



ELSEVIER

Data & Knowledge Engineering 22 (1997) 159–189

**DATA &
KNOWLEDGE
ENGINEERING**

Data schema design as a schema evolution process

H.A. Proper*

Cooperative Information Systems Research Centre, Faculty of Information Technology, Queensland University of Technology, GPO Box 2434, Brisbane, 4001, Australia

Received 13 February 1996; revised 30 August 1996; accepted 25 September 1996

Abstract

In an information system a key role is played by the underlying data schema. This article starts out from the view that the entire modelling process of an information system's data schema can be seen as a schema transformation process. A transformation process that starts out with an initial draft conceptual schema and ends with an internal database schema for some implementation platform. This allows us to describe the transformation process of a database design as an evolution of a schema through a universe of data schemas. Doing so allows a better understanding of the actual design process, countering the problem of 'software development under the lamppost'. Even when the information system design is finalised, the data schema can evolve further due to changes in the requirements on the system. We present a universe of data schemas that allows us to describe the underlying data schemas at all stages of their development. This universe of data schemas is used as a case study on how to describe the complete evolution of a data schema with all its relevant aspects. The theory is general enough to cater for more modelling concepts, or different modelling approaches. To actually model the evolution of a data schema we present a versioning mechanism that allows us to model the evolutions of the elements of data schemas and their interactions, leading to a better understanding of the schema design process as a whole. Finally, we also discuss the relationship between this simple versioning mechanism and general-purpose version-management systems.

Keywords: Conceptual modelling; Version management; Schema evolution; Object role modelling; Entity relationship modelling

1. Introduction

When designing an information system, one usually starts out by specifying a so-called conceptual schema. After such a conceptual schema is finalised, this schema is then implemented on a target platform such as a relational system, a hierarchical system or an object-oriented system. To this end, the conceptual schema is transformed (mapped) to a schema on the chosen target platform. In this article we take the view that the entire (data) modelling process can be seen as a transformation process of data schemas, where a data schema can be an Entity Relationship schema [3,21] an Object-Role Modelling schema [26], a relational schema, or any other internal representation schema.

*Corresponding author: email: erikp@icis.qut.edu.au

This view is in line with transformational approaches to software engineering in general [2,7,45] which take a high-level software specification and convert this by a sequence of transformations to concrete code. Even more, after the data schema of an information system has been implemented, the requirements on the system may have changed, requiring the information system to evolve as well. This means that the underlying data schema has to undergo yet another set of transformations to incorporate the new requirements. This observation has led to the development of so-called evolving information systems (see, for example, [6,28,43,46,47,54]).

1.1. Schema design through transformations

To further motivate our view on a schema modelling process as a transformation process, we first take a closer look at the involved transformations. A conceptual schema design process usually starts out with an initial draft conceptual schema which is then subject to a process of refinement and quality improvements, resulting in the final conceptual schema of the universe of discourse. This modelling process is usually guided by some design procedure. For example, in [26] a design procedure for Object-Role Modelling techniques is discussed. In [40,48] design procedures for object-oriented techniques are provided, while [3,21] discuss some loose guidelines for ER schema design.

After the conceptual schema has been finalised, one can sometimes perform small (equivalence preserving) transformations on these schemas which result in a schema that allows for a more efficient implementation [21,23–25,30,44]. These transformations typically utilise the rich semantics and clarity of conceptual schemas. Performing such transformations after these schemas have already been mapped to a target platform usually becomes too complicated, since these schemas use less concise modelling concepts. This makes it both harder to define the transformations and harder for the information system designers to track the transformations.

Following the optimisation transformations, the schema is consequently mapped to a target platform. For the different conceptual modelling techniques there are different mapping algorithms following varying styles and strategies; see, for instance, [4,8,10,26,42,49–51]. Once a conceptual schema has been transformed to some sort of representation for a target platform, this schema can sometimes be optimised even further. Some of these transformations are discussed in, for example, [3,19,39].

So far we have distinguished five key classes of schema transformations. They are:

1. Conceptual schema refinements
2. Conceptual schema quality improvements
3. (Conceptual) schema optimisations
4. Conceptual to internal (logical) schema mapping
5. Internal schema optimisations

For a more elaborate discussion on the classes of schema transformations, refer to, for example, [23]. Presently, the reverse process of the transformations in 4 and 5 also receives a lot of attention in the database and information systems research community. This reversed process is referred to as reverse engineering. For this fairly new area also a wide range of strategies and algorithms exists [11,22,27,34,35,37,53]. These algorithms all operate on the base of a set of possible schema transformations and heuristics to best apply them. The above discussed five classes either operate on a conceptual data schema or an internal schema of a given implementation platform. The modelling process up until the start of the internal schema mapping can be regarded as a journey through a

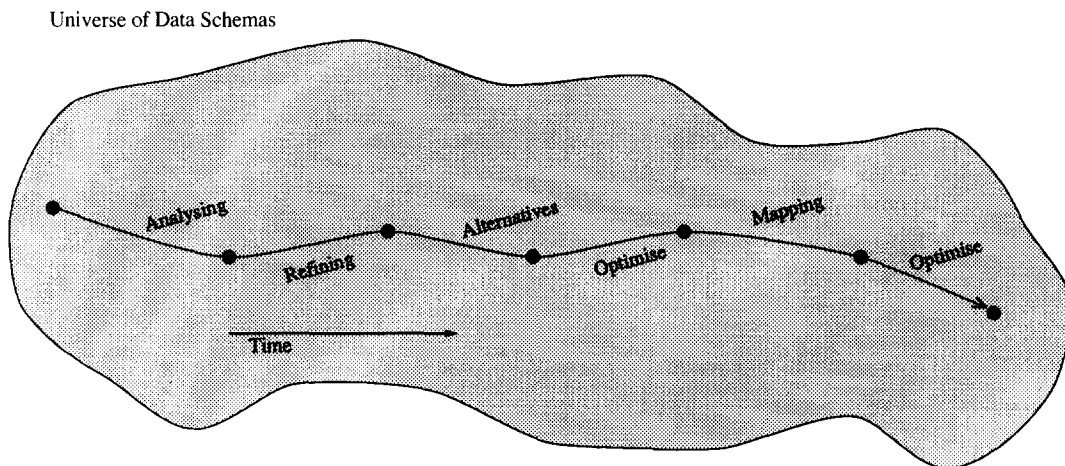


Fig. 1. Evolution of a data schema in time.

universe of data schemas. This is illustrated in Fig. 1. This universe of data schemas must be rich enough such that the data schemas can include details relevant for a conceptual presentation as well as the internal representation.

The aim of this article is to develop a mechanism by which we can describe the evolution of a data schema undergoing the above discussed schema transformations. Such a mechanism can then be used in the context of a CASE-Tool or an evolving information system. This evolution-describing mechanism consists of two key components, the first describing the state space of the universe of data schemas, and the second the historical aspects. In this article, focus is on the latter aspects. Nevertheless, we will have to introduce a state space that is rich enough for our requirements.

1.2. Overview of the version-management system

The state space is built from existing well publicized components [1,10,17,36]. Therefore, it is only discussed briefly in Section 2. What makes the universe of data schemas unique is that it is rich enough to describe (at multiple levels of abstraction) ER and ORM data models *simultaneously* on a conceptual and on an internal level. This is done by combining the CDM (Conceptual Data Modelling) Kernel [17] with a specification technique for internal representations [9,10]. This latter technique extends flat data schemas with a special notation that allows us to group roles together in trees. These trees then directly correspond to internal representations. Below we will see an example of this tree representation.

This combination of the CDM Kernel and the tree representations leads to the view shown in Fig. 2. The actual conceptual schema is provided by the flat data schema and the clustering information, which provides the distinction between attribute types, entity types, complex types and the layers of abstraction. The internal schema is provided by the flat data schema in conjunction with a tree representation of the roles in the conceptual schema, where each tree corresponds to one record structure of the internal schema. For convenience, the resulting data modelling technique is still referred to as the CDM Kernel in the remainder. This article can therefore actually be seen as the third and final part in a trilogy on the CDM Kernel. The first article [17] defining the kernel in full detail,

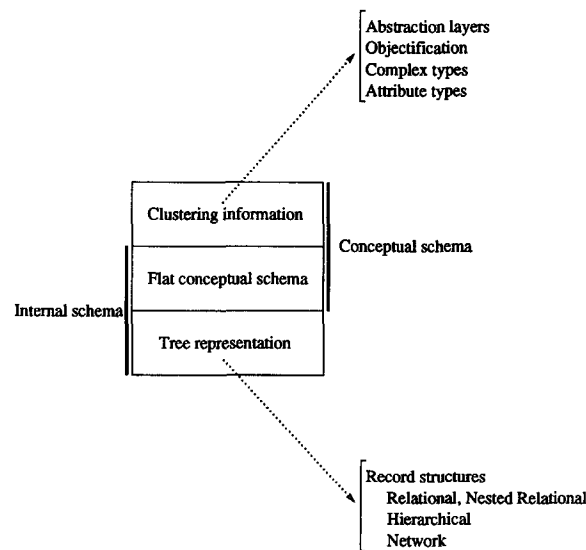


Fig. 2. A complete data schema.

the second one defining an algorithm for bottom-up abstraction of flat conceptual models [13] in the kernel, and, finally, this article, which adds a version-management system.

The second component of the versioning mechanism, introduced in Fig. 3, provides us with a way to model the evolution of the schema design. In our approach we are able to actually model any interactions schema components may have with each other during a schema design process. An example of this is shown in Fig. 3, where schema elements $e1$ and $e2$ fuse by reaction R to become $e3$, for instance, two relationship types that merge. This allows modellers to trace the evolution of schema components through a modelling or evolution process. This should provide the modeller (and even the end user) with more insight into the evolution of the schema design.

This latter issue should not be underestimated. Already in [16] it was stated that most existing software process models, including the (then) new approaches like prototyping and program transformation, “*focus on the series of artifacts that exist at the end of the phases of the process, rather than on the actual processes that are conducted to create the artifacts*”. Curtis et al, compared this phenomenon to the old story of the man who lost his wallet across the street, but searches for it on this side under the lamppost, simply because the light is better here. We expect that by providing

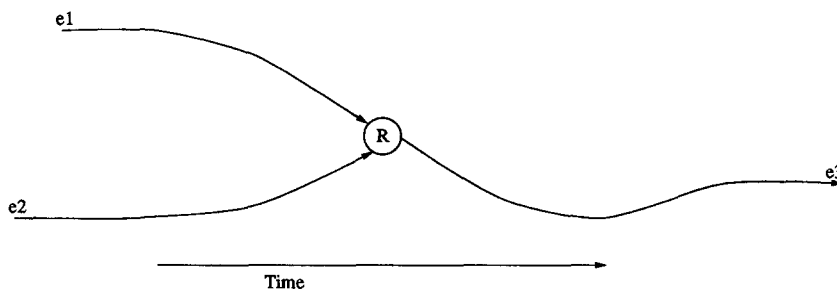


Fig. 3. A reaction between schema elements

better insight into the interactions between model elements during a schema design process, the light will also reach the other side of the street.

In [15,16], it was found that more insight into the actual modelling process, e.g. decision points, was crucial for a better understanding of the project as a whole. We believe that our version-management system provides a start for this in the context of schema design. By being able to explicitly model interactions of schema elements in the course of a schema design, and adding appropriate explanations to each of these interactions and design steps, a better understanding of the design process should result.

The integration between the schema universe and the version-management components is provided in Section 4. This yields a complete specification for an evolution mechanism that allows us to describe the evolution of a data schema through all its stages of design, including that of evolution of the information system as a whole due to changed requirements. The resulting evolution mechanism is referred to as the CDM Evolver.

Existing approaches to the modelling of data schema evolution (e.g. [6,46,54]) are only able to describe the evolution of either the conceptual level, or the internal level, but not both as a unit. For example, [6,54] describe the evolution at the actual database level and not the conceptual level. EVORM [46], on the other hand, only allows us to describe the evolution of a conceptual schema. In this article, we re-apply some of the principles used in the development of EVORM, which was derived from the Object-Role Modelling (ORM) variation PSM [36], while extending it to meet our new requirement to include abstraction layers and internal mapping aspects. Furthermore, the historical dimension of EVORM is further refined in the sense that we are, as stated before, able to explicitly model interactions between data schema elements, for example, the merger of two relationship types into one.

