

# Conceptual Schemas with Abstractions –

Making flat conceptual schemas more comprehensible

L.J. Campbell<sup>1</sup> and T.A. Halpin<sup>2</sup> and H.A. Proper<sup>3</sup>

Version of August 22, 1998 at 22:48

PUBLISHED AS:

*Conceptual Schemas with Abstractions – Making flat conceptual schemas more comprehensible*, L.J. Campbell and T.A. Halpin and H.A. Proper, *Data & Knowledge Engineering*, 20(1), 39–85, 1996

## Abstract

Flat graphical, conceptual modeling techniques are widely accepted as visually effective ways in which to specify and communicate the conceptual data requirements of an information system. Conceptual schema diagrams provide modelers with a picture of the salient structures underlying the modeled universe of discourse, in a form that can readily be understood by and communicated to users, programmers and managers. When complexity and size of applications increase, however, the success of these techniques in terms of comprehensibility and communicability deteriorates rapidly.

This paper proposes a method to offset this deterioration, by adding abstraction layers to flat conceptual schemas. We present an algorithm to recursively derive higher levels of abstraction from a given (flat) conceptual schema. The driving force of this algorithm is a hierarchy of conceptual importance among the elements of the universe of discourse.

KEYWORDS:

Conceptual Data Modelling, Schema Abstraction, ORM, ER, NIAM

## 1 Introduction

Conceptual schemas play an important, and recognized role in the development life cycle of an information system [28]. They serve both as a means by which the salient structures of the underlying universe of discourse (UoD) can be captured, and as a communication tool among the designers, programmers, users and managers [32]. Conceptual schema modeling techniques, such as Entity Relationship (ER) modeling [1], [12] and Object Role Modeling (ORM) [17] are widely acknowledged as being visually effective ways in which to specify and communicate the conceptual data requirements of an information system.

However, as database application requirements increase in size and complexity, the comprehensibility and maintainability of the specification degrades rapidly [28]. Simson identified the problem of representing large data models as ‘one of the most serious limitations of data modeling in practice’ [30]. It is claimed in Feldman and Miller [13] that the ‘usefulness of any diagram is inversely proportional to the size of the model depicted’. This problem, which has been referred to as the *Database Comprehension Problem* [8], is shared by all *flat* data models. In the specification of a *flat* conceptual schema, each object type is viewed at only one level of abstraction in a single diagram and all object types are considered to be of equal importance within the application [7]. While this is satisfactory for small, academic examples, when

<sup>1</sup>Department of Computer Science, University of Queensland, Australia 4072

<sup>2</sup>Asymetrix Corporation, Bellevue WA, USA

<sup>3</sup>Cooperative Information Systems Research Centre, Faculty of Information Technology, Queensland University of Technology, GPO Box 2434, Brisbane, 4001 Australia, E.Proper@acm.org

conceptual schemas of a moderate to large size are involved, this feature reduces the rate of comprehension of the application.

The Database Comprehension Problem in flat data models has motivated several authors to try to find successful methods by which to form abstractions on a given flat conceptual schema [8], [10], [13], [14], [15], [26], [27], [28], [29], [31], [32], [33]. In 1983, Vermeir [32] described a number of abstraction techniques including *viewpoint relative abstraction*, which displays the portion of the schema within a particular distance of a focal object type, and the *absolute abstraction hierarchy*, which iteratively removes non-key concepts from successive layers of abstraction.

An abstraction technique which is quite popular in the literature is *Entity-Relationship Model Clustering*. An *ER Model Cluster Diagram* is ‘a hierarchy of successively more detailed Entity Relationship diagrams, with a lower-level diagram appearing as a single entity type on the next higher level diagram’ [13]. Martin developed an ER clustering procedure based on 1:m relationship between entities, which is simple, but rather arbitrary and judgmental [27]. ER model clustering was then applied to the Whitbread Corporate Data Architecture in June 1983 in order to test the theory on a significant practical application. As an outcome, in 1986, Feldman and Miller proposed a semi-algorithmic approach to ER model clustering - one which still relies heavily on human direction and judgment.

In 1989, Carlson and Ji [8] proposed the Nested Entity Relationship (NER) model as an extension to the multi-level ER clustering techniques of Feldman and Miller. NER supports traditional abstraction techniques such as aggregation, generalization and association, as well as allowing ER diagrams at one level to be abstracted into either complex entities or complex relationships at the next higher level. In that same year, Teorey et. al. [31] introduced a set of ER model clustering rules. In this publication, entities were grouped recursively, based on a list of grouping operations prioritized according to the *cohesion* (or internal strength) among the entities involved. Once again, the algorithm proposed is largely based on arbitrary human judgment. This work was taken further by Huffman and Zoeller in 1989, who confirmed the feasibility of using a rule-based system to automate the ER clustering process of Teorey et. al [26].

A number of other abstraction techniques have been introduced in the years since. In 1991, Czejdo and Embley proposed the management of large complex data models using views and a number of functions to manipulate those views [10]. In [28] Moody proposed a representation scheme for abstraction based on the organization of a street directory, using various levels of detail and intermap references and overlap between scopes. A new abstraction mechanism for typed graphs, called *Caves*, which allows the designer to selectively ‘amplify or diminish’ parts of the conceptual model, was presented by Walko in 1992 [33]. The techniques available through *Caves* include: *Filtration*, which removes extraneous details and constraints; *Perspective* which presents only the local vicinity of a selected focal point; and *Comprehension*, which folds the schema into a smaller version.

The common element of many of the abstraction techniques throughout the literature is the selection of a set of important elements within a conceptual schema. Many different names have been given to these objects which are considered to be of importance within the application domain, including *key concepts* [32], *major entity types* [13], *maximal objects* [8] and *dominant objects* [26], [31]. In this paper (as in [6], [7], [5]) we call these the *major object types*. In order to briefly highlight the differences between these similar concepts throughout the literature, we will refer to a small example conceptual schema shown in figure 1 ([17] p. 402) (This diagram uses Object Role Modeling notation which is explained in section 2).

The abstraction techniques described by Vermeir in [32] are based on the notion of a *key concept*. Key concepts are those objects within a conceptual schema which are considered to be of higher semantic importance because they keep the graph connected. In the example in figure 1, therefore, the object types ‘Movie’, ‘Person’, ‘MoneyAmt’ and ‘Country’ would be considered, in [32], to be the key concepts. It is quite apparent to a human, however, that ‘MoneyAmt’ is not one of the most semantically important object types in the example Universe of Discourse. In fact, Vermeir himself observes that the definition of ‘key concept’ is too simplistic, and many cases arise in which the abstractions produced are not intuitive [32].

Feldman and Miller consider the most important entity types to be those that appear in more than one branch at any particular level of their clustering hierarchy and call these objects *major entity types* [13]. Using Feldman and Miller’s algorithm [13], every object type in the schema (Movie, Title, Person, MoneyAmt, Country, Date) would be classified as a candidate to be a *major entity type*, after which human intuition is required to narrow down the choice. Nested Entity Relationship models give a stronger theoretical foundation than previous multi-level ER approaches, though the concept of *maximal objects* [8].

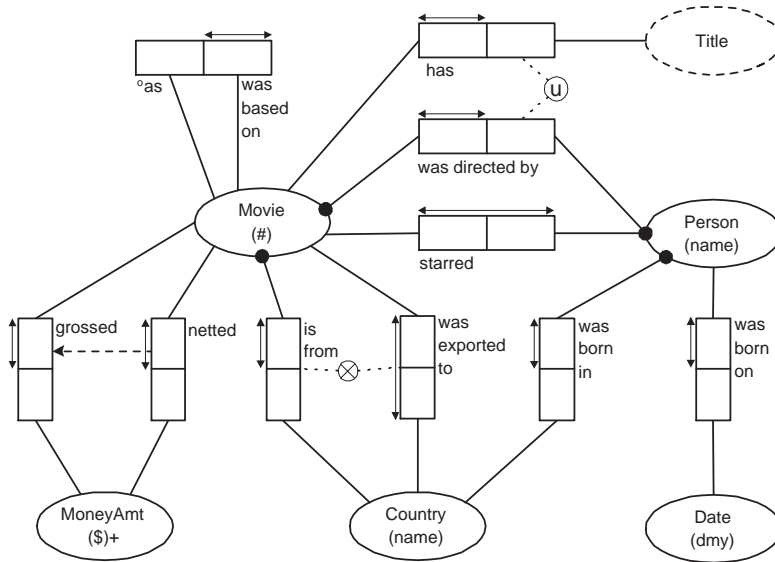


Figure 1: Example Conceptual Schema

However, no attempt is made to automate the selection of *maximal objects* in a conceptual schema. Teorey et al. [31] and Huffman and Zoeller [26] base their ER model clustering rules around what they call the *dominant objects*. The only *dominant object* that results from the automatic algorithm presented in [26] is ‘Movie’. Identifying other dominant objects appears to, once again, necessitate the use of human judgment. While Huffman and Zoeller’s results were loosely comparable to those produced using human intuition, it was acknowledged by the authors that some aspects of the clustering algorithm were too simplistic for complex cases. Referring back to the conceptual schema in figure 1, it is most likely that a human would intuitively decide that the *major object types* in this Universe of Discourse are ‘Movie’ and ‘Person’. None of the algorithms reviewed in the literature arrive at this result automatically.

