different combinations of modeling features which are present in one concrete modeling element. In earlier works, e.g., in our interdisciplinary research on conceptual modeling of chemical engineering processes [BaJa99, BMM*08], this has been shown to lead, even in practice, to a combinatorial explosion of subclasses.

GeRoMe's role-based approach represents each modeling feature by a separate role class. Model elements are plain objects which do not have any semantics by their own. By decorating a model element with role objects, a model element gains the modeling features of these role objects. This allows an arbitrary combination of modeling features.

The implementation of model management operators is simplified as it can focus on the roles which are relevant for a specific operator. Roles provide a view, i.e., a subset of the features of a model element. For example, when elements should be matched by name, the match operator needs to consider only the role providing the name of the element, all other roles can be ignored. Another advantage of the role based modeling approach is that roles – and thereby modeling features – can be easily added to or removed from a model element without changing its identity. This characteristic is in particular important for model transformation.

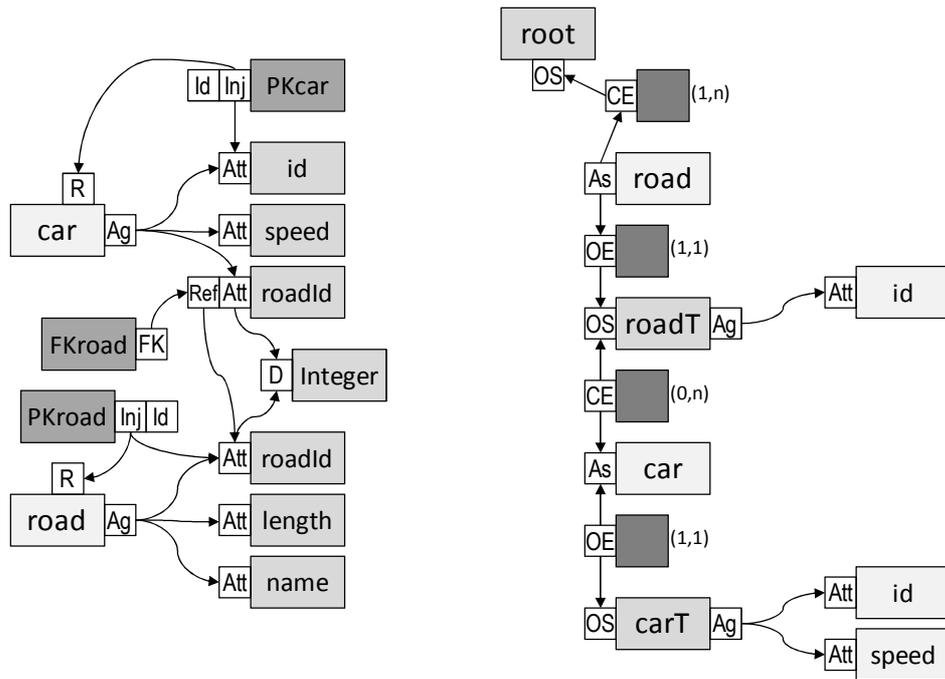


Figure 5. Simplified *GeRoMe* models for the relational and XML schema from Figure 4

Consider again the example in Figure 4. Simplified *GeRoMe* representations for the relational model R and the XML Schema T are shown in Figure 5. Gray boxes denote model elements, white boxes attached to them are role objects. The left part of the figure illustrates the *GeRoMe* model for the relational schema. The light gray elements `car` and `road` represent the model elements for the two relations `car` and `road`. Both elements play the role *Aggregate* (abbreviated by *Ag*), meaning that they can have attributes. In addition, they play the role *Referable* (*R*), which allows them to be the target of a key constraint. The attributes are represented by the medium gray model elements; all of them play the role *Attribute* (*Att*). In addition, `roadId` in `car` plays the role *Reference* (*Ref*) as it is a foreign key attribute. For simplicity, we show only the domain for the `roadId` attributes; it is the model element `Integer` which plays a *Domain* role (*D*).

Constraints are represented by dark gray elements. The primary key constraints `PKcar` and `PKroad` play the roles *Injective* (*Inj*) and *Identifier* (*Id*) as they represent uniqueness constraints which, in addition, identify the referenced type. Furthermore, the foreign key `FKroad` points to the reference role of the corresponding attribute.

The representation of the XML schema is more complex, as it contains more structural information. Complex types in XML-Schema are similar to classes in UML or entity types in the ER model, as they can have attributes and participate

in associations. Elements in XML represent associations, i.e. a relationship between two complex types (a nested and a nesting type). Because of XML's tree structure, an instance of a complex type is always nested into exactly one parent element. In the right part of Figure 5, the complex types `roadT` and `carT` also play `Aggregate` roles as the corresponding constructs in the relational model. In addition, they play `ObjectSet (OS)` roles as they also may participate in associations. In the example, `road` and `car` play the `Association role (AS)` and represent the XML elements. They link the complex types (and the document root) using association end roles, as indicated by unlabeled dark gray boxes. There are two different types of association ends in this example: `CompositionEnds (CE)` state the relationship to the parent element and `ObjectAssociationEnds (OE)` represent the link to the child element. These elements also maintain the information about the cardinality constraints of this association. Associations in UML or relationship types in ER are modeled in *GeRoMe* using the same set of roles.

4.3 Logical Foundation for the Generic Metamodel

Based on our experience with `ConceptBase` [JGJS, JJN*09], the formalization of *GeRoMe* represents a model as a set of logical facts. Each fact defines a model element, a role, a property, or a relationship between role objects. This logical representation allows declarative specifications and sound and efficient implementation of some model management operators.

A generic metamodel must provide a transformation between the concrete modeling languages and its generic representation. *GeRoMe* enables this import and export of models in a declarative way: equivalence rules state that a combination of modeling constructs in a particular modeling language are equivalent to a set of model elements and role objects in *GeRoMe* [KeQu07].

As in `ConceptBase`, the translation between *GeRoMe* and concrete modeling languages is based on Datalog rules [KeQu07]. However, the generation of modeling constructs – a key requirement in model transformation – requires a language that is more expressive than pure Datalog.

Figure 6 shows a simplified version of the rule for translating SQL columns into the corresponding *GeRoMe* elements. As it is an equivalence rule, it can be used

to import SQL schemas into *GeRoMe*, and to export *GeRoMe* models back into SQL schemas.

```
[
  sql_column(ID),
  sql_column_table(ID, TableID)
] <=> [
  modelElement(ID),
  attribute(ID),
  property(aggregate(TableID), attribute(ID)),
  max(attribute(ID), :val(1))
].
```

Figure 6. Simplified rule for translating SQL columns to *GeRoMe* model elements

The first part of the rule refers to the elements of the SQL schema; the corresponding facts are generated by traversing the SQL schema. The second part of the rule specifies a fragment of the *GeRoMe* model: a model element with the identifier `ID`; an attribute role for this element; then, the aggregate role of the element `TableID` (specified by another rule) gets the attribute as an additional property; finally, the maximum cardinality of the attribute is set to one as attributes in the relational model are single-valued.

Our representation for the data instance level again extends ideas from ConceptBase. In ConceptBase, a base relation *proposition* expresses the relationships between objects: $P(o,x,l,y)$ states that the object x has a relationship with label l to object y . The variable o represents the object identifier of this link. We can make further statements about the type of link, e.g., whether it is an attribute link, an instantiation, or a specialization. With only a few predefined classes and only few constraints for this base relation, ConceptBase supports the representation of models as well as of its instances, but also metamodels, metametamodels, ...

In *GeRoMe*, an instance of a model is described by a set of objects; each object may have links to other objects (associations) or atomic values (attributes). In contrast to ConceptBase, *GeRoMe* thus distinguishes objects and values. Furthermore, each object has a type (an instantiation link to a model element). Syntactically, *GeRoMe* instances can be represented as a set of logical facts using a limited set of predefined predicates. This enables the use of logical languages to

express mappings on the instance level. Even if data instances have complex, nested structures, the Datalog representation just uses ‘flat’ facts.