2. Data schema universe

The data schema universe is the set of valid data schemas. This set of data schemas is defined by the data schema language. As stated before, this language must be rich enough to allow us to model both the conceptual aspects of a data schema as well as internal representation aspects in one single model. In this section we therefore extend the CDM Kernel with an existing technique [10] to cater for this. The extension allows us to group roles into trees. Each tree corresponds directly to a record structure; it conforms to the (flat) relational model, the nested relational model, the network model or the hierarchical model. This technique was developed as a means to study alternative internal representations of the Predicator Model [5]. The Predicator Model is one of the first formalisations of Object-Role Modelling, and it was one of the ancestors of the CDM Kernel.

As the CDM Kernel is the result of a series of extensions and refinements of formalisations of ORM, each of which have been published before [1,5,17,36,46], we allow ourselves the luxury of only briefly discussing the current CDM Kernel. We start with a concise description of flat conceptual schemas in Section 2.1, followed in Section 2.2 by the abstraction layers that can be built on top of that. In Section 2.3, the tree representation of flat conceptual schemas is discussed. These trees are used to determine the internal representation of the conceptual schema. We also discuss briefly how these trees can be interpreted for different implementation platforms.

2.1. Flat conceptual schemas

In Fig. 4 an example, taken from [26], of a flat conceptual schema is shown. The domain is concerned with martial arts, and the schema should be self explanatory. The notation used is the style of Object-Role Modelling (ORM). For a flat conceptual schema we cannot directly use the ER

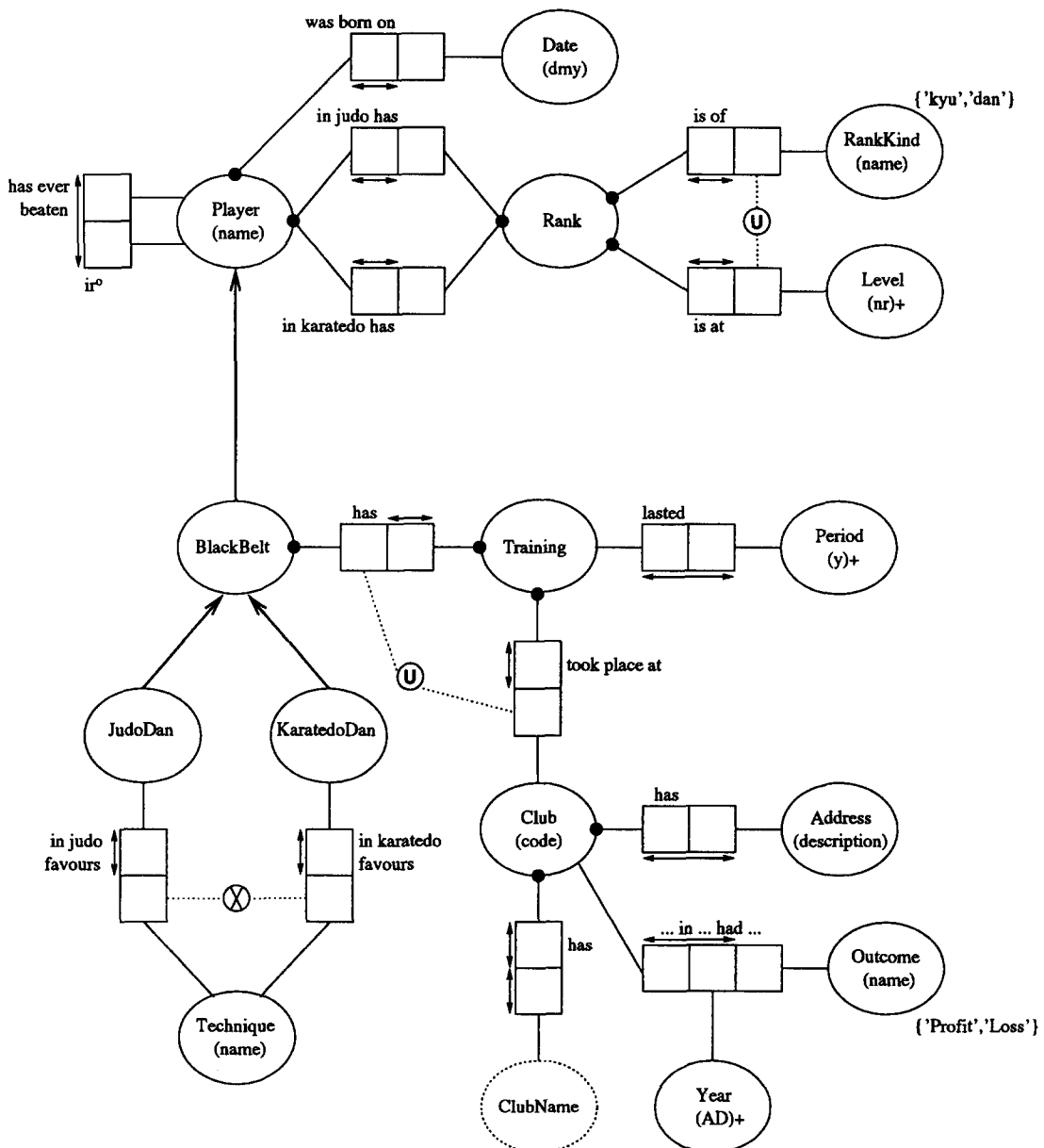


Fig. 4. Martial arts example.

notation, as no information is available to distinguish between entity types and attribute types. For more details on this refer to [17].

Formally, a flat conceptual schema version at point in time t consists of the following components:

1. A set \mathcal{OB}_t of object types.
2. A subset $\mathcal{VL}_t \subseteq \mathcal{OB}_t$ of value types. Value types are types that have instances which are directly denotable on some communication medium. For example, numerals, strings, sound and video. The function $\text{Dom}_t: \mathcal{VL}_t \rightarrow \mathcal{DO}_t$ is used to associate a domain of values to each value type.
3. A set \mathcal{RL}_t of relationship types. The roles of the relationship types are provided as \mathcal{RO}_t , while the function $\text{Roles}_t: \mathcal{RL}_t \rightarrow \wp^+(\mathcal{RO}_t)$ partitions the roles of the relationship types. The players of the roles are given by the function $\text{Player}_t: \mathcal{RO}_t \rightarrow \mathcal{OB}_t$. The set of all types is defined as $\mathcal{TP}_t \triangleq \mathcal{OB}_t \cup \mathcal{RL}_t$.
4. The subtyping hierarchy is captured by the predicate $\text{SubOf}_t \subseteq \mathcal{OB}_t \times \mathcal{OB}_t$.
5. The identification schemes of object types are given by the function $\text{Ident}_t: \mathcal{OB}_t \rightarrow (\wp(\mathcal{RO}_t \times \mathcal{RO}_t))^+$. The CDM Kernel allows identification schemes to be inherited from supertypes onto subtypes, but subtypes may also provide their own identification scheme.
6. A set \mathcal{CN}_t of constraints.
7. Some types in the conceptual schema may be derivable (most notably subtypes). Therefore, $\text{DerRule}_t: \mathcal{TP}_t \mapsto \text{DerivationRules}$ may associate a derivation rule to types. Furthermore, for each derivable type an update rule may be specified defining how to deal with updates on the derived types: $\text{UpdRule}_t: \mathcal{TP}_t \mapsto \text{UpdateRules}$.

In this article we do not elaborate on correctness rules for single schema versions. For a detailed discussion on these rules see [17]. For the martial arts domain we have the following excerpt:

\mathcal{OB}_t	=	{Date, dmy, Player, PlayerName, Rank,..., Outcome}
\mathcal{VL}_t	=	{dmy, PlayerName, RankKindName,..., OutcomeName}
$\text{Dom}_t(\text{dmy})$	=	{ $'d-m-y'$ $1 \leq d \leq 31 \wedge 1 \leq m \leq 12 \wedge 0 \leq y$ }
$\text{Dom}_t(\text{PlayerName})$	=	String
$\text{Dom}_t(\text{RankKindName})$	=	{ $'kyu'$, $'dan'$ }
$\text{Dom}_t(\text{LevelNr})$	=	\mathbb{N}
\mathcal{RL}_t	=	{has ever beaten, was born on, in judo has,..., in karatedo favours}
$\text{Roles}_t(\text{has ever beaten})$	=	{has ever beaten-1, has ever beaten-2}
$\text{Roles}_t(\text{was born on})$	=	{was born on-1, was born on-2}
$\text{Roles}_t(\text{is at})$	=	{is at-1, is at-2}
$\text{Player}_t(\text{was born on-1})$	=	Player
$\text{Player}_t(\text{was born on-2})$	=	Date
$\text{Player}_t(\text{has ever beaten-1})$	=	Player
$\text{Player}_t(\text{has ever beaten-2})$	=	Player
$\text{BlackBeltSubOf}_t(\text{Player})$	=	JudoDanSubOf _t BlackBelt, KaratedoDanSubOf _t BlackBelt

The derivation rules for this domain are the three subtype defining rules as specified in the graphical representation.

2.2. Abstraction layers

In the past, a much heard critique on Object-Role Modelling based techniques was that they were too detailed. Where ORM models show too much detail, (E)ER models lack detail. The reason being that (E)ER, by virtue of its attribute types, already provides a rudimentary means to introduce a single abstraction level. This allows modellers to focus on what they experience as key object types (which then become the entity types). ORM forces modellers to initially regard all object types as equals. This means that an ORM diagram looks initially much more complex than an ER diagram would. For large applications, however, the problem of complex and incomprehensible schemas also haunts the (E)ER modeller [14].

These observations have sparked the development of bottom-up abstraction algorithms that allow for (manual, semi-manual or automatic) generation of higher abstraction layers by identification of so-called major object types. Examples of such algorithms can be found in, for example, [13,14,52,55,56].

The CDM Kernel provides modelling constructs that allow for top-down abstraction as well as bottom-up abstraction. In, for example, [18,20], similar extensions of ORM are discussed. However, the ideas presented there remain informal, and do not establish a connection to (E)ER modelling. Furthermore, the fully automatic bottom-up abstraction algorithm described in [13] generates these abstraction layers for a given flat conceptual model in the CDM Kernel. This makes the CDM Kernel highly suitable for both bottom-up and top-down modelling.

Besides the above discussed notion of abstraction, traditional ORM and ER schemas already (implicitly) feature a different notion of abstraction. In both traditional ORM and ER, but in particular in the many extensions, it was found natural to present complex types as an undividable entity. This has led to such notions as: objectification, collection types, sequence types, aggregate types, bag types, schema types, etc. As an example of this, consider the **Training** object type in Fig. 4. From a conceptual point of view, it may be more natural to treat **Training** as an objectification of a relationship between a **BlackBelt** and **Club**, so each relation between a **BlackBelt** and a **Club** is treated as if it is an individual instance of an object type. This objectified view is shown in Fig. 5. Which view is more natural depends very much on the underlying universe of discourse. In ORM this is usually detected by studying the way example facts from the universe of discourse are verbalised. The objectified view of Fig. 5 hides the underlying relationship types **has** and **took place at**, and by doing so it provides a form of abstraction.

To return to the martial arts example, after identifying the abstractions to hide the details of complex types, we can perform abstractions based on the importance of object types. For the martial arts domain, this may lead to the representation as shown in Fig. 6. This representation focuses on the major object types in the domain, and is likely to be more comprehensible by modellers and their discussion partners. It is now interesting to see that a level-one abstraction from a flat object-role model directly corresponds to an (E)ER schema. In Fig. 7 we have depicted the (E)ER version of the schema in Fig. 6. This example also illustrates how (E)ER [3] schemas are represented in the context of the CDM Kernel. They are treated as flat conceptual schemas with a predetermined level-one abstraction layer.