The first goal of this paper, therefore, is to formalize a method for the strictly automatic selection of major object types. What sets our approach apart from others is that our approach considers the detailed conceptual semantics hidden in the constraints and also the manner in which the facts within the domain are verbalized. In particular, our approach utilizes the detailed constraint specifications and verbalizations provided by Object Role Modeling. It is believed that a lot of the human intuition (conceptual semantics) is contained in these constraints and verbalizations. We, therefore, claim that our approach more accurately imitates human intuition than previous methods. As a second goal, this paper also aims to utilize these selected major object types in an algorithm to derive abstractions for a flat conceptual schema.

We begin, in section 2, by introducing a formal description of Object Role Modeling, which will be used as the foundation of the algorithms presented. Section 3 extends the semantics of Object Role Modeling by introducing the notion of conceptual *anchors*, which are required for the detection of major object types. An automated method for selecting anchors is presented. The selection of anchors is based on the semantics of constraints defined on surrounding relationship types. The semantics of these constraints, in terms of populations, allows us to make this selection. The notion of major object types and abstraction levels is then introduced in section 4, together with a method for automatically determining them. Section section 5 illustrates how this automated abstraction process is performed on a small case study; and conclusions are reached in section six.

## 2 Object Role Modeling

Object Role Modeling (ORM) views the world as a collection of objects which play roles and, unlike Entity-Relationship Modeling, makes no initial use of the *attribute* construct. Every elementary type of

fact which occurs between an object type in the Universe of Discourse (UoD) is verbalized and displayed on a conceptual schema diagram. Object Role Modeling also allows a wide variety of data constraints to be specified on the conceptual schema, including mandatory role, uniqueness, exclusion, equality, subset, and occurrence frequency.

The high level of detail displayed on an ORM diagram allows Object Role Modeling to offer a correspondingly high level of expressiveness. Unfortunately, this high level of detail also tends to promote the degradation in comprehensibility and communicability in large conceptual schemas. An ER diagram, through its use of attributes, can already be thought of as an abstraction (or summary) of a corresponding ORM diagram. In this way, traditional Entity Relationship modeling can postpone the immediate effects of the *Database Comprehension Problem* until a larger Universe of Discourse is required. It is not uncommon in practice, however, for abstractions (or summaries) of ER diagrams themselves to be required. While the scale of the problem, therefore, differs slightly between ORM and ER, the *Database Comprehension Problem*, nonetheless, is universally shared by all flat modeling techniques.

For the purposes of this paper, we plan to consider the more detailed of the two most common data modeling techniques (Object Role Modeling) and introduce a method to control the schema's visual complexity during the information system development. As argued before in [7], [5] and [17], an Entity Relationship model can be considered comparable with the first of the abstraction levels on an ORM model.

The following subsections outline a formalization of some fundamental ORM structures and constraints which will be required in sections 3 and 4 to describe our abstraction methods. The formalization of ORM as presented in this article inherits a rich and well published history, full of constant refinements and additions. The evolution of this particular ORM formalization started out from the PM/PSM version of ORM [2], [20], [23]. More 'modernized' versions of ORM formalizations can be found in [20], [3] and the most recent developments are discussed in [9]. Alternative formalizations can be found in [11] and [16].

While formalizations of ORM have been published before, this paper needs to describe the formalization again in order to be self contained. In this formalization, we limit ourselves to syntactical issues only. Issues regarding the associated semantics can be found in the referenced publications. Furthermore, the formalization presented in this paper is based on a limited number of basic concepts to provide us with only what is needed for the purposes of abstraction. For a detailed description of the methodology associated with Object Role Modeling, refer to [17].

## 2.1 Information Structure

The cornerstone of a conceptual schema is formed by the so called 'Information Structure'. This structure is concerned with the object types and their interrelationships in the modeled Universe of Discourse. The information structure of a conceptual schema is described in the following subsection. In doing so, we assume that the reader has some basic working knowledge of the concepts underlying ORM or ER.

### 2.1.1 Flat Conceptual Modeling

In [9] an ORM version is proposed which extends ORM with both top-down abstraction mechanisms as well as aspects from object oriented conceptual modeling techniques. The relation between that article and this article is that here we are concerned with an algorithm to 'reverse engineer' the abstraction layers from an existing flat conceptual model, whereas [9] provides the extensions to ORM needed to add abstractions in a top down way, which is necessarily a manual process. The output of the algorithm presented in this article can indeed be seen as a 3-dimensional ORM model fitting the top-down abstraction framework.

As a warning to readers of [9], it should be noted that what we call a flat conceptual schema in this article is in fact not a primitive, flat conceptual schema, as described in [9]. In this article our starting point is a schema consisting of object types, relationship types, and objectified relationship types. Objectified relationship types are also referred to as *nested object types*, and as its second name suggests, it already introduces depth into a conceptual schema. It is well known that objectifications can be replaced by so-called co-referenced object types. As an example of this, consider figure 2(a). This schema fragment is equivalent to the fragment depicted in figure 2(b). Depending on the universe of discourse that is being modeled, it may be more natural to use either one of the co-referenced or objectified representations. In

the abstraction algorithms, the choice between a co-referenced object type and an objectified relationship type is honored by treating them slightly differently.

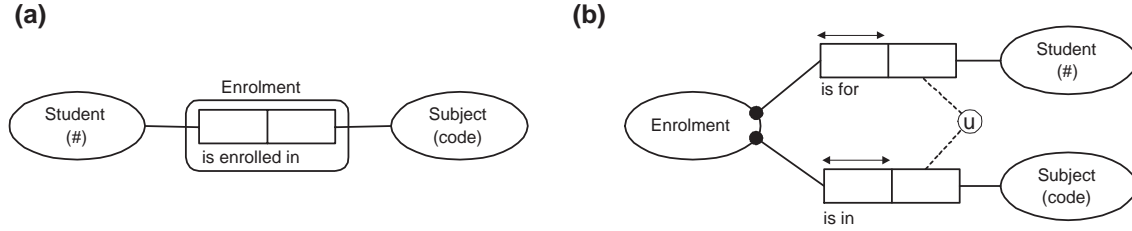


Figure 2: (a) Nested object type (b) Co-referenced object type

### 2.1.2 Typing Scheme

An ORM conceptual schema,  $\mathcal{CS}$ , is presumed to consist of a set of conceptual types,  $\mathcal{TP}$ . These types are divided into three main subclasses. The first class is the set of object types,  $\mathcal{OB}$ . Within this class a subclass of value types,  $\mathcal{VL}$ , can be distinguished. Instances from value types usually originate from some underlying domain such as strings, natural numbers, audio, video, etc. A separate class of types, the relationship types  $\mathcal{RL}$ , contain those types used to describe a relationship between one or more object types. Those object types which are not value types are called non-value types:  $\mathcal{NV} \triangleq \mathcal{OB} - \mathcal{VL}$ . In the formalization used in this paper, we allow types to belong to both the set of object types, and the set of relationship types. We refer to these later types as *nested object types* or *objectified relationship types*. Relationship types which are objectified do not belong to the set of value types. That is:

$$[\mathbf{O1}] \quad \mathcal{VL} \cap \mathcal{RL} = \emptyset$$

These types also have a number of structural properties which we now consider.

### 2.1.3 Roles in Relationship Types

Each relationship type in  $\mathcal{RL}$  contains a collection of roles. We refer to the set of all roles, in an ORM conceptual schema, as  $\mathcal{RO}$ . The roles in  $\mathcal{RO}$  are distributed among the relationship types by the function  $\text{Roles} : \mathcal{RL} \rightarrow \wp^+(\mathcal{RO})$ , which should provide a partition of the set of roles. (Note that  $\wp^+(\mathcal{RO})$  yields all non-empty subsets of  $\mathcal{RO}$ .) Each role has exactly one object type participating in it. This object type can be obtained by applying the function  $\text{Player} : \mathcal{RO} \rightarrow \mathcal{OB}$  to the relevant role. If the participants of a set of role is required, we use a generalization of this function,  $\text{Players} : \wp(\mathcal{RO}) \rightarrow \wp(\mathcal{RO})$ , which is defined as:

$$\text{Players}(v) \triangleq \{ \text{Player}(p) \mid p \in v \}$$

To determine the relationship type to which a given role belongs, the inverse of the Roles function ( $\text{Rel} : \mathcal{RO} \rightarrow \mathcal{RL}$ ) is used:

$$\text{Rel}(s) \triangleq p \text{ such that } s \in \text{Roles}(p)$$

In this formalization, the collection of roles contained in a relationship type is considered to have a pre-defined default order. This order, embedded in the verbalization of the relationship types, is provided by the domain experts during the initial analysis phase. As such, the PosN function is one of the knowledge sources from which we will try to mine the conceptual semantics hidden in the schema. The function  $\text{PosN} : \mathcal{RO} \rightarrow \mathbb{N}^+$  is used to assign a position to each role.