The logical representation of *GeRoMe* uses the predicates *inst*, *attr*, *value*, and *part* to describe an instance of a model. The model is defined by a set of model elements \mathcal{M} and an instance is described by a set of objects O . Furthermore, we have a set A of atomic values (literals).

- *inst*(o,m) denotes that the object $o \in O$ is an instance of the model element $m \in \mathcal{M}$.
- *value*(o,v) denotes that the object $o \in O$ has the value $v \in A$. This is only possible, if o is an instance of a model element that plays a domain role.
- *attr*(o,a,o_v) denotes that the object $o \in O$ has an attribute of type $a \in \mathcal{M}$, and the value for this attribute is represented by the object $o_v \in O$.
- *part*(o,ae,o_p) denotes that the object $o \in O$ is an association, and the object $o_p \in O$ is a participator in this association for the association end $ae \in \mathcal{M}$.

As the predicates *value* and *attr* are often used in combination for simple typed attributes, we use the predicate *av*(o,a,v) as a shortcut to denote that the object o has a value v for the attribute a . As an example, Figure 7 shows instances for the models in Figure 5.

Concrete Instance	<i>GeRoMe</i> Instance
<i>car</i> (1,76,45)	<i>inst</i> (t_1,car), <i>av</i> ($t_1,id,1$), <i>av</i> ($t_1,speed,76$), <i>av</i> ($t_1,roadId,7$)
<i>road</i> (7,766, “Ahornstr.”)	<i>inst</i> ($t_2,road$), <i>av</i> ($t_2,name$, “Ahornstr.”), <i>av</i> ($t_2,roadId,7$), <i>av</i> ($t_2,length,766$)
<road id="7"> <car id="1" speed="76"/> </road>	<i>inst</i> ($x_r,root$), <i>inst</i> ($x_0,road$), <i>inst</i> ($x_1,roadT$), <i>part</i> ($x_0,parent,x_r$), <i>part</i> ($x_0,child,x_1$), <i>av</i> ($x_1,id,45$), <i>inst</i> (x_2,car), <i>inst</i> ($x_3,carT$), <i>part</i> ($x_2,parent,x_1$), <i>av</i> ($x_3,id,1$), <i>av</i> ($x_3,speed,76$), <i>part</i> ($x_2,child,x_3$)

Figure 7. Example instances for the models of Figure 5

The description of the relational instance in *GeRoMe* is straightforward. For each tuple, there is a corresponding object (t_1 and t_2), and each attribute value of these tuples is defined by an *av* predicate. In the XML instance, we first define an object x_r for the (invisible) document root. Then, x_0 denotes the `road` element which is an association between the document root x_r and the instance x_1 of the complex type `roadT`. Attribute values are defined in the same way as in the relational instance. The instance of the `car` element x_2 is also an association; it links the parent element x_1 with the child element x_3 , which is an instance of `carT`.

In practice, we never explicitly create instances of *GeRoMe* models in this verbose representation; it is only the formal basis for the definition of mappings. Our model-based code generation transforms data during mapping execution directly from one native representation into another according to the specified mapping.

4.4 Languages for Schema Mappings

The choice of GAV and LAV mappings mentioned in section 2 is only one aspect of schema mappings. In addition, we have tasks like explicit data exchange or even physical generation of new data.

The choice of a mapping query language depends obviously on the choice for the metamodel. In the literature, the Relational Data Model is most frequently used for data integration systems, and a relational query language chosen for representing mappings. However, the full expressive power of a relationally complete query language (such as Relational Algebra) makes reasoning over mappings undecidable. Query containment (which has often to be proven during query rewriting) is only decidable for conjunctive, i.e. Select-Project-Join (SPJ) queries [Shmu93].

In Clio, Fagin et al. [FKMP05] initially proposed to use tuple-generating dependencies (tgds) [BeVa84] for representing mappings between relational schemas, because tgds can be easily translated into executable queries.

A *source-to-target tuple-generating dependency (s-t tgd)* has the form:

$$\forall \mathbf{x} (\varphi_S(\mathbf{x}) \rightarrow \exists \mathbf{y} \psi_T(\mathbf{x}, \mathbf{y}))$$

\mathbf{x} and \mathbf{y} are sets of variables, φ_S and ψ_T are conjunctive queries over the relational schemas S and T , respectively. A mapping M is then defined as $M = \langle S, T, \Sigma \rangle$

where S and T are the mapped schemas, and Σ is a set of s-t tgds. Their simplicity combined with reasonable expressive power is a strength of s-t tgds.

Clio [HMH01] creates mappings over a nested relational model to support mappings between relational databases and XML data. Using a set of value correspondences as input, Clio is able to generate queries which transform source data into the desired target data structure. However, it would still be difficult to extend this mapping representation to express a mapping between other models, such as UML models, because there is simply no appropriate query language.

Another drawback of these *basic mappings* in Clio is pointed out in [FHH*06]: the mappings do not reflect the nested structure of the data. This leads to an inefficient execution of the mappings and redundant mapping specifications as parts of the mapping have to be repeated for different nesting levels. Furthermore, the desired grouping of the target data cannot be specified using basic mappings, as they would cause redundancy in the target. A nested mapping language introduced in [FHH*06] addresses these problems. Furthermore, they provide an algorithm to compute the nested mappings from simple morphisms which can be executed more efficiently than basic mappings.

Because executable mappings usually drive the transformation of instances of models, Melnik et al. [MBHR05] specified the semantics of model management operators by relating the instances of the operator's input and output models. They also implemented two model management system prototypes to study the specifying and manipulating of executable mappings. In the first implementation, they modified Rondo's [MRB03b] language to define path morphisms and showed that it is possible to generate executable mappings in a form of relational algebra expressions. On the positive side, this system works correctly whenever the input is specified using path morphisms, and the input is also closed under operators which return a single mapping. However, the expressiveness of path morphisms is very limited. To overcome this limitation, they developed a new prototype [MBHR05] in which mappings are specified using embedded dependencies. The expressiveness is improved, but it suffers from the problem that embedded dependencies are not closed under composition. Because of this problem, the output of the very important *Compose* operator cannot be represented as an embedded dependency; thus, a sequence of model management operators may not be executable.

4.4.1 Mapping Composition

In general, the problem of composing mappings has the following definition: Given a mapping M_{12} from model S_1 to model S_2 , and a mapping M_{23} from model S_2 to model S_3 , derive a mapping M_{13} from model S_1 to model S_3 that is equivalent to the successive application of M_{12} and M_{23} [FKPT05].

So far, mapping composition has been studied only for mappings which use the Relational Data Model as a basis [BGMN06, FKPT05, MaHa03].

In [MaHa03], the semantics of the *Compose* operator is defined relative to a class \mathcal{Q} of queries over the model S_3 . For every query $q \in \mathcal{Q}$, the certain answers for q wrt. M_{13} are the same as the certain answers for q wrt. M_{12} and M_{23} . This provided a solid basis for further research on mapping composition, but suffers from certain drawbacks caused by the fact that the semantics is defined relative to a class \mathcal{Q} of queries. Fagin et al. [FKPT05] proposed a different semantics which is defined over instance spaces of schema mappings. M_{13} is the composition of M_{12} and M_{23} if the instance space of M_{13} is the set-theoretic composition of the instance spaces of M_{12} and M_{23} . Under this semantics the mapping composition M_{13} is unique up to logical equivalence.

Another approach to define composition uses relational algebra expressions as mappings [BGMN06]. An incremental algorithm tries to replace as many symbols as possible from the “intermediate” model. Because the result of mapping composition cannot be always expressed as relational algebra expressions, the algorithm may fail under certain conditions; this limitation is in line with the results of [FKPT05].