Formally, the abstraction layers for the CDM Kernel are provided by the clustering function:

$$\text{Cluster}_i: \mathbb{N} \times \mathcal{OB}_i \rightarrow \wp(\mathcal{TP}_i)$$

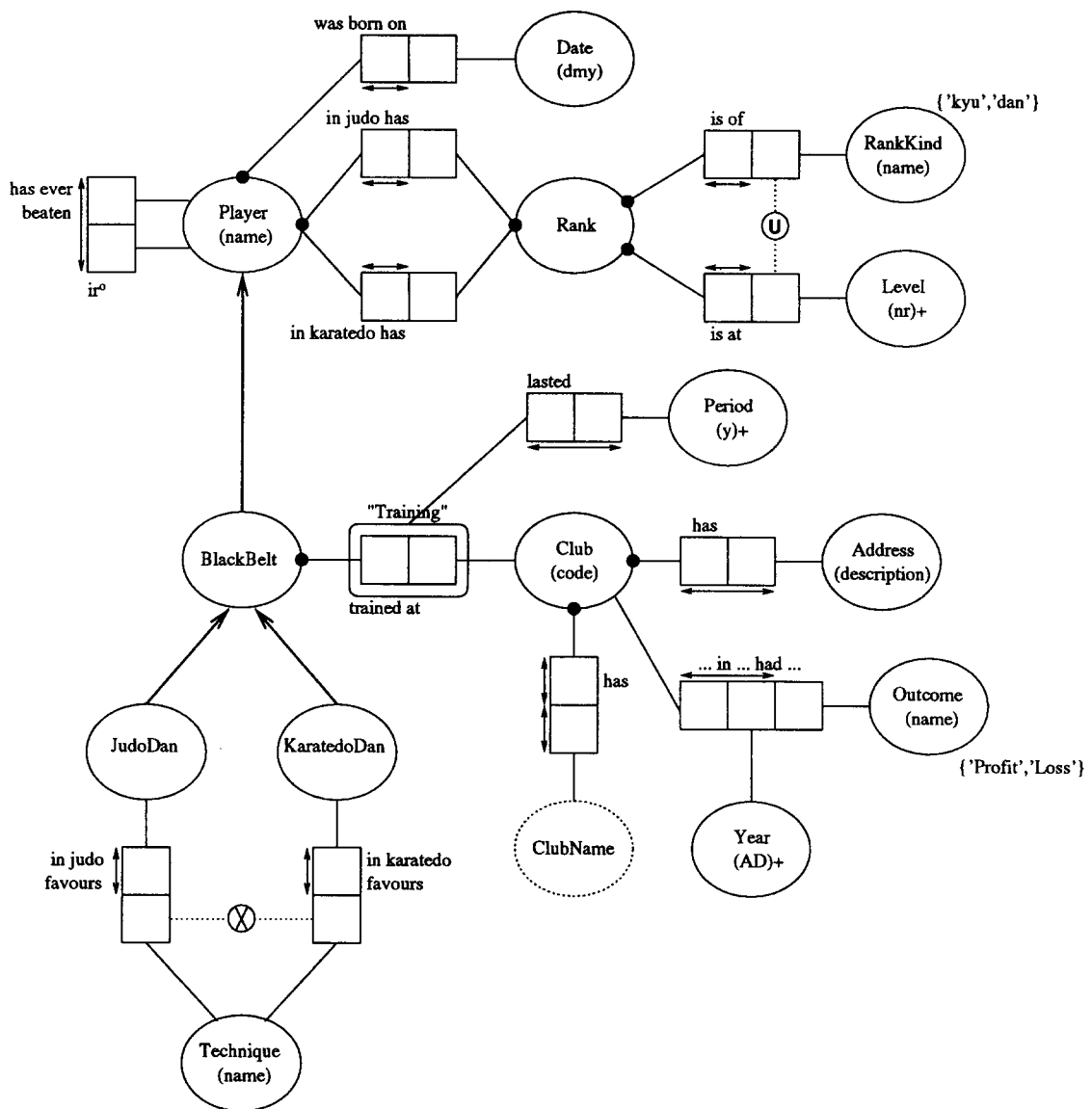


Fig. 5. Objectification flavour of abstraction.

The intuitive meaning is that if $\text{Cluster}_i(i, x) = Y$, then the types in Y are clustered to x on level i . In [17] some general completeness rules for these clusterings are provided. For the martial arts domain we would for example have:

$$\begin{aligned}
 \text{Cluster}_i(0, \text{Training}) &= \{ \text{Training, has, BlackBelt, took place at, Club} \} \\
 \text{Cluster}_i(1, \text{Player}) &= \{ \text{Player, was born on date, has ever beaten, ..., dmy} \} \\
 \text{Cluster}_i(1, \text{Rank}) &= \{ \text{Rank, is of, RankKind, ..., LevelNr} \}
 \end{aligned}$$

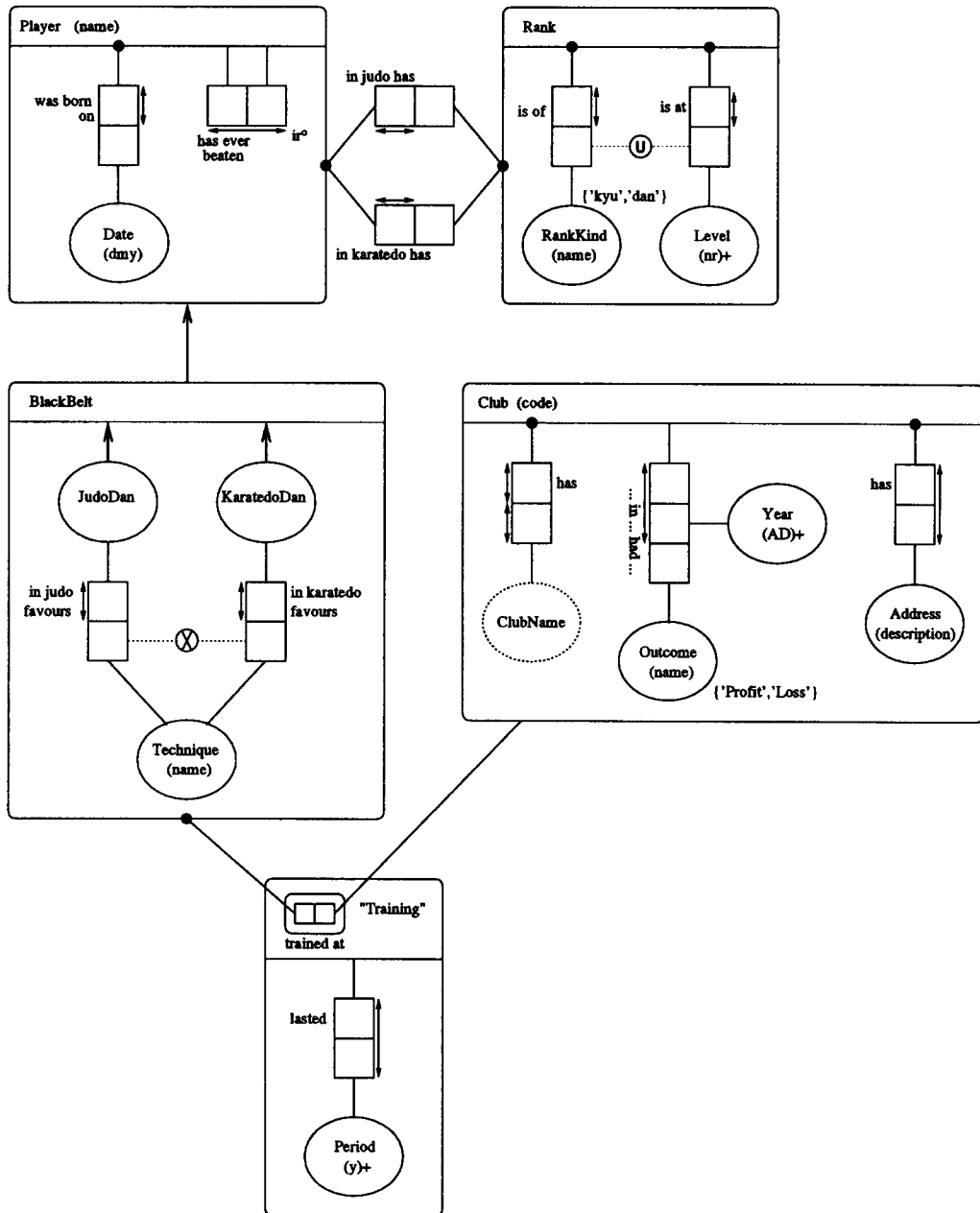


Fig. 6. Abstractions.

To distinguish between the different flavours of abstraction we also introduce the function:

$$\text{CFlavour}_i: \mathcal{OB}_i \times \mathcal{TP}_i \rightarrow \text{Flavours}$$

also, using intuition, if $\text{CFlavour}_i(x, y) = f$, then the clustering of object type x to y has flavour f . For example:

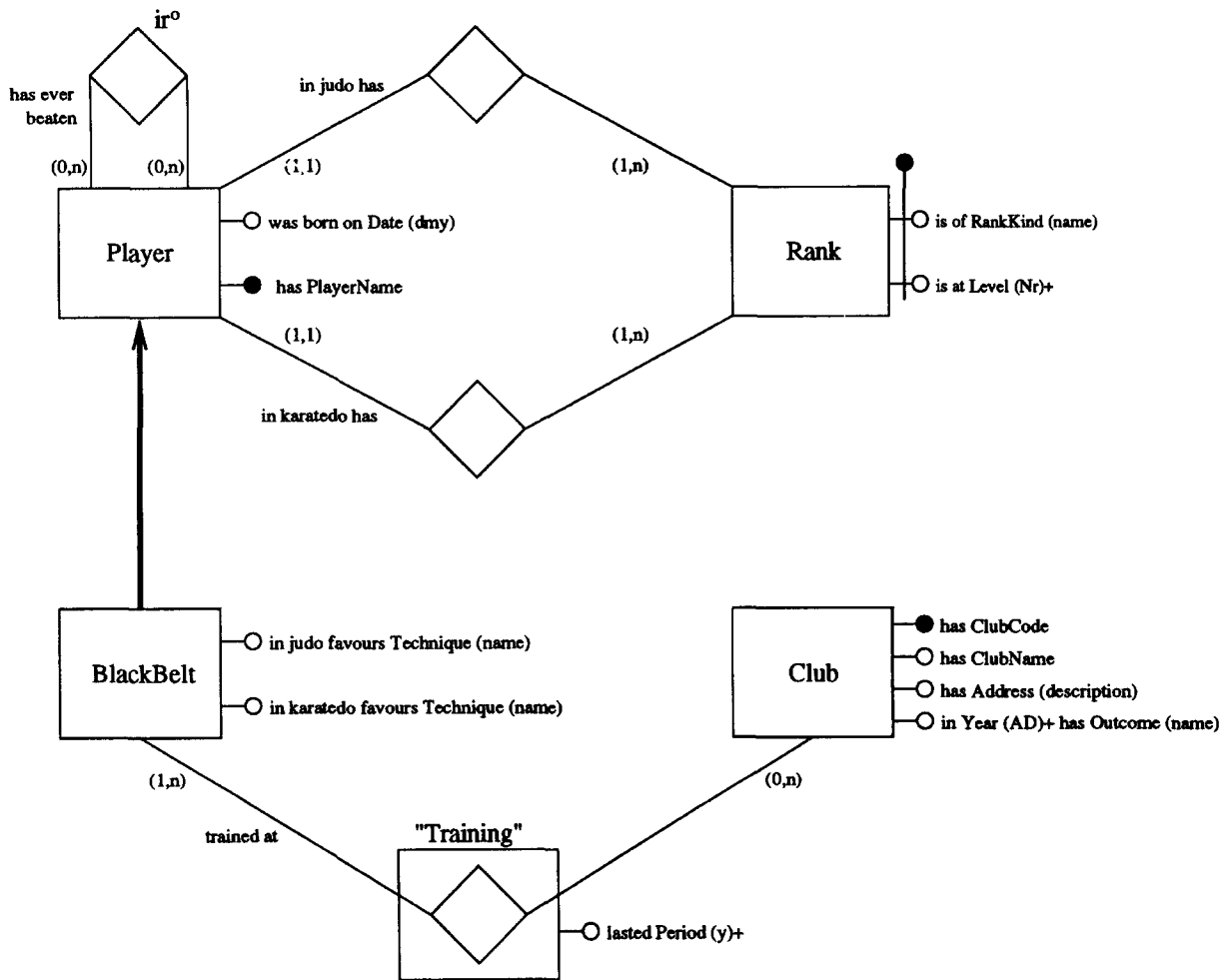


Fig. 7. An ER view on the martial arts domain.

$$\begin{aligned}
 \text{CFlavour}_i(\text{has}, \text{Training}) &= \text{objectification} \\
 \text{CFlavour}_i(\text{took place at}, &= \text{objectification} \\
 &\quad \text{Training})
 \end{aligned}$$

The clusters resulting from CFlavour can also be interpreted as object classes in an object-oriented sense. Each of the clusters in Fig. 6 can be seen as an object class in an object-oriented approach. This view also allowed us to introduce some object-oriented aspects into traditional data modelling [17]. In [12,20,41] the effects of extending ORM or ER with object-oriented features are also discussed. The important aspects with which ORM is extended are: the overriding of inherited relationship types, the association of methods to object types, and the encapsulation of methods and clustered object types. These aspects are discussed in full detail in [17]. In this article we only state the way in which we have 'implemented' these features in the theory. To capture overriding of inheritance, we introduce the function $\text{RoleLim}_i: \mathcal{OB}_i \times \mathcal{RO}_i \rightarrow \mathcal{OB}_i$. If $\text{RoleLim}_i(x, p) = y$, then, in the context of the relationship

types clustered to x , the population of $\text{Player}_i(p)$ should be limited to the population of object type y . The second object-oriented aspect is concerned with the association of methods to object types. To each schema a set $\mathcal{OP}_i \subseteq \mathcal{OB}_i \rightarrow \text{Methods}$ of operations can be associated, accompanied by a signature function $\sigma_i: \mathcal{OP}_i \rightarrow \mathcal{TP}_i^+$. If $o \in \mathcal{OP}_i$, we have an operation with signature $\sigma_i(o)$, while the operation o itself is a function assigning a specific method for different object types (within one single subtyping hierarchy). Obviously, the methods $o(x)$ may differ for different types x in the same type hierarchy, but inheritance is the default. The way in which we use the terms operation and method is borrowed from [48]. The operations themselves are also introduced on different levels of abstraction. These abstraction levels are provided by the function: $\text{Ops}_i: \mathbb{N} \rightarrow \wp(\mathcal{OP}_i)$.