The predicate  $\text{Rels} : \wp(\mathcal{RO}) \rightarrow \wp(\mathcal{RL})$  is a generalization of the Rel function to sets of roles. It returns all relationship types involved in a given set of roles:

$$\text{Rels}(v) \triangleq \{ \text{Rel}(p) \mid p \in v \}$$

Similarly, the Roles function can be extended to accept a set of relationship types and return all roles involved in any of the given relationship types ( $\text{Roles} : \wp(\mathcal{RL}) \rightarrow \wp(\mathcal{RO})$ ):

$$\text{Roles}(v) \triangleq \{\text{Roles}(p) \mid p \in v\}$$

To conveniently access the roles involved in the same relationship type as a given role, we define the function  $\text{CoRoles} : \mathcal{RO} \rightarrow \wp(\mathcal{RO})$ , as:

$$\text{CoRoles}(v) \triangleq \text{Roles}(\text{Rel}(v))$$

Similarly to Roles and Rel, CoRoles can be extended to perform the same operation on a set of roles:

$$\text{CoRoles}(v) \triangleq \text{Roles}(\text{Rels}(v))$$

If we require those roles in the same relationship type as a given role, excluding the given role, we use the function  $\text{OtherRoles} : \mathcal{RO} \rightarrow \wp(\mathcal{RO})$ , which is defined as:

$$\text{OtherRoles}(p) \triangleq \text{CoRoles}(p) - \{p\}$$

### 2.1.4 Subtyping

The specialization relationship between a subtype and a supertype is captured by the relationship  $\text{SubOf} \subseteq \mathcal{OB} \times \mathcal{OB}$ . The intuition is that when  $x \text{ SubOf } y$ , the population of  $x$  is a definable subset of the population of  $y$ . Each subtype hierarchy (defined by SubOf) corresponds to a directed acyclic graph which adheres to the laws of transitivity and irreflexivity: The relation  $\text{Top}(x, y)$  is defined such that  $y$  is a top of  $x$  in the associated subtype hierarchy. The hierarchies we consider must always have one single top; so we can write  $\text{Top}(x) = y$  to refer to that unique top.

Given a set of object types in a subtype hierarchy, we can try to find the common supertypes in this hierarchy that are closest to these object types. To this end, we first need to find all common supertypes. This is done using the function  $\text{CommonSup} : \wp(\mathcal{OB}) \rightarrow \wp(\mathcal{OB})$ , which is defined as:

$$\text{CommonSup}(w) \triangleq \{x \in \mathcal{OB} \mid \forall y \in w [y \text{ SubOf } x]\}$$

The next step is to select those common supertypes that are closest to the given set of object types. We, therefore, introduce the notion of a lowest common supertype. The lowest common supertypes are those common supertypes which do not have any other common supertype as a subtype. A set of object types can actually have more than one lowest common supertype. The function  $\text{LowestCSup} : \wp(\mathcal{OB}) \rightarrow \wp(\mathcal{OB})$  is defined by:

$$\text{LowestCSup}(w) \triangleq \{x \in \text{CommonSup}(w) \mid \neg \exists y \in \text{CommonSup}(w) [y \text{ SubOf } x]\}$$

Given an object type  $x$  in a subtype hierarchy, we can determine the set of subtypes of this object type. This is done using the function  $\text{SubHierarchy} : \mathcal{OB} \rightarrow \wp(\mathcal{OB})$ , which is defined as:

$$\text{SubHierarchy}(x) \triangleq \{y \mid y \text{ SubOf } x\}$$

Consider, for example, the subtype hierarchy defined in figure 3. The supertypes of ‘Bicycle’ are ‘Car or Bicycle’ and ‘Vehicle’. The supertypes of ‘Car’ are ‘Car or Bicycle’, ‘Motorized Vehicle’ and ‘Vehicle’. The common supertypes (CommonSup) of ‘Car’ and ‘Bicycle’ are ‘Car or Bicycle’ and ‘Vehicle’ and the lowest common supertype (LowestCSup) of ‘Car’ and ‘Bicycle’ is ‘Car or Bicycle’.

### 2.1.5 Type Relatedness

Intuitively, object types may, for several reasons, have values in common in some populations. Two types are considered type related if their populations may share instances. Type relatedness, which we denote by  $x \sim y$ , is a property held only by object types which are in the same subtype hierarchy. For more detailed rules on type relatedness, refer to [18], [23].

Two roles are type related if their players are type related; so if  $p, q \in \mathcal{RO}$ , then:

$$p \sim q \triangleq \text{Player}(p) \sim \text{Player}(q)$$

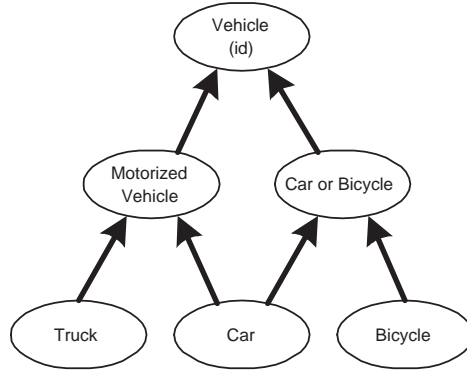


Figure 3: Example of a Specialization Hierarchy

### 2.1.6 Complete Information Structure

We can now define the basic information structure  $\mathcal{IS}$  of a conceptual schema  $\mathcal{CS}$  in terms of the following components:

$$\mathcal{IS} = \langle \mathcal{OB}, \mathcal{VL}, \mathcal{RL}, \mathcal{RO}, \text{Roles}, \text{PosN}, \text{SubOf}, \text{Player} \rangle$$

## 2.2 Conceptual Schema

Besides the information structure, a conceptual schema consists of constraints and derivation rules. For this article, only a limited class of constraints is of interest. The constraint classes we discuss, together with PosN, will be used as a source of information to decide which object types are major.

### 2.2.1 Mandatory Constraint

To specify the requirement that instances of a particular object type must always participate in at least one of some set of roles, we use the mandatory constraint,  $\text{Mand} \subseteq \wp(\mathcal{RO})$  (also referred to as ‘total role constraint’ in [23]). A mandatory constraint specifies that the union of the populations of the constrained set of roles must equal the total population of their player(s). All roles contained in a mandatory constraint must be type related. Therefore, we should have:

$$[\mathbf{O2}] \quad \text{Mand}(v) \Rightarrow \forall_{p,q \in v} [p \sim q]$$

A basic rule for ORM models (as defined in [17]) states that every instance of an object type must participate in at least one (fact type) role. In subsection 2.2.6 we will see that the only exception to this rule are the so-called lazy object types. This results in a mandatory role being implied over each set of type related roles. We identify the mandatory constraints which can be inferred in this way,  $\text{InferMand} \subseteq \wp(\mathcal{RO})$ , with the following derivation rule:

$$\text{InferMand}(V) \triangleq \forall_{p \in V} [p \sim q \iff q \in V \cap \mathcal{FR}] \wedge V \neq \emptyset$$

Note that  $\mathcal{FR}$  (the fact type roles) is the subset of  $\mathcal{RO}$  which is not used in the identification of any object type in the schema. This is more formally defined later in the paper. For abstraction purposes, we only consider those mandatory constraints which are not directly inferable.

### 2.2.2 Uniqueness Constraint

To introduce the concept of uniqueness, we use the predicate  $\text{Unique} \subseteq \wp(\mathcal{RO})$ . A uniqueness constraint requires each tuple in the projection of the join of the given roles (based on asserted join conditions) to appear only once. A uniqueness constraint which involves roles from only one predicate is referred to as an internal uniqueness constraint ( $\mathcal{IU}$ ).

$$\mathcal{IU} \triangleq \{ \text{Unique}(v) \mid \forall_{p,q \in v} [\text{Rel}(p) = \text{Rel}(q)] \}$$

An internal uniqueness constraint simply specifies that each tuple instance for that predicate will have a unique value combination for the constrained roles. If more than one predicate is involved in the uniqueness constraint, then the constraint is classified as an *external* uniqueness constraint ( $\mathcal{EU}$ ).

$$\mathcal{EU} \triangleq \{ \text{Unique}(v) \mid \forall_{p,q \in v} [\text{Rel}(p) \neq \text{Rel}(q)] \}$$

In this case, the predicates involved must be *joinable via common object types* [19]. The general interpretation of a uniqueness constraint is formulated in the *Uniquet Algorithm* provided in [34].

### 2.2.3 Primary Uniqueness Constraint

For every object type in the data schema, there must be some way in which to uniquely identify each instance of that object type. In other words, we insist that every object-type is *identifiable*.