4.4.2 Towards Composable, Executable, Generic Schema Mappings

Fagin et al. [FKPT05] proved that the language of s-t tgds is not closed under composition. To illustrate the problem, we adapt an example from [FKPT05] to our traffic scenario. Suppose we have a schema S_1 with one relation $Travels_1(C,R)$ which means that a car with the id C is traveling on a road named R . Another schema S_2 has two relations: $Travels_2(C,R)$ which is a copy of $Travels_1(C,R)$, and $Car_2(C,N)$ which states that a car C is owned by a person with name N . The mapping M_{12} between these schemas can be expressed using the following s-t tgds:

$$\forall c \forall r \text{ Travels}_1(c,r) \rightarrow \text{Travels}_2(c,r)$$

$$\forall c (\forall r \text{ Travels}_1(c,r)) \rightarrow \exists n \text{ Car}_2(c,n)$$

When we execute this mapping, transform data from S_1 to S_2 , we do not have data for n as this information is not contained in S_1 . However, we know that a car can be owned only by one person; thus, the value of n depends on the car id c . If we execute the mapping, we can thus create ‘labeled null values’ for n .

Assume there is a third schema S_3 with a single relation $\text{DrivesOn}_3(N,R)$ which holds information about persons driving on specific roads. The mapping M_{23} between S_2 and S_3 can be also expressed as a s-t tgd:

$$\forall n \forall r (\forall c \text{ Car}_2(c,n) \wedge \text{Travels}_2(c,r)) \rightarrow \text{DrivesOn}_3(n,r)$$

If we now want to compose the mappings to a mapping $M_{13} = M_{12} \circ M_{23}$, a correct formula in first-order logic would be

$$\forall c \exists n \forall r \text{ Travels}_1(c,r) \rightarrow \text{DrivesOn}_3(n,r) \quad (1)$$

However, this is not a valid s-t tgd as existential quantification is only allowed on the right-hand side of the implication. Note that the s-t tgd

$$\forall c \forall r \text{ Travels}_1(c,r) \rightarrow \exists n \text{ DrivesOn}_3(n,r)$$

is not a composition of M_{12} and M_{23} as n depends now on both, c and r , which is not correct with respect to composition semantics. To ameliorate the problem, we can skolemize formula (1) and replace n with a Skolem function $f(c)$:

$$\exists f (\forall c \forall r \text{ Travels}_1(c,r) \rightarrow \text{DrivesOn}_3(f(c),r))$$

This is a second-order formula as we quantify over the function symbol f . Fagin et al. [FKPT05] showed that second-order tgds (SO tgds) are the smallest class of formulas which can be used to represent the result of the composition of any finite s-t tgds. Thus, SO tgds are closed under composition, and mappings expressed as SO tgds can be executed in polynomial time. Therefore, SO tgds are a good formalization of mappings, albeit only for mapping relational schemas.

By extending SO tgds to *GeRoMe* we have enabled the definition of generic schema mappings [KQL*09]. We will use the running example to illustrate that SO tgds and the logical representation of *GeRoMe* instances fit nicely together and enable generic mappings which are expressive, executable, and composable.

Suppose, we want to map the relational schema R from Figure 4 to the target XML schema T , in which cars are nested into `road` elements. We can use the example instances in Figure 7 as templates for the conjunctive queries that specify the mapping:

$$\begin{aligned}
& \forall I \forall S \forall R \ (\forall T_1 \forall T_2 (\\
& \quad inst(T_1, car) \wedge av(T_1, id, I) \wedge av(T_1, speed, S) \wedge av(T_1, roadId, N) \wedge \\
& \quad inst(T_2, road) \wedge av(T_2, roadId, N)) \Rightarrow \\
& \quad \exists X_r \exists X_0 \exists X_1 \exists X_2 \exists X_3 \\
& \quad inst(X_r, root) \wedge inst(X_0, road) \wedge \\
& \quad part(X_0, parent, X_r) \wedge part(X_0, child, X_1) \wedge \\
& \quad inst(X_1, roadType) \wedge av(X_1, id, N) \wedge inst(X_2, car) \wedge \\
& \quad inst(X_3, carType) \wedge part(X_2, parent, X_1) \wedge part(X_2, child, X_3) \wedge \\
& \quad av(X_3, id, I) \wedge av(X_3, speed, S))
\end{aligned}$$

The concrete values for attributes and object identifiers have been replaced with variables. As the target schema does not contain information about length and name of roads, we do not use the corresponding predicates on the source side. This formula is a valid s-t tgd, but it does not specify how to structure the data on the target side. In this case, there is one **road** element for each tuple in the **road** relation, as **roadId** is the key of the **road** relation and all cars driving on this road should be represented as nested elements. We cannot infer this constraint from the mapping definition above, as all object variables are existentially quantified, i.e. they can have different values for each matching pair of car and road tuples on the source side.

To overcome this problem, we can skolemize the formula in a similar way as in the relational example. The variables X_i on the target side will be replaced by corresponding functions $f_i(I,S,R)$. Note that I , S , and R are the arguments of these functions as they are the universally quantified variables which appear on both sides of the implication. The resulting formula is

$$\begin{aligned}
& \exists f_r \exists f_0 \exists f_1 \exists f_2 \exists f_3 \ (\forall I \forall S \forall R \ (\forall T_1 \forall T_2 \\
& \quad inst(T_1, car) \wedge av(T_1, id, I) \wedge av(T_1, speed, S) \wedge av(T_1, roadId, N) \wedge \\
& \quad inst(T_2, road) \wedge av(T_2, roadId, N)) \Rightarrow \\
& \quad inst(f_r(I,S,R), root) \wedge inst(f_0(I,S,R), road) \wedge \\
& \quad part(f_0(I,S,R), parent, f_r(I,S,R)) \wedge part(f_0(I,S,R), child, f_1(I,S,R)) \wedge \\
& \quad inst(f_1(I,S,R), roadType) \wedge av(f_1(I,S,R), id, N) \wedge inst(f_2(I,S,R), car) \wedge \\
& \quad inst(f_3(I,S,R), carType) \wedge part(f_2(I,S,R), parent, f_1(I,S,R)) \wedge \\
& \quad part(f_2(I,S,R), child, f_3(I,S,R)) \wedge av(f_3(I,S,R), id, I) \wedge \\
& \quad av(f_3(I,S,R), speed, S))
\end{aligned}$$

This is now a valid SO tgd, but it still would not allow generating correctly structured data: All object identifiers depend on I , S , and R , such that for each combination of values for these variables, there will be a road element with a nested car element. As I is the key for cars, a road element will be generated for each car, and this is not the desired structured.

To address this problem, we need to consider the schema information in the *GeRoMe* model to figure out the correct data structuring. We will see that the constraints allow only one instance of the root element for the whole document; thus, the corresponding Skolem function f_r should have no arguments (i.e., it is a constant). Roads are identified by `roadId`; thus, the Skolem functions for the road element (f_0) and the road type (f_1) should have only R as argument. The car elements are identified by the car id I . Therefore, the Skolem functions for the instances of the car element (f_2) and the car type (f_3) have to include the variable I in their argument list. In addition, these functions need also the identifying variables of the parent element, as the car elements are nested under road elements. Consequently, the resulting functions are $f_2(I,R)$ and $f_3(I,R)$. The revised formula is given here:

$$\begin{aligned}
& \exists f_r \exists f_0 \exists f_1 \exists f_2 \exists f_3 (\forall I \forall S \forall R (\forall T_1 \forall T_2 \\
& \quad inst(T_1, car) \wedge av(T_1, id, I) \wedge av(T_1, speed, S) \wedge av(T_1, roadId, N) \wedge \\
& \quad inst(T_2, road) \wedge av(T_2, roadId, N)) \Rightarrow \\
& \quad inst(f_r(), root) \wedge inst(f_0(R), road) \wedge \\
& \quad part(f_0(R), parent, f_r()) \wedge part(f_0(R), child, f_1(R)) \wedge \\
& \quad inst(f_1(R), roadType) \wedge av(f_1(R), id, N) \wedge inst(f_2(I, R), car) \wedge \\
& \quad inst(f_3(I, R), carType) \wedge part(f_2(I, R), parent, f_1(I, R)) \wedge \\
& \quad part(f_2(I, R), child, f_3(I, R)) \wedge av(f_3(I, R), id, I) \wedge \\
& \quad av(f_3(I, R), speed, S)))
\end{aligned}$$

By using the generic mapping language and Skolem functions, we can define mappings between arbitrarily structured models in various modeling languages. To execute the mappings, we do not have to provide an interpretation of the Skolem functions. The Skolem functions will be used by the mapping compiler to generate appropriate queries or update statements which take the defined structure into account.