Finally, as a third object-oriented aspect, the CDM Kernel offers encapsulation. The CDM Kernel offers two flavours of encapsulation. The most liberal one is encapsulation on the type level, which is provided by the function $\text{Encap}_i: \mathcal{OB}_i \rightarrow \wp(\mathcal{OP}_i \cup \mathcal{TP}_i)$. This flavour of abstraction allows us to encapsulate operations and relationship types within definitions of clusterings. Only relationship types can be encapsulated, since the players of relationship may be shared among relationship types. Similarly, in most object-oriented modelling techniques attributes can be encapsulated, but the underlying domains of these attributes can not be encapsulated. If relationship type $r \in \text{Encap}_i(x)$, then r is only visible from instances of x .

2.3. Internal representation

In [10] a representation technique is introduced that allows us to specify the internal representation of a conceptual schema as a forest of trees, where each tree corresponds to one record structure and each node in the tree consists of a set of roles. An example of this is given in Fig. 8. The top part of this figure displays one tree corresponding to one single relational table. This table is shown in the bottom part. Please note that one-on-one relationships, like the one between **Rank** and the combination of **RankKindName** and **Level**, are ignored in the relational table representation. This can be done due to the fact that **Rank** is identified by means of a **RankKind** and **Level**, which, in turn, receive their identification from **RankKindName** and **LevelNr**. As each **Rank** is thus identified by a **RankKindName** and a **LevelNr**, we can replace **Rank** completely by this combination in the relational table representation. As argued in [9,10], the tree representation can be used for non-relational target platforms as well.

Formally, a forest representation of a flat data schema is given by a set of nodes \mathcal{NO}_i , a partition $\text{Roles}_i: \mathcal{NO}_i \rightarrow \wp(\mathcal{RO}_i)$ yielding the roles grouped to the given node (rendering the Roles_i function symbol overloaded), and a set of labelled edges $\text{Edge}_i \subseteq \mathcal{NO}_i \times \mathcal{NO}_i \times \mathcal{RL}_i$. In [10] a number of correctness rules for the resulting forest are given. These rules can actually be simplified slightly as we now treat objectification as a form of abstraction, leading to a simplification of the flat conceptual schema level, whereas in [10] exceptions were needed to cope with objectifications.

Please note that it is quite easy to extend the data modelling language with additional constructs to include, for example, indexing options for the internal representation. For instance, a predicate $\text{Index}_i \subseteq \wp(\mathcal{RO}_i)$ where $\text{Index}_i(P)$ signifies that there is an index defined on the combination of roles P . Obviously, the roles in P must be part of the same tree (relational table) to make it useful to have an index defined on them. Another relevant extension would be to cater for distribution of the conceptual schema over different sites. This could, for instance, (simplistically) be modelled by a function $\text{Site}_i: \mathcal{TP}_i \rightarrow \wp(\text{Sites})$.

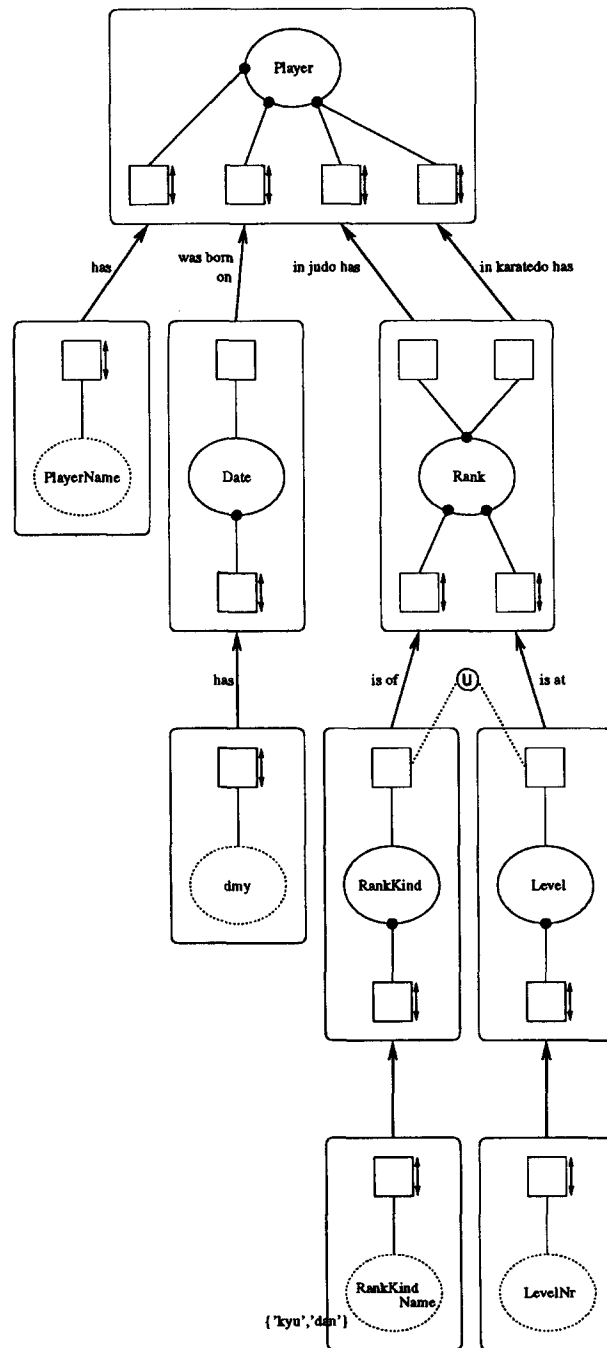


Fig. 8. Internal representation by clustering of roles.

2.4. Summary of components

This completes the overview of the data schema universe. We can now put the different components together. The structural aspects of a data model are fully determined by the components of the following tuple:

$$\mathcal{IS}_i = \langle \mathcal{RL}_i, \mathcal{OB}_i, \mathcal{VL}_i, \mathcal{RO}_i, \mathcal{NO}_i, \text{SubOf}_i, \text{Roles}_i, \text{Player}_i, \text{Cluster}_i, \\ \text{CFlavour}_i, \text{RoleLim}_i, \text{Ident}_i, \text{Edges}_i \rangle$$

The first 5 components provide the types, roles and nodes present in the information structure, and the last 9 components describe their mutual relationships providing the ‘fabric’ of the information structure (together with the abstraction layers and internal representation).

A complete conceptual schema over a set of concrete domains \mathcal{DO} , is then identified by the following components:

$$\mathcal{CS}_i \triangleq \langle \mathcal{IS}_i, \mathcal{CN}_i, \mathcal{OP}_i, \sigma_i, \text{Ops}_i, \text{DerRule}_i, \text{UpdRule}_i, \text{Encap}_i, \text{Dom}_i \rangle$$

At the moment work is underway in establishing a similar universe of models for the process aspects of an information system. The CDM Evolver as it will be presented in the next section is flexible enough to be applied directly to such a process modelling technique as well.

In [17] a series of well-formedness rules for the CDM Kernel is provided. In this article we simply presume that the predicate $\text{IsSch}(\mathcal{CS}_i)$ determines whether a given conceptual schema \mathcal{CS}_i is correct.

3. Model evolution

The previous section briefly discussed the substance of transformation/evolution, i.e. the corpus evolutionis. This section focuses on a way to model the actual evolution of the corpus evolutionis. We try to do this in a generalised way such that it could also be applied to domains other than data schemas (for example, process models). The next section provides the actual coupling between the components of a data schema in the CDM Kernel and the versioning mechanism. In Section 3.6, we return to the earlier mentioned relationship between the way we model evolution of data schemas and generic version-management systems.

The model evolution framework we use consists of the following components:

$$\langle \mathcal{TS}_i, \mathcal{ELV}, \mathcal{ELC}, \text{VClass} \rangle$$

A linear time axis is provided by \mathcal{TS}_i . To cater for the fact that model elements can have different versions, we introduce a set \mathcal{ELV} of \mathcal{E} lement \mathcal{V} ersions. The elements of this set are abstract representations of the underlying version of model elements. A simple classification system, e.g. distinguishing between types and constraints, is provided for these elements by a set of \mathcal{E} lement \mathcal{C} lasses: \mathcal{ELC} . The function $\text{VClass}: \mathcal{ELV} \rightarrow \mathcal{ELC}$ provides the class for each element version. Below we discuss the aim and pragmatics of these components in more detail. In the next section we provide a concrete definition of \mathcal{ELV} and \mathcal{ELC} for the CDM Kernel.

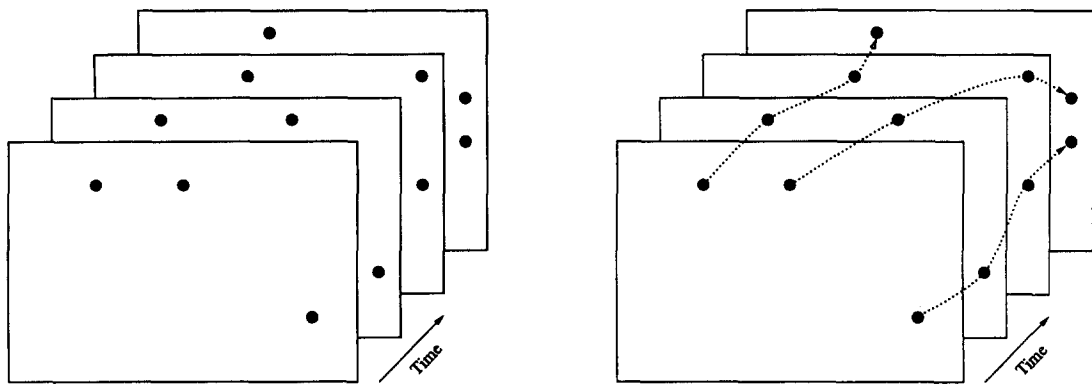


Fig. 9. Granularity of versioning.

3.1. Granularity of versioning

A data schema, or any model for that matter, can be seen as a set of elements. Some of these elements may be atomic while some of them may be composed, e.g. relations between other elements. In a version-management system, one can decide to maintain different versions of models as a whole, or maintain different versions of the elements. We opt for the latter approach. The advantage is that the version history is not just a series of snapshots of the complete model, but one that can actually track the version history (evolution) of single elements, for instance the evolution of a relationship type during the modelling process. As an illustration, consider Fig. 9. The left hand side of this figure depicts a sequence of snapshots, whereas the right hand side shows the same sequence of snapshots, but this time we can see the evolution of the individual elements.