To identify the instances of non-value types ( $\mathcal{NV}$ ), one uniqueness constraint must be selected to be the primary means of identification for that object type. We call the set of such uniqueness constraints ‘PUnique’ and require that  $\text{PUnique}(v) \Rightarrow \text{Unique}(v)$ . If this uniqueness constraint only involves one role, the identification scheme is often collapsed into a reference mode for graphical convenience. The reference mode of an object type is placed in brackets under the object type name. For example, figure 4(b) shows the graphical abbreviation for the explicit identification scheme represented in figure 4(a).

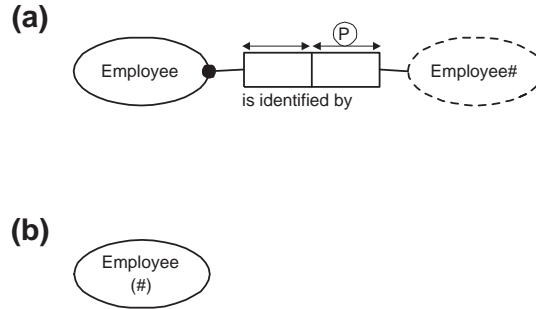


Figure 4: (a) Explicit identification scheme (b) Implicit identification scheme

The algorithms in this paper do not consider this graphical abbreviation. Instead, they presume that all reference schemes are explicitly represented through uniqueness constraints. For more information about primary uniqueness constraints refer to [17]. For more detailed formal requirements on identification in ORM schemas, refer to [19], [23]. Every non-value object type must have exactly one primary identification scheme. That is:

$$[\mathbf{O3}] \quad \forall_{x \in \mathcal{NV}} \exists!_v [\text{PUnique}(v) \wedge x \in \text{Identifies}(v)]$$

where *Identifies* is defined as:

$$\text{Identifies}(v) \triangleq \text{SubHierarchy}(\text{LowestCSup}(\{ \text{Player}(p) \mid p \in \text{OtherRoles}(v) \}))$$

The set of roles which are contained in the primary uniqueness constraint of a given non-value type is given by  $\text{PIdRoles} : \mathcal{NV} \rightarrow \wp(\mathcal{RO})$  such that:

$$\text{PIdRoles}(x) \triangleq \{ p \mid \exists_v [p \in v \wedge x \in \text{Identifies}(v)] \}$$

The set of predicates which are used to identify a given non-value type is given by the function  $\text{PIdRels} : \mathcal{NV} \rightarrow \wp(\mathcal{RL})$  such that:

$$\text{PIdRels}(x) \triangleq \{ \text{Rel}(r) \mid r \in \text{PIdRoles}(x) \}$$



## 2.2.4 Occurrence Frequency Constraint

Uniqueness constraints are used to specify that instances of object types may play a certain combination of roles at most once. Occurrence frequency constraints specify the more general condition that the number of times that object instances may play a combination of roles is restricted to within a fixed range. The condition that the instances of a set of roles  $\sigma$  must occur at least  $n$  and at most  $m$  times is denoted by  $\text{Frequency}(\sigma, n, m)$ . The semantics of  $\text{Frequency}$  are fully defined in [19].

The function  $\text{MaxFreq}$  returns the maximum number of times an object type instance may participate in the given role.  $\text{MaxFreq} : \mathcal{RO} \rightarrow \mathbb{N}$  is defined as:

$$\text{MaxFreq}(r) \triangleq \begin{cases} 1 & \text{if } \text{Unique}(\{r\}) \\ m & \text{if } \text{Frequency}(\{r\}, n, m) \\ \infty & \text{otherwise} \end{cases}$$

Note that when taking set types, sequence types, etc. into consideration,  $\text{Unique}(\sigma)$  should be replaced by  $\text{Unique}(\sigma) \vee \text{ExUnique}(\sigma)$ , where  $\text{ExUnique}(\sigma)$  is the class of *existential uniqueness constraints*. This later constraint class is crucial in defining complex types such as set types [24], [25].

## 2.2.5 Set-comparison Constraints

Set-comparison constraints (which we will refer to as ‘set constraints’) are used to specify conditions which apply between the population sets of two role sequences. If  $X$  is a set, then  $X^+$  denotes the set of sequences built from elements of  $X$ . For sequences, we presume that the operation  $z[i]$  returns the  $i$ -th element of sequence  $z$ .  $\text{Set}(z)$  coerces a sequence  $z$  into a set of elements, so:

$$\text{Set}(z) \triangleq \{x \mid \exists_i [z[i] = x]\}$$

We use  $p \in z$  as an abbreviation for  $\exists_i [z[i] = p]$  while  $|z|$  denotes the length of sequence  $z$ . To determine which position a particular element occupies in a given sequence, we use the function  $\text{Pos}$ . For sequences where no two elements appear more than once in the sequence, we can define  $\text{Pos}$  as:

$$\text{Pos}(p, z) \triangleq i \text{ such that } z[i] = p$$

The relations  $\text{Subset}$ ,  $\text{Equality}$ ,  $\text{Exclusion}$  each apply to an ordered pair of role sequences (we do not consider an  $n$ -ary form of exclusion constraint in this paper). The subset constraint (defined by relation  $\text{Subset} \subseteq \mathcal{RO}^+ \times \mathcal{RO}^+$ ) specifies that the population of the first role sequence is necessarily a subset of the population of the second role sequence; the equality constraint ( $\text{Equality} \subseteq \mathcal{RO}^+ \times \mathcal{RO}^+$ ) specifies that the population of the first role sequence must be exactly equal to the population of the second role sequence; while the exclusion constraint ( $\text{Exclusion} \subseteq \mathcal{RO}^+ \times \mathcal{RO}^+$ ) specifies that the population of the first role sequence does not contain any tuple which is in the population of the second role sequence. For a more formal definition of these constraints, please refer to [19] and [23].

From their definitions, it is easy to infer an implied subset constraint between every optional role and every mandatory role played by the same object type. Similarly it is possible to infer an *implied* equality constraint between every mandatory role played by the same object type. We will not consider set constraints which are inferable in this manner. That is:

$$\forall_{p,q:\text{Subset}(\langle p \rangle, \langle q \rangle) \vee \text{Equality}(\langle p \rangle, \langle q \rangle)} [\neg \text{Mand}(q) \wedge (\text{Mand}(p) \Rightarrow \text{Player}(p) \text{SubOf} \text{Player}(q))]$$

From  $\text{Subset}$ ,  $\text{Equality}$  and  $\text{Exclusion}$ , we derive the more generic predicate  $\text{SetCon}$  using the following rules:

$$\begin{aligned} \text{SC}(v, w) &\triangleq \text{Subset}(v, w) \vee \text{Equality}(v, w) \vee \text{Exclusion}(v, w) \\ \text{SetCon}(v, w) &\triangleq \text{SC}(v, w) \vee \text{SC}(w, v) \end{aligned}$$

The underlying intuition is that if  $\text{SetCon}(v, w)$ , then some set constraint exists which involves the roles in  $v$  and  $w$ .

From these rules, we can specify an even more generic definition for  $\text{SetCon}$  with only a single parameter. If  $\text{SetCon}(v)$  then some set constraint exists which involves the roles in  $v$ .

$$\text{SetCon}(v) \triangleq \exists_w [\text{SetCon}(v, w)]$$

## 2.2.6 Refinements to the Type Classification

Relationship types can now be partitioned into two important subclasses - the *fact types* and the *reference types*. Reference types ( $\mathcal{RF}$ ) are those relationship types which are used within the primary identification scheme of some non-value type:

$$\mathcal{RF} \triangleq \{ \text{Rel}(p) \mid \exists x \in \text{ab03nvtty} [p \in \text{PIdRoles}(x)] \}$$

Fact types ( $\mathcal{FT}$ ) are those relationship types which are not used within the primary identification scheme of a basic entity type or subtype.

$$\mathcal{FT} \triangleq \mathcal{RL} - \mathcal{RF}$$

The set  $\mathcal{FR} \subseteq \mathcal{RO}$  is used to refer to those roles which are contained within a fact type. That is:

$$\mathcal{FR} \triangleq \{ p \in \mathcal{RO} \mid \text{Rel}(p) \in \mathcal{FT} \}$$

Most object types can only be instantiated by instances which participate in some fact type ( $\mathcal{FT}$ ). Instances of *lazy object types*, however, can exist without participating in any fact type. We graphically represent a *lazy object type* by concatenating an exclamation mark to the end of the object type name (x!). As an example, consider figure 5(a). Only countries which border another country can be recorded. In figure 5(b), however, countries may be recorded even if they do not border (by land) any other country (e.g. Australia).