5 Schema Matching and Merging

The formal definition of GeRoMe and its mapping language enable the formal definition and verification of model management operators. We discuss briefly the realization of the operators *Match* and *Merge*. Schema merging and schema matching are two related operators, in that the output of matching can be used as input for schema merging.

5.1 Schema Matching

Schema matching is the task of identifying a set of correspondences (also called a morphism) between schema elements. Many aspects have to be considered during the process of matching, such as data values, element names, constraint information, structure information, domain knowledge, cardinality relationships, and so on. All this information is useful in understanding the semantics of a schema, but it can be a very time consuming problem to collect and apply this information.

Therefore, automatic methods are required for schema matching [RaBe01, ShEu05]. *Element-Level Matchers* separately take the information of each schema element into account, using linguistic information (name of the element) or constraint information (data type, key constraints). *Structure-Level Matchers* use graph matching to measure the similarity of the structures implied the schema. *Instance-Level Matchers* employ data instances to match schema elements: If the instance sets of two elements are similar, or have a similar value distribution, this might indicate a similarity of the schema elements. *Machine-Learning Matchers* use either instance data or previously identified matches as training data to detect similar matches in new schema matching problems.

It is widely agreed, that no single method can solve the schema matching problem. Therefore, matching frameworks such as COMA++ [ADMR05], Protoplasm [BMPQ04], or YAM [DCBM09] combine multiple individual matching methods to achieve a better result. For the heterogenous case, we have developed an extensible and flexible matching framework in our model management system *GeRoMeSuite* [KQLL07], which is also able to combine several matching methods in entirely configurable matching strategies. We could show that our system in particular good for heterogeneous matching tasks, e.g., matching of XML schemas and OWL ontologies [QKL07]. The field of ontology

matching had more attention in the recent years than schema matching because of the structured evaluation of ontology matching systems in the Ontology Alignment Evaluation Initiative (OAEI, [EMS*11]).

Such logical methods include, for example, the validation of correspondences. A computed alignment should be consistent with the information present in the matched ontologies (or schemas). If it is possible to derive an inconsistency from a correspondence, the identified correspondences may be wrong. For example, in ASMOV [JMSK09], validation of the computed alignment is a key concept in their ontology matching system and greatly improves the quality of the match result.

Our matching framework also incorporates validation methods for schema matching. We could show that such methods also improve the quality of matching for our generic approach (e.g., by comparing the results in [QGK*08] and [QGK*09]). This also demonstrates another advantage of the rich generic metamodel: due to the exact representation of models in *GeRoMe*, we are able to apply these logical methods not only to ontologies, but to other modeling languages as well (e.g., by exploiting inheritance relationships or foreign key constraints).

Although the field of schema and ontology matching has made significant improvements in the recent years, there are still a lot of challenges to be addressed. For example, [ShEu12] mention efficiency, matching with background knowledge, matcher selection and tuning as important requirements for ontology matching systems. We have developed a method for automatically retrieving background knowledge in form of ontologies from the web [QRK11]. For a matching task, we inspect the source and target models and extract a few descriptive keywords from the models, in order to characterize the domain of the models to be matched. Using these keywords, we employ traditional search engines or specific ontology search engines (e.g., Swoogle [DFJ*04] or Watson [dAMo11]) to find ontologies on the web. In addition to matching the source and target models directly, we thus match the input models also with the background ontology. For example, the ontology O in Figure 4 can be seen as an ontology bridging the semantic gap between the input models S and T . By composing the computed matches, we can then infer more matches between the models S and T and thereby get a more accurate mapping. In an extensive evaluation of the

approach [QRK11], we could show that the ‘noise’ which might be introduced by inappropriate background knowledge is low and that it is more useful to use several ontologies as background knowledge instead of just one.

5.2 Schema Merging

The challenging part in schema merging is a formal definition of the desired outcome. How can we characterize the requirements for the integrated schema formally? How can we prove that a schema integration method actually produced the correct result?

In [PoBe03], requirements for a *Merge* operator include the *preservation of the original semantics* of the input schemas (e.g., elements, relationships, and equalities) and the *minimality of the merged schema* (i.e., no extra information should be added). We present a method based on the previous formalisms that achieves these goals under certain assumptions.

In the previous sections, we have addressed extensional mappings between instances of schemas for tasks such as data translation or query rewriting. The semantics of extensional mappings (and their composition) is defined with respect to the instances of schemas.

For schema merging, we need *intensional mappings*. For example, consider relational databases from two different cities A and B which maintain information about roads, e.g., both databases have one relation of the form $road(id, length, name)$. We could state the relationship between the two databases using the following mapping expression:

$$\forall I \forall L \forall N \quad road_A(I,L,N) \rightarrow road_B(I,L,N).$$

This mapping is not correct if we consider the extensions of the relations. A road in city A is not a road in city B and vice versa. Nevertheless, at the intensional level, the mapping is useful because it states that the two relations have the same intended semantics. If additionally both databases would contain information for the same domain (i.e., for the same city), then the mapping would be correct. This difference in the semantics of mappings has been characterized in [CGL*98].

Schema merging is about integrating models according to their intensional semantics. It has the goal to construct a “duplicate-free” union of the input models and mapping, with respect to the real world concepts described by the model

elements. The integrated model should describe each real world concept only once. Thus, we need intensional mappings for schema integration.

In our schema merging approach [LQKG10, LiQu11], tuple-generating dependencies (tgds) are used to express the constraints of the input schemas as well as the inter-schema constraints (Figure 8).

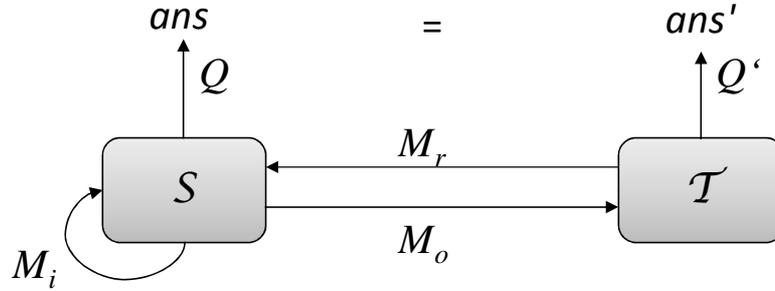


Figure 8. The basic idea for Schema Merge using logical mappings

Assume we want to merge a set of schemas S_1, \dots, S_n . We first construct S as the duplicate-free union of S_1, \dots, S_n (if there are elements with identical names, they need to be renamed). The mappings between the input schemas as well as the constraints are defined as tgds in the input mapping M_i . When we create the integrated schema T , we also produce two mappings: M_o is the output mapping responsible for translating data from the sources to the new target schema; M_r is a recovery mapping (or witness mapping) which is basically an inversion of M_o . The integrated schema T is created by a step-by-step procedure in which we remove incrementally elements from the input schema schema and check whether the obtained schema still fulfills the requirements of the integrated schema.

In our approach, the main requirement is that queries over the integrated schema should have the same result as over the source schemas. We can prove this property by using the generated mappings M_o and M_r . If we can prove that the answers for queries are the same when we evaluate them directly over S and when we evaluate them over S through the mapping M_o and M_r , then we know that no information has been lost in the integrated schema T .