Below we will see how we are able to model the interaction of different elements in the course of time. For instance, consider the schema transformation shown in Fig. 10. This transformation is taken from [44], and is an example of an equivalence-preserving schema transformation. In discussing reactions between element histories, we adopt the terminology used in quantum physics. For the example transformation (left to right) we can say that the type of the reaction is a *fact type reaction*, where *teaches* and *advises* *fuse* to become... *performed... for...*. This reaction *absorbs* the uniqueness constraints on the two original fact types. The reaction *emits* the uniqueness constraint on

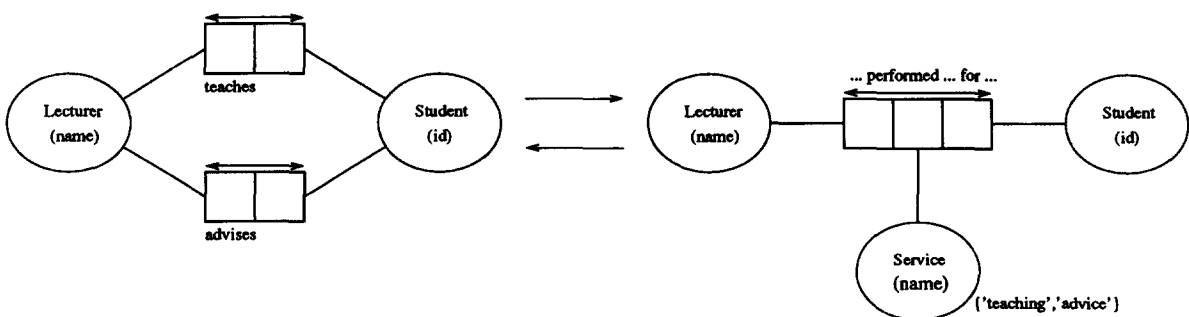


Fig. 10. Example schema transformation.

the resulting fact type, together with the **Service** object type and its associated value type (**name**) and value constraint {'teaching', 'advice'}. We can model this behaviour explicitly in our versioning mechanism, providing deeper insight into the actual design processes, moving away from the traditional snapshot view (i.e. searching under the lamppost).

3.2. Time axis

To model the evolution of a model, some linear time axis is required. In Section 3.6 we discuss how to deal with alternative evolution versions. The time axis we use is provided as: $\mathcal{TL}_s \triangleq \langle \mathcal{TI}, \{<\} \rangle$, where \mathcal{TI} is a set of points in time, and $<$ provides a complete total order on \mathcal{TI} . The $<$ relation allows us to define a one-step increment function $\triangleright: \mathcal{TI} \rightarrow \mathcal{TI}$ on the time axis:

$$\triangleright t_1 = t_2 \Leftrightarrow t_1 < t_2 \wedge \neg \exists_s [t_1 < s < t_2]$$

From the order on the time axis we can also derive:

$$t_1 = t_2 \Leftrightarrow \neg(t_1 < t_2 \vee t_2 < t_1)$$

$$t_1 \leq t_2 \Leftrightarrow \neg(t_2 < t_1)$$

3.3. Element evolutions

As stated before, we model the versioning history of each element separately. To elegantly model the version history of these elements, and their interactions, we introduce the notion of an extra-temporal (time independent) *element identifier*. This concept is indeed similar to the object identifier notion of object-oriented approaches. Let \mathcal{EID} be the set of element identifiers. At each point in time, an element may have a version. To remain as general as possible, we do not yet want to elaborate on what such a version actually is. Therefore, we use the set \mathcal{ELV} of version elements to identify these versions, and for now we treat these versions as abstract objects. This means that we can now associate to each element identifier a function $\mathcal{TI} \mapsto \mathcal{ELV}$ describing the evolution of the element in terms of its versions in the course of time. The history of all element identifiers, the complete model history, can, therefore, be seen as a function: $H: \mathcal{EID} \rightarrow (\mathcal{TI} \mapsto \mathcal{ELV})$. The set of all possible *Model HI* stories is then: $\mathcal{MHI} \triangleq \mathcal{EID} \rightarrow (\mathcal{TI} \mapsto \mathcal{ELV})$.

A model history as such does not capture the full spectre of the evolution of a model. What is still missing are interactions that have taken place between model elements in the course of time. This is the extra spice in the meal that makes a model history into a true model evolution. The interactions of elements in the course of their evolution can be modelled as a tuple $\langle E, c, t \rangle$, where $E \subseteq \mathcal{EID}$ is a set of element identifiers which are together involved in some form of *reaction* of type $c \in \mathcal{ELC}$ at point in time t . The general set of such *REACTIONS* is: $\mathcal{REA} \triangleq \wp(\mathcal{EID}) \times \mathcal{ELC} \times \mathcal{TI}$. On reactions, the following access functions can be defined:

$\text{RTIME}(\langle E, c, t \rangle)$	\triangleq	t
$\text{Reagents}(\langle E, c, t \rangle)$	\triangleq	E
$\text{RClass}(\langle E, c, t \rangle)$	\triangleq	c

A complete evolution of models built from version identifiers \mathcal{ELV} with element identifiers \mathcal{EID} can now be described by a tuple:

$$\langle R, H \rangle \in \wp(\mathcal{REA}) \times \mathcal{MHI}$$

where R provides the reactions and H provides the historical component. The set of all possible Model EVolutions is then:

$$\mathcal{MEV} \triangleq \wp(\mathcal{REA}) \times \mathcal{MHI}$$

3.4. Correctness of model evolutions

The CDM Kernel provides a number of rules to which a conceptual schema should adhere. Some of these rules are actually formulated as optional rules (*electronic switches*) to cater for well-formedness variations based for the different data modelling schools. In this subsection we introduce some general well-formedness rules on model evolutions $\langle R, H \rangle \in \mathcal{MEV}$. It should be noted that also for the model evolutions one may chose to add extra rules based on one's philosophical stance with respect to conceptual data modelling. Here we can do nothing more but provide a general framework.

If all of the MEW (Model Evolution Well-formedness) axioms we introduce in this section hold, then a model evolution $\langle R, H \rangle \in \mathcal{MEV}$ is deemed correct: $\text{IsModEvol}(R, H)$. In the definitions provided in this section we shall use the following abbreviation $h \downarrow t \triangleq t \in \text{dom}(h)$, where $\text{dom}(h)$ is the domain of function h . So if $h \downarrow t$, then function h is defined for t .

An element evolution cannot be empty:

$$[\text{MEW1}] e \in \mathcal{EID} \Rightarrow \exists, [H(e) \downarrow t]$$

The nature of the relationship between reactions and element histories is captured by the following two axioms:

[MEW2] If $r \in R$, then:

$$\forall_{e \in \text{Reagents}(r)} [H(e) \downarrow \text{RTIME}(r) \vee H(e) \downarrow \triangleright \text{RTIME}(r)]$$

[MEW3] If $e \in \mathcal{EID}$ such that $H(e) \downarrow t \wedge H(e) \downarrow \triangleright t$, then:

$$H(e)(t) \neq H(\triangleright t) \Rightarrow \exists_{r \in R} [e \in \text{Reagents}(r) \wedge \text{RTIME}(r) = t]$$

The first axiom states that all elements involved in some reaction must be alive at, or immediately after, the reaction takes place. In plain words this means that dead elements cannot partake in any reaction. The second axiom requires all changes in an element's evolution to be the result of some reaction.

The next rule states that evolution of elements is bound to classes. For example, a type may not evolve into a method, and a constraint may not evolve into a derivation rule. In other words there is no magic; we do not allow for frogs to turn into princesses. This is formalised in the following axiom:

[MEW4] (history separation) If $e \in \mathcal{EID}$, then:

$$H(e) \downarrow t \wedge H(e) \downarrow t' \Rightarrow \text{VClass}(H(e)(t)) = \text{VClass}(H(e)(\triangleright t))$$

In concrete instances one may sometimes debate whether this rule should be enforced as a hard rule,

or as a *deontic* rule. In the next section we will see that, although we distinguish between different types in a conceptual schema, an element evolution describing the evolution of a type is allowed to ‘roam’ in the set of all types.

The above axiom allows us to introduce a classification of element identifiers. This classification is defined in the context of a version history H using the existing classification provided by **VClass**:

$$\text{EClass}_H: \mathcal{EID} \rightarrow \mathcal{ELC}$$

$$\text{EClass}_H(e) \triangleq \text{the unique } c \text{ such that } \forall_t [H(e) \downarrow t \Rightarrow \text{VClass}(H(e)(t)) = c]$$

Using the classification on element identifiers, we are also able to take a closer look at the kinds of reactions that can occur between element histories. The input and output of a reaction are identified by the functions:

$$\text{Input}_H: \mathcal{REA} \rightarrow \wp(\mathcal{EID})$$

$$\text{Input}_H(r) \triangleq \{e \in \text{Reagents}(r) \mid H(e) \downarrow \text{RTime}(r)\}$$

$$\text{Output}_H: \mathcal{REA} \rightarrow \wp(\mathcal{EID})$$

$$\text{Output}_H(r) \triangleq \{e \in \text{Reagents}(r) \mid H(e) \downarrow \triangleright \text{RTime}(r)\}$$

A catalyst is an element that partakes in a reaction without being changed itself. For our reactions we can identify catalysts by:

$$\text{Catalyst}_H: \mathcal{REA} \rightarrow \wp(\mathcal{EID})$$

$$\text{Catalyst}_H(r) \triangleq \{e \in \text{Input}_H(r) \cap \text{Output}_H(r) \mid H(e)(t) = H(e)(\triangleright t)\}$$

If a reaction r is stated to have class **RClass**(r) then this must be visible in the input and output of the reaction. So there must be some input or output of this class:

[MEW5] If $r \in R$, then:

$$c \in \text{RClass}(r) \Rightarrow \exists_{h \in \text{Input}_H(r) \cup \text{Output}_H(r)} [\text{EClass}_H(h) = c]$$

In a reaction some element histories of the same class may fuse, and some may be split. These respective elements are identified by:

$$\text{Fusion}_H: \mathcal{REA} \rightarrow \wp(\mathcal{EID})$$

$$\text{Fusion}_H(r) \triangleq \{h \in \text{Input}_H(r) - \text{Catalyst}_H(r) \mid \text{HClass}_H(h) = \text{RClass}(r)\}$$

$$\text{Fission}_H: \mathcal{REA} \rightarrow \wp(\mathcal{EID})$$

$$\text{Fission}_H(r) \triangleq \{h \in \text{Output}_H(r) - \text{Catalyst}_H(r) \mid \text{HClass}_H(h) = \text{RClass}(r)\}$$

Sometimes an element history may absorb element histories of another class, or vice versa. For instance a constraint may be absorbed when an objectification is transformed into a ternary. To identify such histories we introduce:

$$\text{Absorption}_H: \mathcal{REA} \rightarrow \wp(\mathcal{ED})$$

$$\text{Absorption}_H(r) \triangleq (\text{Input}_H(r) - \text{Fusion}_H(r)) - \text{Catalyst}_H(r)$$

$$\text{Emission}_H: \mathcal{REA} \rightarrow \wp(\mathcal{ED})$$

$$\text{Emission}_H(r) \triangleq (\text{Output}_H(r) - \text{Fission}_H(r)) - \text{Catalyst}_H(r)$$

As an example of what can happen during a schema transformation, consider the two Olympic Games schemas depicted in Fig. 11. The transformation from the first schema to the second schema is based on the observation that there are exactly three kinds of medals, so the ternary fact type may be specialised into three binaries, one for each medal kind. You may visualise the transformation from schema (a) into schema (b) thus:

Fusion of:

won medals of... in, MedalKind (code)

Fission to:

won gold in, won silver in, and won bronze in

Absorption of:

uniqueness constraint over first two roles of won medals of... in,

enumeration constraint f'G', 'S', 'B' g

Emission of:

uniqueness constraints on first roles of won gold in, won silver in, and won bronze in

What we hope to illustrate with this example, despite its simplicity, is that explicit modelling of the interaction between elements during their evolution provides more insight into the evolution process. One could even explicitly include modelling decisions by storing a description of the schema

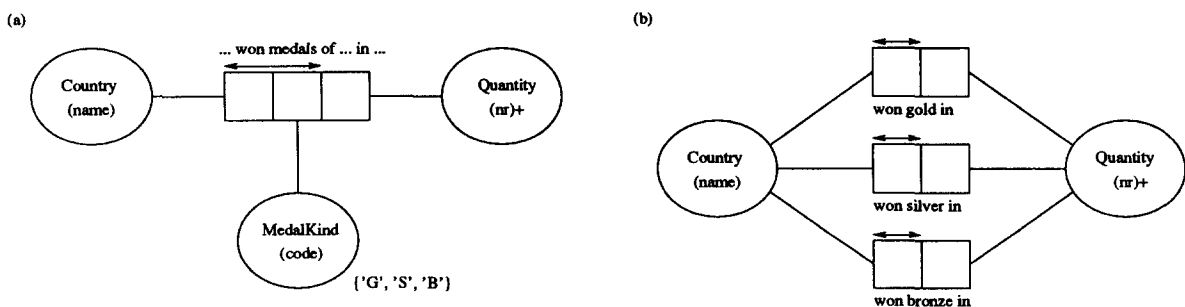


Fig. 11. Olympic Games universe of discourse.

transformation that caused the 'reaction', a short explanation of why it was applied, and what role each of the involved elements plays in the reaction. So:

ResultsFrom: $\mathcal{REA} \rightarrow \text{SchemaTransformation}$

Reason: $\mathcal{REA} \rightarrow \text{Text}$

Role: $\mathcal{REA} \times \rightarrow \text{Text}$

This, we believe, should provide modellers (and project managers for that matter) with the additional insight into the modelling process as needed according to [15,16]. We hope that this way, the light will start shining on the other side of the street as well.