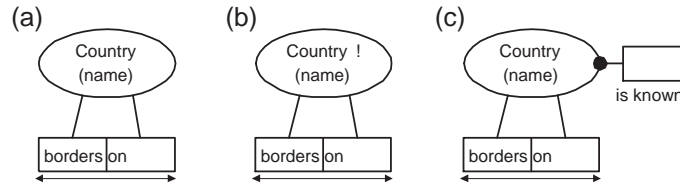


Figure 5: 'Country' object type represented as (a) non-lazy (b) lazy and (c) implied lazy

A lazy entity type's behavior can be compared to that of an object type that mandatorily participates in a unary (one roled) fact type which represents its existence (as depicted in figure 5(c)). For the purposes of this paper, we consider a lazy entity type to be a graphical simplification to conveniently represent those entity types which participate in a single, mandatory unary role. No special consideration is therefore necessary for lazy entity types in the ensuing algorithms.

As stated before, other complex types like set types and sequence types are not discussed in full detail in this article, however, we will briefly return to this issue.

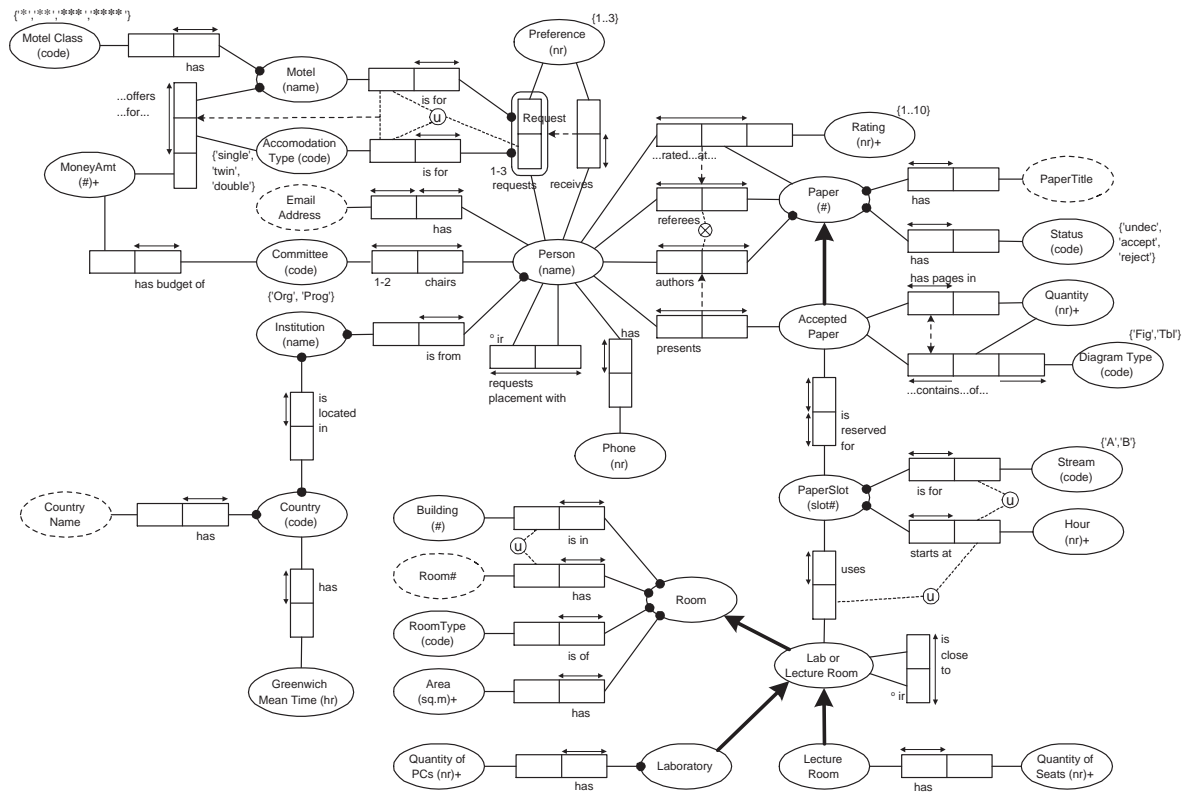
## 2.3 Summary

A conceptual schema  $\mathcal{CS}$  can now be defined in terms of both the information structure  $\mathcal{IS}$  and the basic constraints which apply to this information structure.

$$\mathcal{CS} = \langle \mathcal{IS}, \text{Mand}, \text{Unique}, \text{PUnique}, \text{Frequency}, \text{Subset}, \text{Equality}, \text{Exclusion} \rangle$$

Conceptual schemas can have many other components, including ring constraints, subtype definitions, derived fact types and other extraneous constraints. None of these, however, will be considered in this paper, because they do not impact on the abstraction algorithms presented.

An example ORM conceptual schema can be found in figure 6. Entity types are depicted as named, solid ellipses. Value types are shown as named, broken ellipses. Predicates are shown as named sequences of role boxes, with the predicate name located in or beside the first role of the predicate. A nested object type is shown as a frame around a predicate (e.g. 'Request'). Arrow-tipped bars over one or more role boxes indicate an internal uniqueness constraint over these roles. A black dot at the base of a connector between an object type and a role indicates a mandatory constraint. Other constraints are represented as defined in [17]. As an example, consider the conceptual schema depicted in figure 6. In this schema we have  $\text{EmailAddress} \in \mathcal{VL}$ ,  $\text{Preference} \in \mathcal{OB}$ , and  $\text{requests} \in \mathcal{RL}$ . This schema is used as the running example throughout this article.



**SUBTYPE DEFINITIONS**

- each Accepted Paper is a Paper that has Status 'accept'
- each LabOrLecRoom is a Room that is of RoomType in {'lab', 'lec'}
- each Laboratory is a Room that is of RoomType 'lab'
- each Lecture Room is a Room that is of RoomType 'lec'

Figure 6: Academic Conference Example Application

### 3 Anchoring Fact Types

In the following sections, we introduce a method by which we can view an ORM conceptual schema at various levels of abstraction. At each subsequently higher level of abstraction, we show only the most conceptually important (major) object types from the previous level, thereby creating a procedure that generates an incremental summary of the schema based on conceptual relevance. To this end, we first provide a mechanism by which the major object types in a given schema can be derived.

Throughout this paper, sets of conceptual objects from the current conceptual schema  $\mathcal{CS}$  will be denoted as belonging to a particular abstraction view by subscripting the set with the abstraction level. For example,  $TP_i$  refers to the set of types in  $\mathcal{CS}$  which appear at abstraction level  $i$ .

#### 3.1 Definitions

To define the notion of a *major* object type (relative to a particular abstraction level  $i$ ) we consider each fact type individually and decide which object type(s) is (are) the most conceptually important participants in this fact type. We say that a role *anchors* a fact type to its player at the current abstraction level, if that player is (one of) the conceptually most important participants in the fact type. Conceptual importance is, to a certain degree subjective. However, a reasonable (and often measurable) indicator of conceptual importance is the proportion of the population of each object type that participates in the fact type.

The conceptual importance of role  $p$  played by object type  $X$  in fact type predicate  $F$  can be indicated by:

$$\begin{cases} \frac{|\text{Pop}(\pi_p F)|}{|\text{Pop}(X)|} & \text{if } |\{q \in \mathcal{FR} \mid \text{Player}(q) = X\}| > 1 \\ 0 & \text{otherwise} \end{cases}$$

where  $\text{Pop}$  is the population function and  $\pi_p F$  indicates the projection on role  $p$  of fact type  $F$ . Consider, for example, a fact type ‘*Subject is lectured by Academic*’. Suppose we know that a greater percentage of ‘Subjects’ are lectured by an ‘Academic’ than the percentage of ‘Academics’ who lecture a ‘Subject’; so:

$$\frac{|\text{Pop}(\pi_1('Subject is lectured by Academic'))|}{|\text{Pop}(Subject)|} > \frac{|\text{Pop}(\pi_2('Subject is lectured by Academic'))|}{|\text{Pop}(Academic)|}$$

It is obvious that, as a result, a particular ‘Subject’ is more likely to be participating in the fact type than a particular ‘Academic’. It can also be observed that the fact type ‘*Subject is lectured by Academic*’ is more likely to be accessed in relation to a particular ‘Subject’ than in relation to a particular ‘Academic’. We therefore consider ‘Subject’ to be the more ‘conceptually important participant’ and consider ‘*Subject is lectured by Academic*’ to be anchored on the role played by ‘Subject’. This reasoning will, in general, only be useful when we have access to a typical population of a conceptual schema. When only the conceptual schema is available, we must rely on the conceptual constraints to derive such information from the type level. This is the approach taken in this paper.