This means that we must be able to prove query equivalence for all queries in a specific query language over the integrated schema and the source schemas. This is only possible if we have a formal mapping from the data sources to the integrated schema and vice versa. As query equivalence is decidable for

conjunctive queries, we chose the class of conjunctive queries over relational schemas as our query language.

The target schema T is minimized by identifying elements which are potentially redundant. Redundant attributes are detected by reasoning over the input mappings and constraints. Our implementation of this merge method [LQK*11] also considers the case of collapsing relations, i.e., we remove also redundant relations, by discovering bi-directional inclusion dependencies [LiQu11].

6 Conclusion and Outlook

In this paper, we have reviewed the evolution of data-centric approaches that emphasize the semantics of information integration. We have shown that this approach has enabled significant progress in automation of many data and model management tasks, especially in the context of the relational data model and its extensions, e.g. to nested relations. With our *GeRoMe* approach, we have also illustrated current research on how to extend these approaches to the case of heterogeneous data models, without losing again the advantages of clear semantics and highly automated tools.

We intentionally limited our discussion to the case of structured or semi-structured data with a schema. The linkage of these approaches to text and media information integration which are also subject to intense research in the last years, still remains to be explored in depth.

But even within schema-based approaches, there is still a long way to go, in order to truly conquer the challenge of heterogeneity. While our algorithms in *GeRoMeSuite* address surprisingly well the aspects of (automatically) executable mapping in a significantly richer model language context than earlier solutions, and provide decent assistance for heterogeneous schema matching, our present solution for merging – information-preserving schema integration including target schema minimization – is unfortunately only provably correct and complete for the case of relational model integration while a tractable case where this also holds for heterogenous merging remains an open problem.

7 References

- [AAB*05] Abiteboul, S., Agrawal, R., Bernstein, P., Carey, M., Ceri, S., Croft, B., DeWitt, D., Franklin, M., Molina, H.G., Gawlick, D., Gray, J., Haas, L., Halevy, A., Hellerstein, J., Ioannidis, Y., Kersten, M., Pazzani, M., Lesk, M., Maier, D., Naughton, J., Schek, H., Sellis, T., Silberschatz, A., Stonebraker, M., Snodgrass, R., Ullman, J., Weikum, G., Widom, J., Zdonik, S.: The Lowell database research self-assessment. *Commun. ACM* 48(5), 111–118 (2005). DOI <http://doi.acm.org/10.1145/1060710.1060718>
- [AHV95] Abiteboul, S., Hull, R., Vianu, V.: *Foundations of Databases*. Addison-Wesley (1995)
- [AAB*08] Agrawal, R., Ailamaki, A., Bernstein, P.A., Brewer, E.A., Carey, M.J., Chaudhuri, S., Doan, A., Florescu, D., Franklin, M.J., Garcia-Molina, H., Gehrke, J., Gruenwald, L., Haas, L.M., Halevy, A.Y., Hellerstein, J.M., Ioannidis, Y.E., Korth, H.F., Kossmann, D., Madden, S., Magoulas, R., Ooi, B.C., O'Reilly, T., Ramakrishnan, R., Sarawagi, S., Stonebraker, M., Szalay, A.S., Weikum, G.: The Claremont report on database research. *SIGMOD Rec.* 37(3), 9–19 (2008). DOI <http://doi.acm.org/10.1145/1462571.1462573>
- [ABLM10] Arenas, M., Barceló, P., Libkin, L., Murlak, F.: *Relational and XML Data Exchange. Synthesis Lectures on Data Management*. Morgan & Claypool Publishers (2010)
- [AHK96] Arens, Y., Hsu, C.N., Knoblock, C.A.: Query Processing in the Sims Information Mediator. In: *Advanced Planning Technology*, pp. 61–69 (1996)
- [ABB*09] Atzeni, P., Bellomarini, L., Bugiotti, F., Gianforme, G.: Mism: A platform for model-independent solutions to model management problems. *Journal of Data Semantics* 14, 133–161 (2009)
- [ACB05] Atzeni, P., Cappellari, P., Bernstein, P.A.: A multilevel dictionary for model management. In: L.M.L. Delcambre, C. Kop, H.C. Mayr, J. Mylopoulos, O. Pastor (eds.) *Proc. 24th International Conference on Conceptual Modeling (ER)*, Lecture Notes in Computer Science, vol. 3716, pp. 160–175. Springer, Klagenfurt, Austria (2005)
- [ACB06] Atzeni, P., Cappellari, P., Bernstein, P.A.: Model-independent schema and data translation. In: Y.E. Ioannidis, M.H. Scholl, J.W. Schmidt, F. Matthes, M. Hatzopoulos, K. Böhm, A. Kemper, T. Grust, C. Böhm (eds.) *Proc. 10th International Conference on Extending Database Technology (EDBT)*, Lecture Notes in Computer Science, vol. 3896, pp. 368–385. Springer, Munich, Germany (2006)

- [ACT*08] Atzeni, P., Cappellari, P., Torlone, R., Bernstein, P.A., Gianforme, G.: Model-independent schema translation. *VLDB Journal* 17(6), 1347–1370 (2008)
- [AtTo96] Atzeni, P., Torlone, R.: Management of multiple models in an extensible database design tool. In: P.M.G. Apers, M. Bouzeghoub, G. Gardarin (eds.) *Proc. 5th International Conference on Extending Database Technology (EDBT)*, Lecture Notes in Computer Science, vol. 1057, pp. 79–95. Springer, Avignon, France (1996)
- [ADMR05] Aumüller, D., Do, H.H., Massmann, S., Rahm, E.: Schema and ontology matching with COMA++. In: F. Özcan (ed.) *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pp. 906–908. ACM, Baltimore, Maryland, USA (2005)
- [BaDa77] Bachman, C.W., Daya, M.: The role concept in data models. In: *Proceedings of the Third International Conference on Very Large Data Bases (VLDB)*, pp. 464–476. IEEE-CS and ACM, Tokyo, Japan (1977)
- [BaJa99] Baumeister, M., Jarke, M.: Compaction of Large Class Hierarchies in Databases for Chemical Engineering. *Proceedings 8. GI-Fachtagung für Datenbanksysteme in Büro, Technik und Wissenschaft (BTW)*, Freiburg, Springer, pp. 343–361 (1999)
- [BLN86] Batini, C., Lenzerini, M., Navathe, S.B.: A comparative analysis of methodologies for database schema integration. *ACM Computing Surveys* 18(4), 323–364 (1986)
- [BeVa84] Beeri, C., Vardi, M.Y.: A proof procedure for data dependencies. *Journal of the ACM* 31(4), 718–741 (1984)
- [BCVB01] Bergamaschi, S., Castano, S., Vincini, M., Beneventano, D.: Semantic integration of heterogeneous information sources. *Data & Knowledge Engineering* 36(3), 215–249 (2001)
- [BBC*98] Bernstein, P., Brodie, M., Ceri, S., DeWitt, D., Franklin, M., Garcia-Molina, H., Gray, J., Held, J., Hellerstein, J., Jagadish, H.V., Lesk, M., Maier, D., Naughton, J., Pirahesh, H., Stonebraker, M., Ullman, J.: The Asilomar report on database research. *SIGMOD Rec.* 27(4), 74–80 (1998). DOI <http://doi.acm.org/10.1145/306101.306137>
- [BDD*89] Bernstein, P.A., Dayal, U., DeWitt, D.J., Gawlick, D., Gray, J., Jarke, M., Lindsay, B.G., Lockemann, P.C., Maier, D., Neuhold, E.J., Reuter, A., Rowe, L.A., Schek, H.J., Schmidt, J.W., Schrefl, M., Stonebraker, M.: Future directions in dbms research -the laguna beach participants. *SIGMOD Rec.* 18(1), 17–26 (1989). DOI <http://doi.acm.org/10.1145/382272.1367994>
- [BGMN06] Bernstein, P.A., Green, T.J., Melnik, S., Nash, A.: Implementing mapping composition. In: U. Dayal, K.Y. Whang, D.B. Lomet, G. Alonso, G.M. Lohman, M.L. Kersten, S.K. Cha, Y.K. Kim (eds.) *Proc. 32nd Intl. Conference on Very Large Data Bases (VLDB)*, pp. 55–66. ACM Press (2006)