3.5. Derivation of model versions

A complete model snapshot at a given point in time t is easily derived from a complete evolution $\langle R, H \rangle$ by determining the set of version identifiers in the snapshot:

$\text{Snapshot}_H: \mathcal{I} \rightarrow \wp(\mathcal{ELV})$

$\text{Snapshot}_H(t) \triangleq \{H(e)(t) | H(e) \downarrow t\}$

3.6. Relation to generic version management

We are now in a position to take a closer look at the relationship between the version-management system proposed in this article, and the more generic version-management systems as discussed in, for example [6,38].

The version-management system proposed in this article focuses on information system modelling processes, and, in particular, the evolution of the resulting models in the course of time. In doing so we initially do not cater for things like alternative versions of models. Initially, we only acknowledge the existence of one time-line along which the models have evolved. Our focus is on an elegant description and construction of this model evolution in the context of information system design.

A generic version-management framework is in our opinion a suitable (and efficient) implementation platform for the version-management system discussed in this article. In doing so, we would obtain things like alternative versions and alternative courses of model evolutions for free. To do this we should regard a model evolution (like that depicted in Fig. 1) as a (possible) configuration in the framework of Katz [38]. Alternative evolution processes then simply correspond to alternative configurations. The framework described in [38] also provides ways to minimise storage of alternative evolutions by allowing for dynamic construction of configurations.

4. Evolution in the CDM kernel

We have now discussed the two main components of the CDM Evolver: the *corpus evolutionis* and the *version manager*. All that remains to be done is to provide the glue to interconnect these two

components. We do this by introducing a universe of data schemas. This universe provides the boundaries of schema evolution. The elements that make up this universe then provide the element versions needed to make the connection to the evolution mechanism. From this we can then also derive what a schema version within the universe is.

4.1. CDM schema universe

The universe of data schemas is built from the extra temporal versions of the schema components as identified in Section 2. The notion of a universe in which data schemas evolve has been used before, in, for example, [46,47].

Formally, the information structure universe is provided by:

$$\mathcal{IS} = \langle \mathcal{RL}, \mathcal{OB}, \mathcal{VL}, \mathcal{RO}, \mathcal{NO}, \text{SubOf}, \text{Roles}, \text{Player}, \text{Cluster}, \\ \text{CFlavour}, \text{RoleLim}, \text{Ident}, \text{Edge} \rangle$$

The base sets, like \mathcal{RL} and \mathcal{OB} , provide the set of all possible relationship types and object types, respectively. The functions defining the fabric of the information structure, like $\text{Player}: \mathcal{RO} \rightarrow$ and $\text{Roles}: \mathcal{RL} \rightarrow \wp(\text{Roles})$, provide the players for any role that may be part of a schema and the set of roles that may ever be associated to any relationship type, respectively.

Similarly, a complete data schema universe over a set of concrete domains \mathcal{DO} , is identified by the following (extra-temporal) components:

$$\mathcal{CS} \triangleq \langle \mathcal{IS}, \mathcal{CN}, \mathcal{OP}, \sigma, \text{Ops}, \text{DerRule}, \text{UpdRule}, \text{Encap}, \text{Dom} \rangle$$

We presume that the atomic base sets of the data schema universe: $\mathcal{RL}, \mathcal{OB}, \mathcal{VL}, \mathcal{RO}, \mathcal{NO}, \mathcal{CN}, \mathcal{OP}$ are exclusive. Previously we stated that they are exclusive for each single version, but we now require this to be the case for the entire universe as a whole. We can now formally glue the CDM Kernel to the evolution mechanism by stating:

$$\begin{aligned} \mathcal{ELV} &\triangleq \mathcal{TP} \cup \mathcal{RO} \cup \mathcal{NO} \cup \mathcal{CN} \cup \mathcal{OP} \\ \mathcal{ELC} &\triangleq \{ \mathcal{TP}, \mathcal{RO}, \mathcal{NO}, \mathcal{CN}, \mathcal{OP} \} \\ \text{VClass}(e) &\triangleq \mathbb{X} \in \{ \mathcal{TP}, \mathcal{RO}, \mathcal{NO}, \mathcal{CN}, \mathcal{OP} \} \text{ such that } e \in \mathbb{X} \end{aligned}$$

Note that we have combined all types into one single set (\mathcal{TP}). This means that, for instance, a relationship type is allowed to evolve into an object type, but not into a constraint.

From the above definition it may seem to follow that the player of a role may not change in the course of time, or that an object type may not change its subtypes. However, the object types in the universe are simply treated as element *versions*. For example, changing a cluster of an object type, or change a subtyping, requires the involved object types *versions* to be replaced by ‘fresh’ object type *versions*. As shown above, and before, in, for example, [47] and [46], this way of describing allows for elegant definitions. From an implementational point of view, one could regard an object type version from \mathcal{OB} as the combination of its properties: name, supertypes, clustered types, etc. Whenever one of these properties changes, we would automatically have a ‘fresh’ object type. A possible class definition in an OO style would be:

```

VersionElement=CLASS
Name          STRING
END CLASS;
Type=CLASS (VersionElement)
  DerRule:    DerivationRules (OPTIONAL);
  UpdRule:    UpdateRules (OPTIONAL);
END CLASS;
ObjectType=CLASS (Type)
  SuperTypes: SET OF ObjectType;
  Encap:      SET OF (Operation UNION Type);
  Ident:      LIST OF RolePair;
  Cluster:    FUNCTION Natno TO SET OF Type;
  Flavour:    FUNCTION Type TO Flavours;
END CLASS;

```

4.2. CDM schema versions

Now that we have glued the evolution mechanism to the CDM schema universe, we can derive versions of the schemas for a given schema evolution $\langle R, H \rangle$. Using $\text{Snapshot}_H(t)$ the schema version at t in history H can be derived by separating out the different base sets. For each of the base sets $\mathbb{X} \in \{\mathcal{TP}, \mathcal{RL}, \mathcal{OB}, \mathcal{VL}, \mathcal{RO}, \mathcal{NO}, \mathcal{CN}, \mathcal{OP}\}$ we can define:

$$\mathbb{X}_t \triangleq \mathbb{X} \cap \text{Snapshot}_H(t)$$

The remaining components (all functions and predicates) of a schema version can now simply be derived by limiting the extra-temporal domains to the current versions of the base sets, so:

$$\begin{aligned}
\text{Subof}_t &\triangleq \{ \langle x, y \rangle \in \text{Subof} \mid x, y \in \mathcal{OB}_t \} \\
\text{Roles}_t &\triangleq \{ \langle x, R \rangle \in \text{Roles} \mid x \in \mathcal{RL}_t \cup \mathcal{NO}_t \wedge R \subseteq \mathcal{RO}_t \} \\
\text{Player}_t &\triangleq \{ \langle r, x \rangle \in \text{Player} \mid r \in \mathcal{RO}_t \wedge x \in \mathcal{OB}_t \}
\end{aligned}$$

A direct result of this definition is, what might be expected intuitively, that for each schema component $\mathbb{X} \in \{\mathcal{RL}, \mathcal{OB}, \mathcal{VL}, \mathcal{RO}, \mathcal{NO}, \text{SubOf}, \text{Roles}, \text{Player}, \text{Cluster}, \text{CFlavour}, \text{RoleLim}, \text{Ident}, \text{Edge}, \mathcal{CN}, \mathcal{OP}, \sigma, \text{Ops}, \text{DerRule}, \text{UpdRule}, \text{Encap}, \text{Dom}\}$ we have:

$$\mathbb{X}_t \subseteq \mathbb{X}$$

One can quite easily see that it is not hard to extend this framework when taking additional aspects of conceptual schemas into consideration, for example, verbalizations (names) of the object types and relationship types in a schema. In this article we have not elaborated on these issues; however, during the evolution of a conceptual schema there is also the need to change the names of object types or the

verbalisation of fact types. For a detailed study of fact verbalisation and extensions needed for the formal model of a conceptual schema, refer to, for example, [32,33].

4.3. Well-formedness of evolution

Finally, besides well-formedness rules on versions, one might want to formulate rules that limit the evolution of a data schema, i.e. transition-oriented constraints. Specifying such rules, however, can be a rather arbitrary process which depends on the kind of evolution one has in mind. In [47] and [46] some example rules have been formulated in the context of schema evolution due to evolution of the universe of discourse. However, for evolution during the design phase of an information system, one would typically like to be more liberal.

Although we do not formulate such rules explicitly here, one might imagine having a set of extra rules EW (Evolution Well-formedness) that govern the evolution of data schemas. Then we can now refine the *IsModEvol* predicate to the CDM Kernel specific definition:

$$\begin{aligned} \text{IsCDMEvol}(R,H) &\Leftrightarrow \text{IsModEvol}(R,H) \wedge \forall_{t \in \mathcal{T}, \mathcal{I}} [\text{IsSch}(\mathcal{CS}_t)] \\ &\wedge \langle R,H \rangle \text{ obeys the EW rules} \end{aligned}$$

5. Case study

In this section we provide a case study as an illustration of the use of the versioning mechanism presented. We will study the effects of a sequence of changes to a part of the martial arts domain. These changes could be part of the design phase or evolution phase of the system. This does, however, have no influence on our case study.

We start with the schema presented in Fig. 12. To describe the evolution of this schema, we will use ten element identifiers for types: t_1, \dots, t_{10} , and ten element identifiers for constraints: c_1, \dots, c_{10} .

Let s_1 be the first point in time at which a change (a reaction) takes place. At this point in time, the schema as shown in Fig. 12 is current, and we would have for the model history $\langle R,H \rangle$:

$$\begin{aligned} H(t_1)(s_1) &= 1:\text{Player} \\ H(t_2)(s_1) &= 1:\text{Rank} \\ H(t_3)(s_1) &= 1:\text{in judo has} \\ H(c_1)(s_1) &= 1:\text{EACH Player in judo has SOME Rank} \\ H(c_2)(s_1) &= 1:\text{EACH Player in judo has A UNIQUE Rank} \\ H(c_3)(s_1) &= 1:\text{FOR EACH Rank THERE IS A Player WHO in judo has THAT Rank} \end{aligned}$$

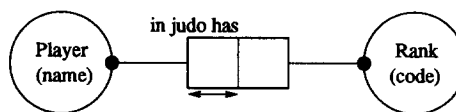


Fig. 12. Martial Arts — Step 1.

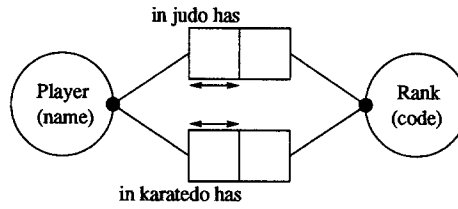


Fig. 13. Martial Arts — Step 2.

The element versions for the object types are identified by referring to the schema number and the name within the schema. For example: 1: Rank. In a real implementation, this identifier will of course have to be some internally generated identifier. The three constraints in this ORM model have been verbalised textually in the style of a conceptual language like FORML [26,29] or LISA-D [31].

The first step in the evolution of the model will take us to the schema as shown in Fig. 13. It was found that besides playing judo, players also played karatedo, requiring us to introduce an extra relationship type. If $s_2 = \triangleright s_1$, then we could model this new version as:

$H(t_1)(s_2)$	=	1:Player
$H(t_2)(s_2)$	=	1:Rank
$H(t_3)(s_2)$	=	1:in judo has
$H(t_4)(s_2)$	=	2:in karatedo has
$H(c_1)(s_2)$	=	2:EACH Player EITHER in judo has SOME Rank OR in karatedo has SOME Rank
$H(c_2)(s_2)$	=	1:EACH Player in judo has A UNIQUE Rank
$H(c_3)(s_2)$	=	2:FOR EACH Rank THERE IS A Player WHO EITHER in judo has THAT Rank OR in karatedo has THAT Rank
$H(c_4)(s_2)$	=	2:EACH Player in karatedo has A UNIQUE Rank

All that has happened is that a new relationship type t_4 and constraint c_4 have been introduced, while the existing constraints c_1 and c_3 are modified to incorporate the new relationship type.