The fact that a given role is an anchor, at abstraction level  $i$ , is captured by the predicate:  $\text{Anchor}_i \subseteq \mathcal{FR}_i$ , where  $\mathcal{FR}_i$  represents the fact type roles which are present at abstraction level  $i$ . For convenience, we also introduce the infix predicate  $\text{AnchoredTo}_i \subseteq \mathcal{FT}_i \times \mathcal{OB}$ , which indicates that the given role is anchored to the given object type (at abstraction level  $i$ ):

$$r \text{ AnchoredTo}_i x \triangleq \text{Anchor}_i(r) \wedge \text{Player}(r) = x$$

When considering anchors, it is important to do so in their proper context, i.e. at a particular level of abstraction. For example a role which receives one hundred percent participation (i.e. a mandatory role), may become implied mandatory at a higher level of abstraction and consequently lose ‘conceptual importance’. As an example, consider the schema fragments in figure 7. In figure 7(a), ‘Subject’ mandatorily participates in ‘*Employee teaches Subject*’ and ‘Department’ mandatorily participates in ‘*Employee works for Department*’. Because these object types have one hundred percent participation, the corresponding roles are therefore considered to be anchors in figure 7(a). In figure 7(b), which shows the next highest level of abstraction, however, the same roles (played by ‘Subject’ and ‘Department’) are only mandatory by

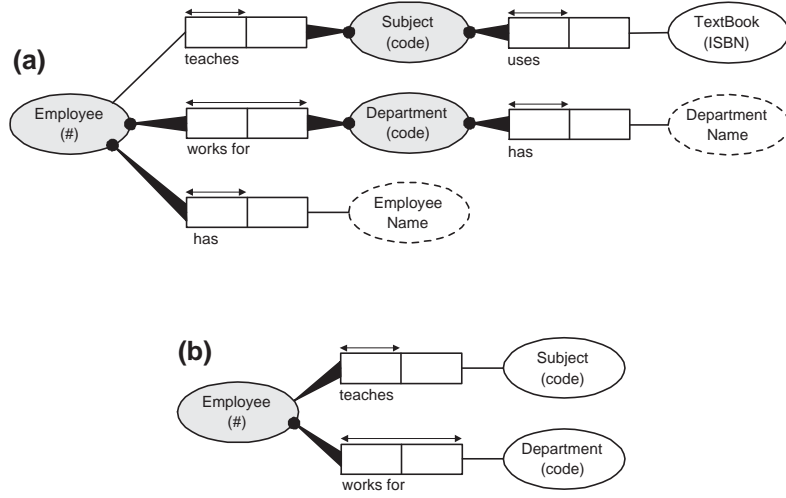


Figure 7: (a) Abstraction Level  $n$  (b) Abstraction Level  $n + 1$

implication. The roles played by ‘Subject’ and ‘Department’, therefore, lose their conceptual importance.

Anchors for fact types are selected by comparing the conceptual importance of the roles involved. To this end, we introduce the notion of the weight of role, to indicate how firmly the role is attached to its player:

$$\text{Weight}_i : \mathcal{FR}_i \rightarrow \mathbb{N}$$

The weight function is used as a rough indicator of both the importance of each role within a fact type and of the relative importance of anchors between fact types. A role is considered to be an anchor if its weight is the highest (or equal highest) weight of any role in the same fact type. An anchor can therefore be defined in terms of the role weights, as such:

$$\text{Anchor}_i(r) \triangleq \text{Weight}_i(r) > 0 \wedge \text{Weight}_i(r) = \max \{ \text{Weight}_i(p) \mid p \in \text{CoRoles}(r) \}$$

A corollary which ensues is:

$$\text{Corollary 3.1} \quad (\text{Weight}_i(r) < \text{Weight}_i(p) \wedge p \in \text{CoRoles}(r)) \Rightarrow \neg \text{Anchor}_i(r)$$

The particular weight associated with each anchor is defined by the `WeightSchema` procedure, which is introduced in the next subsections.

We consider a fact type predicate to be anchored ( $\text{Anchored}_i \subseteq \mathcal{FT}_i$ ) if it contains a role which anchors that fact type predicate to an object type:

$$\text{Anchored}_i(s) \triangleq \exists_{r \in \text{Roles}(s)} [\text{Anchor}_i(r)]$$

or in other words, if the sum of the weights of each of its roles is greater than zero:

$$\text{Corollary 3.2} \quad \text{Anchored}_i(s) \triangleq \sum_{r \in \text{Roles}(s)} \text{Weight}_i(r) > 0$$

We refer to an object type, which has at least one fact type anchored to it, as an  $\text{AnchorPoint}_i \subseteq \mathcal{OB}_i$ :

$$\text{AnchorPoint}_i(x) \triangleq \exists_{r \in \mathcal{FR}_i} [\text{Player}(r) \sim x \wedge \text{Anchor}_i(r)]$$

For an abstraction level  $i$  to be completely anchored, every fact type within the schema must be anchored:

$$\text{AnchoredSchema}(i) \triangleq \forall_{s \in \mathcal{FT}_i} [\text{Anchored}_i(s)]$$

This is the goal for this section.

## 3.2 Weighting a Schema

The procedure called `WeightSchema`, shown below, automatically assigns default weights to each fact-type role, based on the given semantic constraints within the associated conceptual domain. `Weight` is a total function. Since there will always be some subjective qualities that can not be captured by such an automatic procedure, it is important that the user has the ability to override some automatic weighting decisions that may be questionable. For this reason, and because the user will usually only want to express such alternative preferences once, we allow our automatic abstraction procedure to take previous user-driven weightings into account.

The automatic weighting is defined by a set of weighting rules that associate a weight to each fact type role, based on the context of these roles. The weighting algorithm works by continuously trying to increase the weight of the roles. This is a repetitive process, as increases of weights in one part of the schema may lead to further increases in other parts of the schema.

We will refer to the weightings, which were explicitly generated by a user decision, through the function `UserWeight` :  $\mathcal{FR} \rightarrow \mathbb{N}$ . The weights that are derived automatically from the underlying schema are provided by the function `AutoWeight` :  $\mathbb{N} \times \mathcal{FR} \rightarrow \mathbb{N}$ . This function is introduced in the next subsection by a set of derivation rules. We employ the notation `AutoWeighti(p)` rather than `AutoWeight(i, p)` because of the fact the index  $i$  is used as a label rather than carrying any specific semantics.

Once `WeightSchema` has been automatically performed, the user would have a further opportunity to alter the `UserWeights` by modifying the set of anchors produced in accordance with an appropriate set of modification rules.

```

WeightSchema : CS → (FR → ℕ )
WeightSchema(CS) ≜
VAR
  Weight : FR → ℕ ;
  p : FR ;
BEGIN
  { Initialize anchors }
  FOR EACH p ∈ FR DO
    Weight(p) := UserWeight(p) ;
  END ;
  WHILE ∃i,p [AutoWeighti(p) > Weight(p)]
    Weight(p) := max { AutoWeighti(p) | ⟨i, p⟩ ∈ dom(AutoWeight) } ;
  END ;
  RETURN Weight ;
END WeightSchema ;

```

This algorithm will always terminate. From the condition on the `WHILE` loop we can see it terminates if  $\neg \exists_{i,p} [\text{AutoWeight}_i(p) > \text{Weight}(p)]$ . In the next subsection we will also see that the maximum value returned by `AutoWeight` is fixed to 10. From the body of the `WHILE` loop follows that we never reduce the `Weight` of a role. This means that for any role  $p$  once  $\text{Weight}(p) \geq 10$  we cannot find a rule labeled  $i$  such that  $\text{AutoWeight}_i(p) > \text{Weight}(p)$ . Which means the loop must eventually terminate.

## 3.3 Rules for Role Weighting

The following paragraphs describe each of the twelve rules that together define `AutoWeight`. The resulting weightings returned by these rules serve as a comparative guide, and should at some stage be refined based on empirical testing in practical situations. The existing rules have been formulated after studying a number of cases to observe the effect of particular constraints on the associated populations and on the conceptual importance of surrounding object types.

### Rule 1 – Mandatory Roles

All non-implied mandatory roles have, by definition, full participation by the population of the player(s). Therefore, any fact type role which is involved in a mandatory role constraint (even if this is a disjunctive

mandatory constraint) should be weighted, *unless* the mandatory role constraint is implied (as described by InferMand). This is the only rule which can cause a fact type predicate to be anchored more than once.

$$\text{AutoWeight}_1(p) \triangleq \max \left( \left\{ \left\lceil \frac{10}{|v|} \right\rceil \mid \text{Mand}(v) \wedge p \in v \wedge \neg \text{InferMand}(v) \right\} \cup \{0\} \right)$$

A role can be involved in a number of mandatory role constraints. The simplest case would be where the role itself is mandatory, which would lead to a weight of 10. However, a role may be involved in a disjunctive mandatory role. This means that the instances of the participating object type must play at least one of the roles involved in the disjunctive mandatory role. In this latter case, the weight of 10 is ‘shared’ among the involved roles. As one role may be involved in a number of mandatory role constraints, we take the weight to be the maximum of the possible weights that would follow from these involvements.

For example, rule 1 would cause each fact type in figure 8 to be anchored towards the non-implied mandatory role played by Employee.

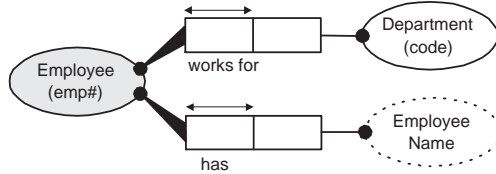


Figure 8: All non-implied mandatory roles are weighted.