- [BeHa08] Bernstein, P.A., Haas, L.M.: Information integration in the enterprise. *Commun. ACM* 51(9), 72–79 (2008)
- [BHP00] Bernstein, P.A., Halevy, A.Y., Pottinger, R.: A vision for management of complex models. *SIGMOD Record* 29(4), 55–63 (2000)
- [BeMe07] Bernstein, P.A., Melnik, S.: Model management 2.0: Manipulating richer mappings. In: L. Zhou, T.W. Ling, B.C. Ooi (eds.) *Proc. ACM SIGMOD Intl. Conf. on Management of Data*, pp. 1–12. ACM Press, Beijing, China (2007). DOI <http://doi.acm.org/10.1145/1247480.1247482>
- [BMPQ04] Bernstein, P.A., Melnik, S., Petropoulos, M., Quix, C.: Industrialstrength schema matching. *SIGMOD Record* 33(4), 38–43 (2004)
- [BiCo86] Biskup, J., Convent, B.: A formal view integration method. In: C. Zaniolo (ed.) *Proc. ACM SIGMOD Intl. Conf. on Management of Data*, pp. 398–407. ACM Press, Washington, D.C. (1986)
- [BMM*08] Brandt, S.C., Morbach, J., Miatidis, M., Theißen, M., Jarke, M., Marquardt, W.: An ontology-based approach to knowledge management in design processes. *Computers & Chemical Engineering* 32(1-2), 320-342 (2008)
- [Brod10] Brodie, M.L.: Data integration at scale: From relational data integration to information ecosystems. In: *Proc. 24th IEEE Intl. Conf. on Advanced Information Networking and Applications (AINA)*, pp. 2–3. IEEE Computer Society, Perth, Australia (2010)
- [CCGL04] Calì, A., Calvanese, D., Giacomo, G.D., Lenzerini, M.: Data integration under integrity constraints. *Information Systems* 29(2), 147–163 (2004). DOI [http://dx.doi.org/10.1016/S0306-4379\(03\)00050-4](http://dx.doi.org/10.1016/S0306-4379(03)00050-4)
- [CGL*98] Calvanese, D., Giacomo, G.D., Lenzerini, M., Nardi, D., Rosati, R.: Description logic framework for information integration. In: A.G. Cohn, L.K. Schubert, S.C. Shapiro (eds.) *Proceedings of the Sixth International Conference on Principles of Knowledge Representation and Reasoning (KR'98)*, pp. 2–13. Morgan Kaufmann, Trento, Italy (1998)
- [CGL*01] Calvanese, D., Giacomo, G.D., Lenzerini, M., Nardi, D., Rosati, R.: Data Integration in Data Warehousing. *International Journal of Cooperative Information Systems (IJCIS)* 10(3), 237–271 (2001)
- [CaVi83] Casanova, M.A., Vidal, V.M.P.: Towards a sound view integration methodology. In: *Proc. 2nd ACM Symposium on Principles of Database Systems (PODS)*, pp. 36–47. ACM, Atlanta, GA (1983)
- [Catt10] Cattell, R.: Scalable SQL and NoSQL data stores. *SIGMOD Record* 39(4), 12–27 (2010)
- [CGT90] Ceri, S., Gottlob, G., Tanca, L.: *Logic Programming and Databases*. Springer (1990)

- [CePe84] Ceri, S., Pelagatti, G.: Distributed Databases: Principles and Systems. McGraw-Hill Book Company (1984)
- [CHS91] Collet, C., Huhns, M.N., Shen, W.M.: Resource integration using a large knowledge base in carnot. *IEEE Computer* 24(12), 55–62 (1991)
- [dAMo11] d'Aquin, M., Motta, E.: Watson, more than a semantic web search engine. *Semantic Web Journal* 2(1), 55–63 (2011). <http://www.semantic-web-journal.net/content/new-submission-watson-more-semantic-web-search-engine>
- [DFJ*04] Ding, L., Finin, T., Joshi, A., Pan, R., Cost, R., Peng, Y., Reddivari, P., Doshi, V., Sachs, J.: Swoogle: a search and metadata engine for the semantic web. In: Proc. CIKM (2004)
- [DoRa02] Do, H.H., Rahm, E.: Coma -a system for flexible combination of schema matching approaches. In: Proc. 28th Intl. Conference on Very Large Data Bases (VLDB), pp. 610–621. Morgan Kaufmann, Hong Kong, China (2002)
- [Dolk88] Dolk, D.R.: Model management and structured modeling: The role of an information resource dictionary system. *Communications of the ACM* 31(6) (1988)
- [DCBM09] Duchateau, F., Coletta, R., Bellahsene, Z., Miller, R.J.: (Not) yet another matcher. In: D.W.L. Cheung, I.Y. Song, W.W. Chu, X. Hu, J.J. Lin (eds.) Proc. 18th ACM Conference on Information and Knowledge Management (CIKM), pp. 1537–1540. ACM, Hong Kong, China (2009)
- [EMS*11] Euzenat, J., Meilicke, C., Stuckenschmidt, H., Shvaiko, P., dos Santos, C.T.: Ontology alignment evaluation initiative: Six years of experience. *Journal on Data Semantics* 15, 158–192 (2011)
- [FHH*09] Fagin, R., Haas, L.M., Hernández, M.A., Miller, R.J., Popa, L., Velegrakis, Y.: Clio: Schema mapping creation and data exchange. In: A. Borgida, V.K. Chaudhri, P. Giorgini, E.S.K. Yu (eds.) *Conceptual Modeling: Foundations and Applications*, Lecture Notes in Computer Science, vol. 5600, pp. 198–236. Springer (2009)
- [FKMP05] Fagin, R., Kolaitis, P., Miller, R.J., Popa, L.: Data exchange: Semantics and query answering. *Theoretical Computer Science* 336, 89–124 (2005)
- [FKPT05] Fagin, R., Kolaitis, P.G., Popa, L., Tan, W.C.: Composing schema mappings: Second-order dependencies to the rescue. *ACM Trans. Database Syst.* 30(4), 994–1055 (2005)
- [FHH*06] Fuxman, A., Hernández, M.A., Ho, C.T.H., Miller, R.J., Papotti, P., Popa, L.: Nested mappings: Schema mapping reloaded. In: U. Dayal, K.Y. Whang, D.B. Lomet, G. Alonso, G.M. Lohman, M.L. Kersten, S.K. Cha, Y.K. Kim (eds.) Proc. 32nd Intl. Conference on Very Large Data Bases (VLDB), pp. 67–78. ACM Press (2006)