The reactions underlying this evolution step can be modelled as:

$$\begin{aligned}
 r_1 &= \langle \{t_3, t_4, c_1, c_3, c_4\}, \mathcal{TP}, s_1 \rangle \\
 \text{Reason}(r_1) &= \text{'Introduction of new martial art karatedo'}
 \end{aligned}$$

This reaction took place at s_1 , and is viewed by us as being of class type (\mathcal{TP}) as it is due to the introduction of a new martial art. It is easy to verify that we now have:

$$\begin{aligned}
 \text{Catalyst}_H(r_1) &= \{t_3\} \\
 \text{Fusion}_H(r_1) &= \emptyset \\
 \text{Fission}_H(r_1) &= \{t_4\} \\
 \text{Absorption}_H(r_1) &= \{c_1, c_3\} \\
 \text{Emission}_H(r_1) &= \{c_1, c_3, c_4\}
 \end{aligned}$$

Please note that c_1 and c_3 both act as input and output to the reaction, but that they are not catalysts as they are *changed* by the reaction.

The only element to remain the same is t_3 , representing the existing *in judo has* relationship type. We could have chosen not to include this type in the reaction. However, we have the *freedom* to do so, and in this case we have chosen to include it as we view the *in judo has* relationship as a cause for later introducing the *in karatedo has* relationship. This is, however, entirely up to the subjective considerations of the modeller.

The next evolution step is brought about when we realise that **Ranks** in martial arts are not simply expressed as codes, but rather as a combination of a rank kind and a level. The rank kinds used are the ‘kyu’ and ‘dan’. So a karatedo could have as rank ‘7th dan’. This evolution step leads to the schema provided in Fig. 14.

The effects on the history are (presuming $s_3 = \triangleright s_2$):

$H(t_1)(s_3)$	=	1:Player
$H(t_2)(s_3)$	=	3:Rank
$H(t_3)(s_3)$	=	3:in judo has
$H(t_4)(s_3)$	=	3:in karatedo has
$H(t_5)(s_3)$	=	3:is of
$H(t_6)(s_3)$	=	3:is at
$H(t_7)(s_3)$	=	3:RankKind
$H(t_8)(s_3)$	=	3:Level
$H(c_1)(s_3)$	=	3:EACH Player EITHER in judo has SOME Rank OR in karatedo has SOME Rank
$H(c_2)(s_3)$	=	3:EACH Player in judo has A UNIQUE Rank
$H(c_3)(s_3)$	=	3:FOR EACH Rank THERE IS A Player WHO EITHER in judo has THAT Rank OR in karatedo has THAT Rank
$H(c_4)(s_3)$	=	3:EACH Player in karatedo has A UNIQUE Rank
$H(c_5)(s_3)$	=	3:EACH Rank is of SOME RankKind
$H(c_6)(s_3)$	=	3:EACH Rank is at SOME Level
$H(c_7)(s_3)$	=	3:EACH Rank is of A UNIQUE RankKind
$H(c_8)(s_3)$	=	3:EACH Rank is at A UNIQUE Level

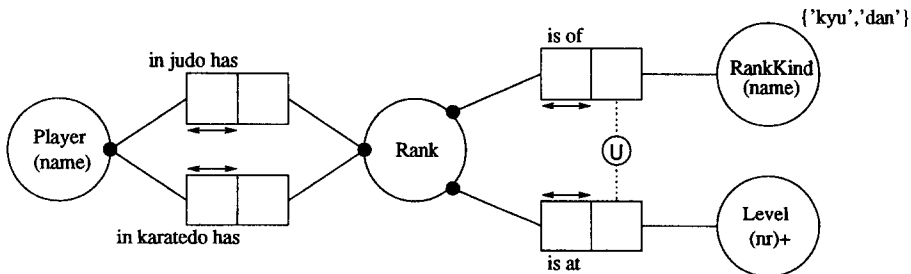


Fig. 14. Martial Arts — Step 3.

$H(c_9)(s_3) = 3:\text{EACH Rank IS UNIQUELY DETERMINED BY THE RankKind IT is of AND THE Level IT is at}$

$H(c_{10})(s_3) = 3:\text{A RankKind IS ONE OF 'kyu', 'dan'}$

As the **Rank** object type changes its identification, all relationship types in which it is involved implicitly change, as well as constraints defined over these relationship types. So, although t_3 , t_4 , c_1 , c_2 , c_3 and c_4 may *appear* to be the same in Figs. 13 and 14, they are in fact different, as an underlying object type (**Rank**) has been changed.

This reaction took place at s_2 , and is again viewed as being of class type (\mathcal{TP}):

$r_2 = \langle \{t_2, \dots, t_8, c_1, \dots, c_{10}\}, \text{TP}_{s_2} \rangle$
Reason(r_2) = 'Introduction of new identification scheme for ranks'

For this reaction we have:

Catalyst $H(r_1) = \emptyset$
Fusion $H(r_1) = \{t_2, t_3, t_4\}$
Fission $H(r_1) = \{t_1, \dots, t_8\}$
Absorption $H(r_1) = \{c_1, \dots, c_4\}$
Emission $H(r_1) = \{c_1, \dots, c_{10}\}$

As the final evolution step in our case study, let us assume that we want to be able to store the players of any arbitrary martial art. This requires us to modify the two binary relationships for judo and karatedo to the more general situation as depicted in Fig. 15. There we have introduced a ternary relationship that allows us to record arbitrary martial arts for players, together with their ranks.

Let $s_4 = \triangleright s_3$, then we can model the effects on the history as:

$H(t_1)(s_4) = 1:\text{Player}$
 $H(t_2)(s_4) = 3:\text{Rank}$
 $H(t_3)(s_4) = 3:\text{is of}$
 $H(t_6)(s_4) = 3:\text{is at}$

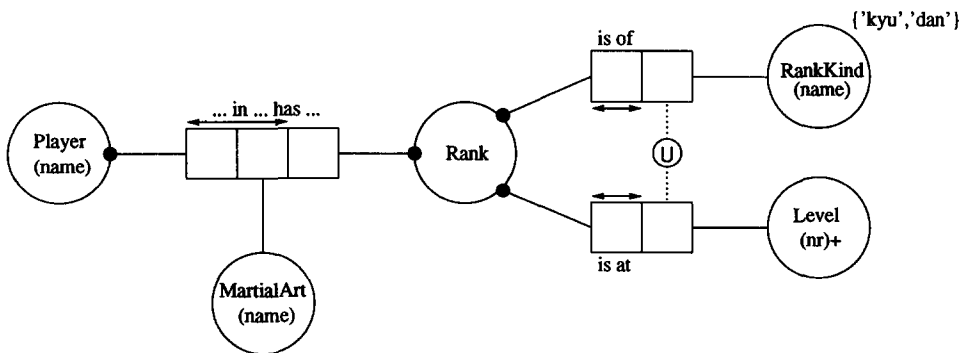


Fig. 15. Martial Arts — Step 4.

$H(t_7)(s_4)$	=	3:RankKind
$H(t_8)(s_4)$	=	3:Level
$H(t_9)(s_4)$	=	4:... in... has...
$H(t_{10})(s_4)$	=	4:MartialArt
$H(c_1)(s_4)$	=	4:EACH Player in SOME MartialArt has SOME Rank
$H(c_2)(s_4)$	=	4:FOR EACH Player AND MartialArt: THAT Player in THAT MartialArt has A UNIQUE Rank
$H(c_3)(s_3)$	=	3:FOR EACH Rank THERE IS A Player WHO in SOME MartialArt has THAT Rank
$H(c_5)(s_4)$	=	3:EACH Rank is of SOME RankKind
$H(c_6)(s_4)$	=	3:EACH Rank is at SOME Level
$H(c_7)(s_4)$	=	3:EACH Rank is of A UNIQUE RankKind
$H(c_8)(s_4)$	=	3:EACH Rank is at A UNIQUE Level
$H(c_9)(s_4)$	=	3:EACH Rank IS UNIQUELY DETERMINED BY THE RankKind IT is of AND THE Level IT is at
$H(c_{10})(s_4)$	=	3:A RankKind IS ONE OF 'kyu', 'dan'

The reaction itself can be modelled as:

$$\begin{aligned} r_3 &= \langle \{t_3, t_4, t_9, t_{10}, c_2, c_4\}, \mathcal{TP}, s_3 \rangle \\ \text{Reason}(r_3) &= \text{'Generalisation to any martial art'} \end{aligned}$$

This leads to:

$$\begin{aligned} \text{Catalyst}H(r_1) &= \emptyset \\ \text{Fusion}H(r_1) &= \{t_3, t_4\} \\ \text{Fission}H(r_1) &= \{t_9, t_{10}\} \\ \text{Absorption}H(r_1) &= \{c_2, c_4\} \\ \text{Emission}H(r_1) &= \{c_2\} \end{aligned}$$

This concludes our small case study. What we hope to have illustrated is that using the described versioning system, we can describe in more detail the evolution of a conceptual schema, be it during the analyses, design or maintenance phase of the information system.

6. Conclusions

In this article we have defined a version-management system for schema evolution: the CDM Evolver. It can be seen as the finishing touch to the existing CDM Kernel. This version-management system allows us to model the evolution of data schemas in the CDM Kernel. Two important features are the ability to describe the evolution of data schemas through the entire modelling process including the internal representation, and the fact that it can model the interaction between model

elements in the course of time. This latter property is believed to be crucial to provide a better insight into schema design processes, both for modellers themselves as well as their managers.

The version-management system has been setup such that it is extendible and easy to adapt to other models. This means that as such it is not only useful in the context of the CDM Kernel, but can be applied to other modelling techniques as well.

Acknowledgments

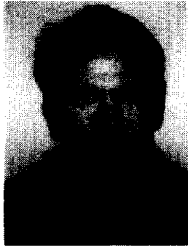
I would like to thank A.H.M. ter Hofstede for his valuable comments on the first drafts of this article.

References

- [1] G.H.W.M. Bronts, S.J. Brouwer, C.L.J. Martens and H.A. Proper, A unifying object role modelling approach, *Information Systems* 20(3) (1995) 213–235.
- [2] F.L. Bauer, M. Broy, H. Partsch, P. Pepper and H. Wössner, Systematics of transformation rules, in F.L. Bauer and M. Broy (eds.), *Program Construction, Vol. 69, Lecture Notes in Computer Science* (Springer-Verlag, Berlin, 1979) pp. 273–289.
- [3] C. Batini, S. Ceri and S.B. Navathe, *Conceptual Database Design — An Entity-Relationship Approach* (Benjamin Cummings, Redwood City, CA, 1992).
- [4] S. Berman, A semantic data model as the basis for an automated database design tool, *Information Systems* 11(2) (1986) 149–165.
- [5] P. van Bommel, A.H.M. ter Hofstede and Th.P. van der Weide, Semantics and verification of object-role models, *Information Systems* 16(5) (October 1991) 471–495.
- [6] J. Banerjee, W. Kim, H.J. Kim and H.F. Korth, Semantics and implementation of schema evolution in object-oriented databases, *SIGMOD Record* 16(3) (December 1987) 311–322.
- [7] F.L. Bauer, B. Möller, H. Partsch and P. Pepper, Formal program construction by transformations — computer-aided, intuition-guided programming, *IEEE Trans. on Software Engineering* 15(2) (February 1989) 165–180.
- [8] P. van Bommel, Implementation selection for object-role models, in T.A. Halpin and R. Meersman (eds.), *Proc. 1st Intern. Conf. on Object-Role Modelling (ORM-1)*, Magnetic Island, Australia (July 1994) pp. 103–112.
- [9] P. van Bommel, Database optimization: An Evolutionary approach, *Ph.D. Thesis* (University of Nijmegen, Nijmegen, The Netherlands, 1995).
- [10] P. van Bommel and Th.P. van der Weide, Reducing the search space for conceptual schema transformation, *Data & Knowledge Engineering* 8 (1992) 269–292.
- [11] R. Chiang, T. Barron and V. Storey, Reverse engineering of relational databases: Extraction of an eer model from a relational database, *Data & Knowledge Engineering* 12(2) (1994) 107–142.
- [12] P.N. Creasy and W. Hesse, Two-level NIAM: A way to get it object-oriented, in T.W. Olle and A.A. Verrijn-Stuart (eds.), *Proc. of the IFIP WG 8.1 CRIS-94 Conf. on Methods and Associated Tools for the Information Life Cycle*, Maastricht, The Netherlands (Sept. 1994) pp. 209–221.
- [13] L.J. Campbell, T.A. Halpin and H.A. Proper, Conceptual schemas with abstractions — making flat conceptual schemas more comprehensible, *Data & Knowledge Engineering* 20(1) (1996) 39–85.
- [14] C.R. Carlson, W. Ji and A.K. Arora, The nested entity-relationship model, in F.H. Lochovsky (ed.), *Proc. 8th Intern. Conf. on Entity-Relationship Approach, Entity-Relationship Approach to Database Design and Querying* (Elsevier Science Publishers, Toronto, Canada, 1990) pp. 43–57.