This rule can also be considered in the broader context of complex types like sequence types, bag types, etc. We can now discuss why these complex types do not require special provisions in our algorithm. In figure 9(a) we show an example of a set type: namely ‘Convoy’. A convoy consists of a set of ships, each of which is commanded by a unique captain. Both a ship and its captain are each individually identified by a name. A convoy, however, is identified by a set of ships.

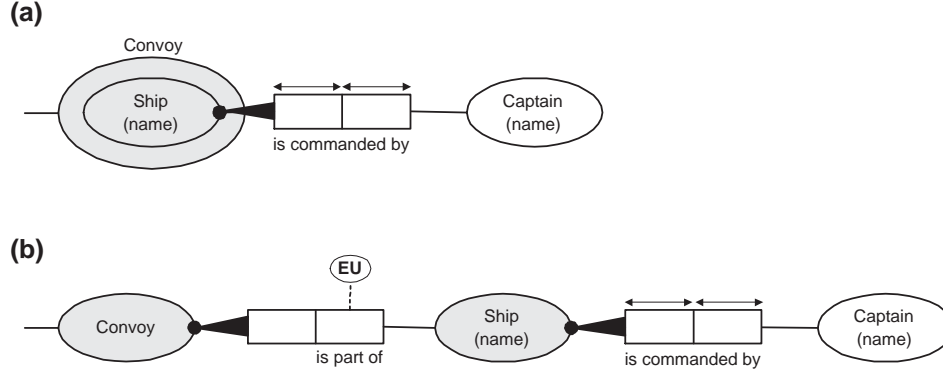


Figure 9: A convoy of ships modeled using (a) a complex set type and (b) elementary relationships

In figure 9(b), this set type is modeled in terms of more elementary relationships using the existensial uniqueness constraint (represented by the encircled EU symbol) [25], [24]. The AutoWeight rules that we are defining, can therefore be directly applied to the elementary representation of the complex types.

## Rule 2 – Unary Roles

The player of the only role in a unary predicate must obviously be ‘the most important participant’ in that predicate. All roles in unary predicates are therefore weighted.

It should be remembered that, for the purposes of this algorithm, lazy object types are treated like non-lazy object types which play a mandatory unary predicate representing the existence of the instances.

$$\text{AutoWeight}_2(p) \triangleq \text{if } \text{CoRoles}(p) = \{p\} \text{ then } 10 \text{ else } 0$$

Figure 10 shows an example subschema in which every unary predicate is anchored on its one role.

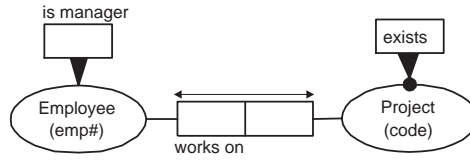


Figure 10: All unary roles are anchored.

### Rule 3 – Non-Leaf Object Types

A leaf facttype role ( $\text{Leaf} \subseteq \mathcal{FR}$ ) is one which has a player that plays only that fact type role. That is:

$$\text{Leaf}(p) \triangleq \neg \exists q \in \mathcal{FR} [\text{Player}(p) = \text{Player}(q) \wedge \text{Rel}(p) \neq \text{Rel}(q)]$$

If only one role in a fact type is played by a non-leaf object type, then this role is considered ‘conceptually important’ enough to be given a strong weighting.

$$\text{AutoWeight}_3(p) \triangleq \text{if } \text{Leaf}(p) \wedge \forall q \in \text{OtherRoles}(p) [\text{Leaf}(q)] \text{ then } 9 \text{ else } 0$$

This rule as such is rather straightforward. The reason to assign only a weight of 9, instead of 10, is that mandatory non-leaf roles are considered to be conceptually more important than optional non-leaf roles.

In the example subschema in figure 11, rule 3 would be fired, causing both fact types to be anchored towards the Employee object type. Notice that ‘Room’ is actually a leaf object type because, while it participates in three roles, it only participates in one fact type role. ‘Room is in Building’ and ‘Room has Room#’ are not considered in the weighting procedure as they are both reference types.

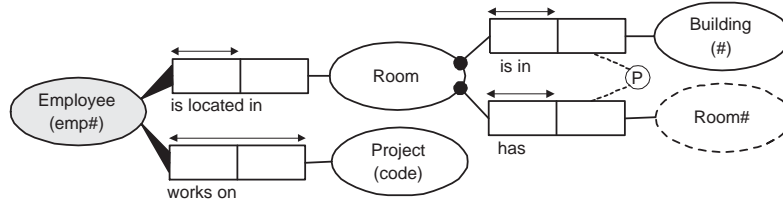


Figure 11: Non-leaf object types may indicate automatic anchorage.

### Rule 4 – Smallest Maximum Frequency

The maximum frequency of the population of a role can be determined from one of two constraints. A single role uniqueness constraint indicates that the role has a maximum frequency of one. Alternatively, an occurrence frequency constraint often explicitly specifies the maximum frequency of a role. If exactly one role within a fact type predicate has a smaller maximum frequency than all other roles in that fact type, then this role should be anchored.

$$\text{AutoWeight}_4(p) \triangleq \text{if } \forall q \in \text{OtherRoles}(p) [\text{MaxFreq}(p) < \text{MaxFreq}(q)] \text{ then } 2 + \left\lceil \frac{6}{\sqrt{\text{MaxFreq}(p)}} \right\rceil \text{ else } 0$$

The closer the maximum frequency of a role is to 1, the higher the weighting. The maximum AutoWeight of 8 is applied in those cases in which a uniqueness constraint holds on the role, causing MaxFreq to be 1. If the maximum frequency is higher than 1, the AutoWeight will become lower and lower, down to a minimum of 2. However, because increments in MaxFreq should have less effect if the frequency is already high, we



have taken the division of the square root of the MaxFreq value. For example, the increment from a MaxFreq of 8 to one of 9 will have less effect on the AutoWeight than an increment from 1 to 2. The result is a curve that drops down quickly from a maximum Weight of 8, but starts to level out when it gets closer to 2.

The example in figure 12 depicts a subschema in which the fact type ‘Project is managed by Employee’ is anchored due to a uniqueness constraint, and the fact type ‘Employee works on Project for Duration’ is anchored due to a frequency constraint.

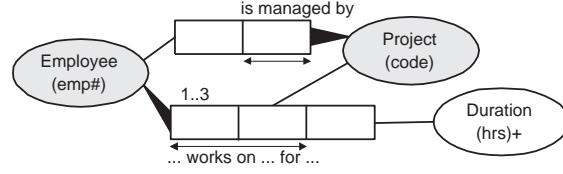


Figure 12: Roles with the smallest maximum frequency may be anchored.

### Rule 5 – Non-Value Types

If exactly one role in a fact type is played by a non-value type, then the fact type should be anchored on this role.

$$\text{AutoWeight}_5(p) \triangleq \text{if } \text{Player}(p) \notin \mathcal{VL} \wedge \forall_{q \in \text{OtherRoles}(p)} [\text{Player}(q) \in \mathcal{VL}] \text{ then } 7 \text{ else } 0$$

The rationale behind this is that value types are by definition conceptually less important than non-value types.

In the example shown in figure 13, rules 1 to 4 fail to determine an appropriate anchorage for either fact type. Rule 5, however, triggers the obvious conclusion that both fact types should be anchored towards ‘Employee’.

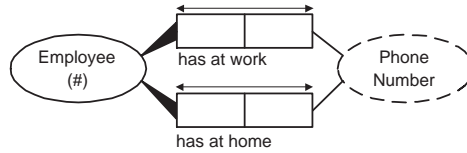


Figure 13: Roles played by non-value-types may become automatically anchored.

### Rule 6 – Anchor Points

As we have already discussed, those object types which serve as anchor points to fact types are considered to possess a relatively high conceptual importance. Therefore, if exactly one role in a given fact type is played by an object type which became an anchorpoint via rules 1 to 5, the fact type is anchored on this role. For this purpose we introduce the notion of a ‘heavy role’ as:

$$\text{HeavyRole}(p) \triangleq \exists_{s: \text{Weight}(s) \geq 7} [p \sim s]$$

The AutoWeight rule then becomes:

$$\text{AutoWeight}_6(p) \triangleq \text{if } \text{HeavyRole}(p) \wedge \forall_{q \in \text{OtherRoles}(p)} [\neg \text{HeavyRole}(q)] \text{ then } 6 \text{ else } 0$$

In the example subschema of figure 14, the uniqueness constraint on ‘Employee is managed by Project’ causes rule 4 to anchor the upper fact type towards ‘Project’. Since ‘Project’ is now the only participant in the ternary fact type which is an anchorpoint, the lower fact type is anchored on ‘Project’.

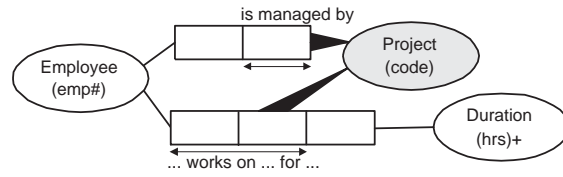


Figure 14: Roles played by anchorpoints may become automatically anchored.