- [GPQ*97] Garcia-Molina, H., Papakonstantinou, Y., Quass, D., Rajaraman, A., Sagiv, Y., Ullman, J.D., Vassalos, V., Widom, J.: The tsimmis approach to mediation: Data models and languages. *Journal of Intelligent Information Systems* 8(2), 117–132 (1997)
- [GQSJ12] Geisler, S., Quix, C., Schiffer, S., Jarke, M.: An evaluation framework for traffic information systems based on data streams. *Transportation Research Part C* 23, 29–55 (2012)
- [GKD97] Genesereth, M.R., Keller, A.M., Duschka, O.M.: Infomaster: An information integration system. In: J. Peckham (ed.) *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pp. 539–542. ACM Press, Tucson, Arizona (1997)
- [Haas07] Haas, L.M.: Beauty and the beast: The theory and practice of information integration. In: T. Schwentick, D. Suciu (eds.) *ICDT, Lecture Notes in Computer Science*, vol. 4353, pp. 28–43. Springer, Barcelona, Spain (2007)
- [HHH*05] Haas, L.M., Hernández, M.A., Ho, H., Popa, L., Roth, M.: Clio grows up: from research prototype to industrial tool. In: F. Özcan (ed.) *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pp. 805–810. ACM, Baltimore, Maryland, USA (2005)
- [Hale01] Halevy, A.Y.: Answering queries using views: A survey. *VLDB Journal* 10(4), 270–294 (2001)
- [HIM*04] Halevy, A.Y., Ives, Z.G., Madhavan, J., Mork, P., Suciu, D., Tatarinov, I.: The piazza peer data management system. *IEEE Transactions on Knowledge and Data Engineering* 16(7), 787–798 (2004). DOI <http://doi.ieeecomputersociety.org/10.1109/TKDE.2004.1318562>
- [HaK110] Haslhofer, B., Klas, W.: A survey of techniques for achieving metadata interoperability. *ACM Comput. Surv.* 42(2) (2010)
- [HMH01] Hernández, M.A., Miller, R.J., Haas, L.M.: Clio: A semi-automatic tool for schema mapping. In: *Proc. ACM SIGMOD Intl. Conference on the Management of Data*, p. 607. ACM Press, Santa Barbara, CA (2001)
- [IRDS90] ISO/IEC: Information technology – Information Resource Dictionary System (IRDS) framework. International Standard ISO/IEC 10027:1990, ISO International Organization for Standardization, (1990)
- [MOF05] ISO/IEC: Information technology -Meta Object Facility (MOF). International Standard ISO/IEC 19502:2005, ISO International Organization for Standardization (2005)

- [JGJS95] Jarke, M., Gallersdörfer, R., Jeusfeld, M.A., Staudt, M.: ConceptBase -a deductive object base for meta data management. *Journal of Intelligent Information Systems* 4(2), 167–192 (1995)
- [JJN*09] Jarke, M., Jeusfeld, M., Nissen, H., Quix, C., Staudt, M.: Metamodelling with datalog and classes: Conceptbase at the age of 21. In: *Proc. 2nd Intl. Conf. Object Databases (ICOODB 09)*, pp. 95–112. Springer-Verlag (2009)
- [JJQV99] Jarke, M., Jeusfeld, M.A., Quix, C., Vassiliadis, P.: Architecture and Quality in Data Warehouses: An Extended Repository Approach. *Information Systems* 24(3), 229–253 (1999)
- [JLVV03] Jarke, M., Lenzerini, M., Vassiliou, Y., Vassiliadis, P. (eds.): *Fundamentals of Data Warehouses*, 2 edn. Springer-Verlag (2003)
- [JMSK09] Jean-Mary, Y.R., Shironoshita, E.P., Kabuka, M.R.: Ontology matching with semantic verification. *Journal of Web Semantics* 7(3), 235–251 (2009)
- [Jeus92] Jeusfeld, M.A.: *Änderungskontrolle in deduktiven Objektbanken*. PhD thesis, Universität Passau (1992)
- [JeJo95] Jeusfeld, M.A., Johnen, U.A.: An executable meta model for reengineering of database schemas. *Intl. Journal of Cooperative Information Systems* 4(2-3), 237–258 (1995)
- [KeQu07] Kensche, D., Quix, C.: Transformation of models in(to) a generic metamodel. In: *Proc. BTW Workshop on Model and Metadata Management*, pp. 4–15 (2007)
- [KQCJ07] Kensche, D., Quix, C., Chatti, M.A., Jarke, M.: GeRoMe: A generic role based metamodel for model management. *Journal on Data Semantics VIII*, 82–117 (2007)
- [KQLL07] Kensche, D., Quix, C., Li, X., Li, Y.: GeRoMeSuite: A system for holistic generic model management. In: C. Koch, J. Gehrke, M.N. Garofalakis, D. Srivastava, K. Aberer, A. Deshpande, D. Florescu, C.Y. Chan, V. Ganti, C.C. Kanne, W. Klas, E.J. Neuhold (eds.) *Proceedings 33rd Intl. Conf. on Very Large Data Bases (VLDB)*, pp. 1322–1325. Vienna, Austria (2007)
- [KQL*09] Kensche, D., Quix, C., Li, X., Li, Y., Jarke, M.: Generic schema mappings for composition and query answering. *Data Knowl. Eng.* 68(7), 599–621 (2009). DOI <http://dx.doi.org/10.1016/j.datak.2009.02.006>
- [KDH*05] Kerner, B., Demir, C., Herrtwich, R., Klenov, S., Rehborn, H., Aleksic, M., Haug, A.: Traffic state detection with floating car data in road networks. In: *Proceedings of the 8th International IEEE Conference on Intelligent Transportation Systems*, pp. 700–705. Daimler Chrysler AG (2005)
- [KLSS95] Kirk, T., Levy, A.Y., Sagiv, Y., Srivastava, D.: The Information Manifold. In: *Proceedings of the AAAI 1995 Spring Symposium on Information Gathering from Heterogeneous, Distributed Environments*, pp. 85–91 (1995)

- [Lenz02] Lenzerini, M.: Data integration: A theoretical perspective. In: L. Popa (ed.) Proc. 21st ACM Symposium on Principles of Database Systems (PODS), pp. 233–246. ACM Press, Madison, Wisconsin (2002). DOI <http://doi.acm.org/10.1145/543613.543644>
- [LiQu11] Li, X., Quix, C.: Merging relational views: A minimization approach. In: M.A. Jeusfeld, L.M.L. Delcambre, T.W. Ling (eds.) Proc. 30th Intl. Conference on Conceptual Modeling (ER 2011), Lecture Notes in Computer Science, vol. 6998, pp. 379–392. Springer, Brussels, Belgium (2011)
- [LQKG10] Li, X., Quix, C., Kensche, D., Geisler, S.: Automatic schema merging using mapping constraints among incomplete sources. In: J. Huang, N. Koudas, G.J.F. Jones, X. Wu, K. Collins-Thompson, A. An (eds.) Proc. 19th ACM Conf. on Information and Knowledge Management (CIKM), pp. 299–308. ACM, Toronto, Ontario, Canada (2010)
- [LQK*11] Li, X., Quix, C., Kensche, D., Geisler, S., Guo, L.: Automatic generation of mediated schemas through reasoning over data dependencies. In: S. Abiteboul, K. Böhm, C. Koch, K.L. Tan (eds.) Proc. 27th Intl. Conf. on Data Engineering (ICDE), pp. 1280–1283. IEEE Computer Society, Hannover, Germany (2011)
- [MaHa03] Madhavan, J., Halevy, A.Y.: Composing mappings among data sources. In: J.C. Freytag, P.C. Lockemann, S. Abiteboul, M.J. Carey, P.G. Selinger, A. Heuer (eds.) Proc. of 29th Intl. Conference on Very Large Data Bases (VLDB), pp. 572–583. Morgan Kaufmann, Berlin, Germany (2003)
- [MBHR05] Melnik, S., Bernstein, P.A., Halevy, A.Y., Rahm, E.: Supporting executable mappings in model management. In: F. Özcan (ed.) Proceedings of the ACM SIGMOD International Conference on Management of Data, pp. 167–178. ACM, Baltimore, Maryland, USA (2005)
- [MRB03a] Melnik, S., Rahm, E., Bernstein, P.A.: Developing metadata-intensive applications with Rondo. *Journal of Web Semantics* 1(1), 47–74 (2003)
- [MRB03b] Melnik, S., Rahm, E., Bernstein, P.A.: Rondo: A programming platform for generic model management. In: Proc. ACM SIGMOD Intl. Conference on Management of Data, pp. 193–204. ACM, San Diego, CA (2003)
- [MHH00] Miller, R.J., Haas, L.M., Hernández, M.A.: Schema mapping as query discovery. In: A.E. Abbadi, M.L. Brodie, S. Chakravarthy, U. Dayal, N. Kamel, G. Schlageter, K.Y. Whang (eds.) Proc. 26th Intl. Conference on Very Large Data Bases (VLDB), pp. 77–88. Morgan Kaufmann, Cairo, Egypt (2000)
- [MBM07] Mork, P., Bernstein, P.A., Melnik, S.: Teaching a schema translator to produce o/r views. In: Proc. 26th Intl. Conf. on Conceptual Modeling (ER'07), LNCS, vol. 4801, pp. 102–119. Springer (2007)