- [15] B. Curtis, H. Krasner and N. Iscoe, A field study of the software design process for large systems, *Communications of the ACM* 31(11) (November 1988) 1268–1287.
- [16] B. Curtis, H. Krasner, V. Shen and N. Iscoe, On building software process models under the lamppost, in *Proc. 9th Intern. Conf. on Software Engineering* (IEEE Computer Society Press, Monterey, CA, 1987) pp. 96–103.
- [17] P.N. Creasy and H.A. Proper, A generic model for 3-dimensional conceptual modelling, *Technical Report 342* (Department of Computer Science, University of Queensland, Australia, July 1995); electronically available at: <http://www.icis.qut.edu.au/~erikp/articles/ConDaMo.ps>; *Data & Knowledge Engineering* 20(2) (1996) 119–162.
- [18] O.M.F. De Troyer, The OO-binary relationship model: A truly object-oriented conceptual model, in R. Andersen, J.A. Bubenko and A. Sølvberg (eds.), *Proc. 3rd Intern. Conf. CAiSE'91 on Advanced Information Systems Engineering, Vol. 498 of Lecture Notes in Computer Science* (Springer-Verlag, Trondheim, Norway, May 1991) pp. 561–578.
- [19] O.M.F. De Troyer, On data schema transformations, *Ph.D. Thesis* (University of Tilburg (K.U.B.), Tilburg, The Netherlands, 1993).
- [20] O.M.F. De Troyer and R. Janssen, On modularity for conceptual data models and the consequences for subtyping, inheritance & overriding, in E.K. Elmagarmid and E.J. Neuhold (eds.), *Proc. 9th IEEE Conf. on Data Engineering (ICDE '93)* (IEEE Computer Society Press, 1993) pp. 678–685.
- [21] R. Elmasri and S.B. Navathe, *Fundamentals of Database Systems*, 2nd Edition (Benjamin Cummings, Redwood City, CA, 1994).
- [22] M.M. Fonkam and W.A. Gray, An Approach to Eliciting the Semantic of Relational Databases, in P. Loucopoulos, editor, *Proc. 4th Intern. Conf. CAiSE'92 on Advanced Information Systems Engineering, Vol. 593, Lecture Notes in Computer Science* (Springer-Verlag, San Mateo, CA, 1991).
- [23] J.-L. Hainaut, Entity-generating schema transformation for entity-relationship models, in *Proc. 10th Intern. Conf. on the Entity-Relationship Approach, Lecture Notes in Computer Science* (Springer-Verlag, San Mateo, CA, 1991).
- [24] T.A. Halpin, Conceptual schema optimization, *Australian Computer Science Communications* 12(1) (1990) 136–145.
- [25] T.A. Halpin, A fact-oriented approach to schema transformation, in B. Thalheim, J. Demetrovics and H.-D. Gerhardt (eds.), *MFDDBS '91, Vol. 495, Lecture Notes in Computer Science* (Springer-Verlag, Rostock, Germany, 1991) pp. 342–356.
- [26] T.A. Halpin, *Conceptual Schema and Relational Database Design*, 2nd edition (Prentice-Hall, Sydney, Australia, 1995).
- [27] J.-L. Hainaut, M. Chadelon, C. Tonneau and M. Joris, Contribution to a theory of database reverse engineering, in *Proc. IEEE Working Conf. on Reverse Engineering* (IEEE Computer Society Press, Baltimore, MD, May 1993).
- [28] J.-L. Hainaut, V. Englebert, J. Henrard, J.-M. Hick and D. Roland, Database evolution: The DB-MAIN approach, in P. Loucopoulos (ed.), *Proc. 13th Intern. Conf. on the Entity-Relationship Approach, Vol. 881, Lecture Notes in Computer Science* (Springer-Verlag, Manchester, United Kingdom, December 1994) pp. 112–131.
- [29] T.A. Halpin and J. Harding, Automated support for verbalization of conceptual schemas, in S. Brinkkemper and F. Harmsen (eds.), *Proc. 4th Workshop on the Next Generation of CASE Tools*, Paris, France (June 1993) pp. 151–161.
- [30] T.A. Halpin and H.A. Proper, Database schema transformation and optimization, in M.P. Papazoglou (ed.), *Proc. of the OOER'95, 14th Intern. Object-Oriented and Entity-Relationship Modelling Conf., Vol. 1021, Lecture Notes in Computer Science* (Springer-Verlag, Gold Coast, Australia, December 1995) pp. 191–203.
- [31] A.H.M. ter Hofstede, H.A. Proper and Th.P. van der Weide, Formal definition of a conceptual language for the description and manipulation of information models, *Information Systems* 18(7) (October 1993) 489–523.
- [32] A.H.M. ter Hofstede, H.A. Proper and Th.P. van der Weide, Grammar based information modelling, *Technical Report CSI-R9414* (Computing Science Institute, University of Nijmegen, Nijmegen, The Netherlands, October 1994); electronically available at: <http://www.icis.qut.edu.au/~erikp/articles/Grammar.ps>.
- [33] A.H.M. ter Hofstede, H.A. Proper and Th.P. van der Weide, Exploring fact verbalizations for conceptual query formulation, in *Proc. 2nd Intern. Workshop on Applications of Natural Language to Databases (NLDB'96)* (IOS Press, Amsterdam, The Netherlands, June 1996) pp. 40–51.
- [34] J.-L. Hainaut, C. Tonneau, M. Joris and M. Chadelon, Schema transformation techniques for database reverse engineering, in *Proc. 12th Intern. Conf. on the Entity-Relationship Approach, Lecture Notes in Computer Science* (Springer-Verlag, Dallas, TX, December 1993).
- [35] J.-L. Hainaut, C. Tonneau, M. Joris and M. Chadelon, Transformation-based database reverse engineering, in *Proc.*

- 12th Intern. Conf. on the Entity-Relationship Approach, Lecture Notes in Computer Science* (Springer-Verlag, Dallas, TX, December 1993).
- [36] A.H.M. ter Hofstede and Th.P. van der Weide, Expressiveness in conceptual data modelling, *Data & Knowledge Engineering* 10(1) (February 1993) 65–100.
 - [37] K. Kalman, Implementation and critique of an algorithm which maps a relational database to a conceptual model, in R. Andersen, J.A. Bubenko and A. Sølvberg (eds.), *Proc. 3rd Intern. Conf. CAiSE'91 on Advanced Information Systems Engineering, Vol. 498, Lecture Notes in Computer Science* (Springer-Verlag, Trondheim, Norway, May 1991) 393–415.
 - [38] R.H. Katz, Toward a unified framework for version modelling in engineering databases, *ACM Computer Surveys* 22(4) (1990) 375–408.
 - [39] I. Kobayashi, Classification and transformations of binary relationship relation schemata, *Information Systems* 11(2) (1986) 109–122.
 - [40] G. Kristen, *Object Orientation, the KISS Method: From Information Architecture to Information System* (Addison-Wesley, Reading, MA, 1994).
 - [41] G. Kappel and M. Schrefl, Local referential integrity, in G. Pernul and A.M. Tjoa (eds.), *11th Intern. Conf. on the Entity-Relationship Approach, Vol. 645, Lecture Notes in Computer Science* (Springer-Verlag, Karlsruhe, Germany, October 1992) pp. 41–61.
 - [42] J.I. McCormack, T.A. Halpin and P.R. Ritson, Automated mapping of conceptual schemas, in C. Rolland, F. Bodart and C. Cauvet (eds.), *Proc. 5th Intern. Conf. CAiSE'93 on Advanced Information Systems Engineering, Vol. 685, Lecture Notes in Computer Science* (Springer-Verlag, Paris, France, 1993) pp. 432–448.
 - [43] J.L.H. Oei, H.A. Proper and E.D. Falkenberg, Evolving information systems: Meeting the ever-changing environment, *Information Systems Journal* 4(3) (1994) 213–233.
 - [44] H.A. Proper and T.A. Halpin, Conceptual schema optimisation — database optimisation before sliding down the waterfall, *Technical Report 341* (Department of Computer Science, University of Queensland, Australia, July 1995); electronically available at: <http://www.icis.qut.edu.au/~erikp/articles/Trans2.ps>.
 - [45] H. Partsch and R. Steinbrüggen, Program transformation systems, *Computing Surveys* 15(3) (1983).
 - [46] H.A. Proper and Th.P. van der Weide, EVORM: A conceptual modelling technique for evolving application domains, *Data & Knowledge Engineering* 12 (1994) 313–359.
 - [47] H.A. Proper and Th.P. van der Weide, A general theory for the evolution of application models, *IEEE Trans. on Knowledge and Data Engineering*, 7(6) (December 1995) 984–996.
 - [48] J. Rumbaugh, M. Blaha, W. Premerlani, F. Eddy and W. Lorensen, *Object-Oriented Modeling and Design* (Prentice-Hall, Englewood Cliffs, NJ, 1991).
 - [49] N. Rishe, A methodology and tool for top-down relational database design, *Data & Knowledge Engineering* 10 (1993) 259–291.
 - [50] P.R. Ritson, Use of conceptual schemas for a relational implementation, *Ph.D. Thesis* (University of Queensland, Brisbane, Australia, 1994).
 - [51] P. Shoval and M. Even-Chaime, ADDS: A system for automatic database schema design based on the binary-relationship model, *Data & Knowledge Engineering*, 2(2) (1987) 123–144.
 - [52] P. Shoval, Essential information structure diagrams and database schema design, *Information Systems*, 10(4) (1985) 417–423.
 - [53] P. Shoval and N. Shreiber, Database reverse engineering: From the relational to the binary relationship model, *Data & Knowledge Engineering* 10 (1993) 293–315.
 - [54] M.T. Tresch and M.H. Scholl, Meta object management and its application to database evolution, in G. Pernul and A.M. Tjoa (eds.), *11th Intern. Conf. on the Entity-Relationship Approach, Vol. 645, Lecture Notes in Computer Science* (Springer-Verlag, Karlsruhe, Germany, October 1992) pp. 299–321.
 - [55] T.J. Teorey, G. Wei, D.L. Bolton and J.A. Koenig, ER model clustering as an aid for user communication and documentation in database design, *Communications of the ACM*, 32(8) (August 1989) 975–987.
 - [56] D. Vermeir, Semantic hierarchies and abstractions in conceptual schemata, *Information Systems*, 8(2) (1983) 117–124.



Dr. H.A. (Erik) Proper is currently a lecturer at the School of Information Systems in the Queensland University of Technology, Brisbane, Australia. He is a member of the Cooperative Information Systems Research Centre from that University. He is also a member of the Distributed Systems Technology Centre (DSTC), one of the Cooperative Research Centres funded by the Australian Government, Australian Universities, and a number of multinational companies.

Dr. Proper received his Master's degree from the University of Nijmegen, The Netherlands in May of 1990, and received his Ph.D. from the same University in April 1994. In his doctoral thesis he developed a theory for conceptual modelling of evolving application domains, yielding a formal specification of evolving information systems. From May 1994 to September 1995 he worked as a Research Fellow at the Computer Science Department of the University of Queensland, Brisbane, Australia. During that period he also conducted research in the Asymetrix Research Lab at that University for Asymetrix Corp., Bellevue, Washington.

Dr. Proper has co-authored several journal papers and conference publications. Currently, his main research interests include information retrieval, knowledge discovery, knowledge based systems, conceptual query languages and conceptual modelling. He has particular interests in the application of logic and linguistics to the above areas.