- [MJBK90] Mylopoulos, J., Borgida, A., Jarke, M., Koubarakis, M.: Telos: Representing Knowledge About Information Systems. *ACM Transactions on Information Systems* 8(4), 325–362 (1990)
- [NaBr03] Nardi, D., Brachman, R.J.: An introduction to description logics. In: F. Baader, D. Calvanese, D.L. McGuinness, D. Nardi, P.F. Patel-Schneider (eds.) *Description Logic Handbook*. Cambridge University Press (2003)
- [NiJa99] Nissen, H.W., Jarke, M.: Repository Support for Multi-Perspective Requirements Engineering. *Information Systems* 24(2), 131-158 (1999)
- [PaSp98] Parent, C., Spaccapietra, S.: Issues and approaches of database integration. *Communications of the ACM* 41(5), 166–178 (1998)
- [PoBe03] Pottinger, R., Bernstein, P.A.: Merging models based on given correspondences. In: J.C. Freytag, P.C. Lockemann, S. Abiteboul, M.J. Carey, P.G. Selinger, A. Heuer (eds.) *Proc. of 29th Intl. Conference on Very Large Data Bases (VLDB)*, pp. 826–873. Morgan Kaufmann, Berlin, Germany (2003)
- [PoHa01] Pottinger, R., Halevy, A.Y.: Minicon: A scalable algorithm for answering queries using views. *VLDB Journal* 10(2-3), 182–198 (2001)
- [Quix09a] Quix, C.: Meta data repository. In: L. Liu, M.T. Ozsu (eds.) *Encyclopedia of Database Systems*, pp. 1718–1722. Springer (2009)
- [Quix09b] Quix, C.: Model management. In: L. Liu, M.T. Ozsu (eds.) *Encyclopedia of Database Systems*, pp. 1760–1764. Springer (2009)
- [QGK*08] Quix, C., Geisler, S., Kensche, D., Li, X.: Results of GeRoMesuite for OAEI 2008. In: *Proc. 3rd Intl. Workshop On Ontology Matching (OM2008)* (2008). URL <http://data.semanticweb.org/workshop/om/2008/paper/main/13>
- [QGK*09] Quix, C., Geisler, S., Kensche, D., Li, X.: Results of geromesuite for oaei 2009. In: P. Shvaiko, J. Euzenat, F. Giunchiglia, H. Stuckenschmidt, N.F. Noy, A. Rosenthal (eds.) *Proc. 4th Intl. Workshop on Ontology Matching, CEUR Workshop Proceedings*, vol. 551. CEUR-WS.org, Chantilly, USA (2009).
- [QKL07] Quix, C., Kensche, D., Li, X.: Matching of ontologies with xml schemas using a generic metamodel. In: R. Meersman, Z. Tari (eds.) *Proc. OTM Confederated International Conf. CoopIS/DOA/ODBASE/GADA/IS, Lecture Notes in Computer Science*, vol. 4803, pp. 1081–1098. Springer, Vilamoura, Portugal (2007)
- [QRK11] Quix, C., Roy, P., Kensche, D.: Automatic selection of background knowledge for ontology matching. In: *Proc. Intl. Workshop on Semantic Web Information Management (SWIM)*, pp. 5:1–5:7. ACM, New York, NY, USA (2011).
- [RaBe01] Rahm, E., Bernstein, P.A.: A survey of approaches to automatic schema matching. *VLDB Journal* 10(4), 334–350 (2001)

- [RaJa01] Ramesh, B., Jarke, M.: Toward Reference Models of Requirements Traceability. *IEEE Transactions on Software Engineering* 27(1), 58-93 (2001)
- [RiSc91] Richardson, J., Schwarz, P.: Aspects: extending objects to support multiple, independent roles. In: *Proc. ACM SIGMOD Intl. Conference on Management of Data*, pp. 298–307. Denver, CO (1991). DOI
- [Shmu93] Shmueli, O.: Equivalence of datalog queries is undecidable. *Journal of Logic Programming* 15(3), 231–241 (1993)
- [SHT*77] Shu, N.C., Housel, B.C., Taylor, R.W., Ghosh, S.P., Lum, V.Y.: EXPRESS: A Data EXtraction, Processing, and REStructuring System. *ACM Trans. Database Syst.* 2(2), 134–174 (1977)
- [ShEu05] Shvaiko, P., Euzenat, J.: A survey of schema-based matching approaches. *Journal on Data Semantics IV*, 146–171 (2005). LNCS 3730
- [ShEu12] Shvaiko, P., Euzenat, J.: Ontology matching: State of the art and future challenges. *IEEE Transactions on Knowledge and Data Engineering* (2012). To appear, preprint available at http://www.dit.unitn.it/~p2p/RelatedWork/Matching/SurveyOMtkde_SE.pdf
- [SSU90] Silberschatz, A., Stonebraker, M., Ullman, J.D.: Database systems: Achievements and opportunities -the "lagunita" report. *SIGMOD Record* 19(4), 6–22 (1990)
- [SSU96] Silberschatz, A., Stonebraker, M., Ullman, J.D.: Database research: Achievements and opportunities into the 21st century. Tech. rep., Stanford University, Stanford, CA, USA (1996)
- [SCJ*97] Singh, M.P., Cannata, P.E., Jacobs, N., Ksiezzyk, T., Ong, K., Sheth, A.P., Tomlinson, C., , Woelk, D.: The carnot heterogeneous database project: Implemented applications. *Distributed and Parallel Databases* 5(2), 207–225 (1997)
- [Smit07] Smith, M.: Toward enterprise information integration. *Software Magazine* (2007). URL <http://www.softwaremag.com/content/ContentCT.asp? P=3034>
- [SpPa94] Spaccapietra, S., Parent, C.: View integration: A step forward in solving structural conflicts. *IEEE Transactions on Knowledge and Data Engineering* 6(2), 258–274 (1994)
- [StJa00] Staudt, M., Jarke, M.: View Management Support in Advanced Knowledge Base Servers. *Journal Intelligent Information Systems* 15(3), 253-285 (2000)
- [Ston10] Stonebraker, M.: SQL databases v. NoSQL databases. *Commun. ACM* 53(4), 10–11 (2010)
- [SBH*10] Stubing, H., Bechler, M., Heussner, D., May, T., Radusch, I., Rechner, H., Vogel, P.: simtd: A car-to-x system architecture for field operational tests. *IEEE Communications Magazine* 48(5), 148–154 (2010).

- [Voge07] Vogels, W.: Data access patterns in the amazon.com technology platform. In: C. Koch, J. Gehrke, M.N. Garofalakis, D. Srivastava, K. Aberer, A. Deshpande, D. Florescu, C.Y. Chan, V. Ganti, C.C. Kanne, W. Klas, E.J. Neuhold (eds.) Proceedings 33rd Intl. Conf. on Very Large Data Bases (VLDB), p. 1. Vienna, Austria (2007)
- [Wied92] Wiederhold, G.: Mediators in the architecture of future information systems. IEEE Computer 25(3), 38–49 (1992)
- [WCL97] Wong, R.K., Chau, H.L., Lochovsky, F.H.: A data model and semantics of objects with dynamic roles. In: A. Gray, P. A. Larson (eds.) Proceedings of the 13th International Conference on Data Engineering (ICDE), pp. 402–411. IEEE Computer Society, Birmingham, UK (1997)