### Rule 7 – Single-Role Set Constraints

If a fact type is involved in exactly one single-role set constraint (i.e. subset, equality or exclusion constraint), and the role at the other end of the set constraint is anchored, then the constrained role in the given fact type is anchored.

$$\text{AutoWeight}_7(p) \triangleq \text{if } \exists_{s:\text{Anchor}(s)} [\text{SetCon}(\langle s \rangle, \langle p \rangle)] \wedge \forall_{q \in \text{OtherRoles}(p)} [\neg \text{SetCon}(\langle q \rangle)] \text{ then } 5 \text{ else } 0$$

In figure 15, the fact type ‘*Employee* is supervisor in *Project*’ is anchored to ‘*Employee*’ by rule 4, as a result of the simple uniqueness constraint. Since the role played by ‘*Employee*’ in the ternary fact type is connected to this anchored role via a single-role subset constraint, this role is consequently anchored by rule 7.

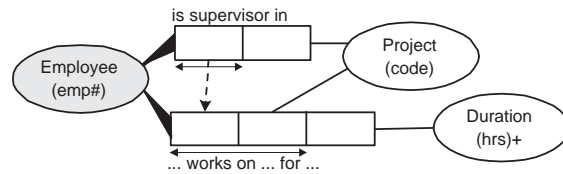


Figure 15: A role connected to an anchored role by single role set constraints is anchored.

It is important to consider the case in figure 16, in which the single role set constraints contradict each other. In this case, rule 7 could not produce a determinant anchorage for the predicate ‘works on’. We therefore ensure that this rule only fires on non-anchored fact types which are only involved in one single-role set constraint.

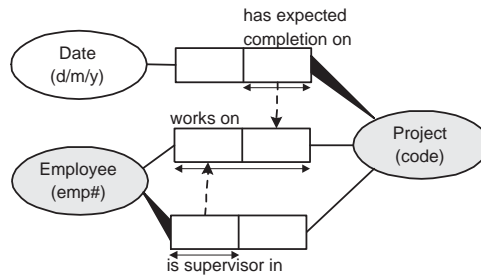


Figure 16: Rule 7 does not consider cases in which a single-role set constraint contradicts another.

It is also important to realize why we chose not to require the contradicting set constraint to necessarily have an anchor assigned. It would have been inadequate to require the following condition on the other roles in  $p$ 's fact type:

$$\forall_{q \in \text{OtherRoles}(p)} \neg \exists_{s:\text{Anchor}(s)} [\text{SetCon}(\langle s \rangle, \langle q \rangle)]$$

We illustrate this, by considering the case shown in figure 17. If we assume that the above condition is adequate, (i.e. that rule 7 is fired as long as no other role in the fact type participates in an anchored single

role subset constraint), then two possible scenarios are possible for the schema fragment below. Firstly, rule 7 could cause ‘*Person* owns *Car*’ to be anchored towards ‘*Person*’; which would then cause ‘*Person* has driven *Car*’ to also be anchored towards ‘*Person*’. Alternatively, rule 7 could first cause ‘*Person* caused crash of *Car*’ to be anchored towards ‘*Car*’; which would then cause ‘*Person* has driven *Car*’ to also be anchored towards ‘*Car*’. As a result, ‘*Person* has driven *Car*’ could be anchored in either direction, depending on the order in which the rule was fired.

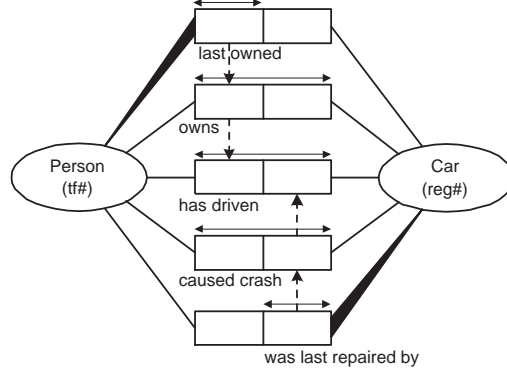


Figure 17: Rule 7 requires that  $p$ 's fact type participates in *no* other single-role subset constraint.

For this reason, we only allow fact types to be anchored on a role,  $p$ , if no other role in its fact type is involved in any kind of single-role set constraint (as defined in  $\text{AutoWeight}_7(p)$ ). The definition of rule 7 will, therefore anchor all fact types in figure 17, except for ‘*Person* has driven *Car*’.

### Rule 8 – Multi-Role Set Constraints

If a fact type is involved in exactly one (possibly multi-role) set constraint (i.e. subset, equality or exclusion constraint), and exactly one of the roles in the fact type is in the corresponding position within the set constraint as an anchored role, then this role is itself anchored.

$$\text{AutoWeight}_8(p) \triangleq \mathbf{if} \exists_{v,w:\text{SetCon}(v,w)} [p \in v \wedge \text{Anchor}(w[\text{Pos}(p,v)])] \wedge \text{SingleSetCon}(v,p) \mathbf{then} 4 \mathbf{else} 0$$

where  $\text{SingleSetCon}(v,p)$  enforces the singularity of the set constraint  $v$  with respect to role  $p$ :

$$\text{SingleSetCon}(v,p) \triangleq \forall_{x:\text{SetCon}(x)} [v \neq x \Rightarrow \text{OtherRoles}(p) \cap \text{Set}(x) = \emptyset]$$

In the example in figure 18, the fact type ‘*Employee* is supervisor in *Project*’ is anchored to ‘*Employee*’ by rule 4, as a result of the simple uniqueness constraint. Since the role played by ‘*Employee*’ in the ternary fact type is connected to this anchored role via a multi-role subset constraint, this role is consequently anchored by rule 8.

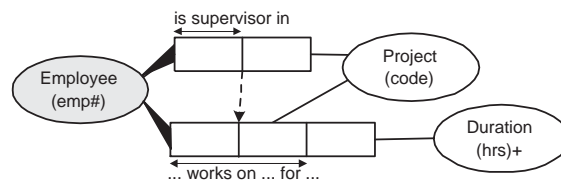


Figure 18: Roles connected by multi-role set constraints are anchored.

Similarly to rule 7, it is important to consider the case in figure 19, in which the multi-role set constraints contradict each other. In this case, rule 8 would not produce a determinant anchorage for the predicate ‘works on’. We therefore only use this rule on non-anchored fact types which are only involved in one multi-role set constraint.

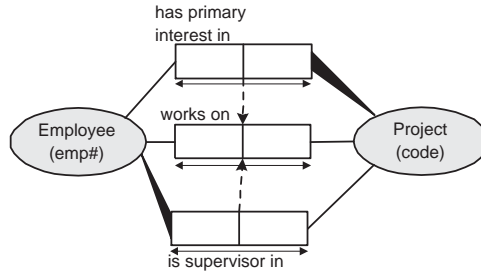


Figure 19: Rule 8 does not consider cases in which a multi-role set constraint contradicts another.

### Rule 9 – Set Constraints and AnchorPoints

If there exists a non-implied set constraint in which one of the roles involved in the constraint is the only involved role in its fact type to be played by an anchorpoint and the fact type of the role corresponding to it in the other role sequence is not anchored, then this role should become an anchor. If a set constraint  $v$  anchors a role  $p$  in this sense, then we refer to this as  $\text{Anchors}(v, p)$ :

$$\text{Anchors}(v, p) \triangleq p \in v \wedge \text{AnchorPoint}_i(\text{Player}(p))$$

The resulting AutoWeight rule is then:

$$\text{AutoWeight}_9(p) \triangleq \text{if } \exists_{v:\text{SetCon}(v)} [\text{Anchors}(v, p) \wedge \forall_{r \in \text{OtherRoles}(p)} [\neg \text{Anchors}(r, v)] \wedge \text{SingleSetCon}(v, p)] \text{ then } 3 \text{ else } 0$$

For example, in figure 20, the firing of rule 1 causes ‘Project has budget of  $\text{MoneyAmt}$ ’ to be anchored towards ‘Project’. Since ‘Project’, consequently, becomes the only player of a role involved in the subset constraint to be an anchorpoint, rule 9 causes both of the other fact types to also be anchored towards ‘Project’.

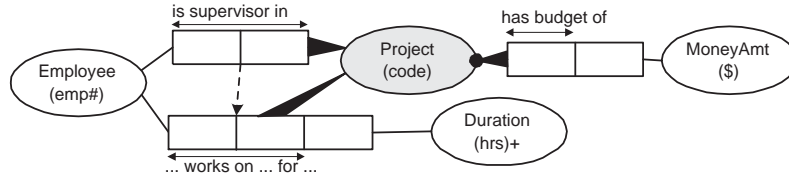


Figure 20: Roles played by major object types and involved in multi-role set constraints are anchored.

### Rule 10 – Joining Roles of Set Constraints

For this rule, we consider each role sequence which is involved in a set constraint and which spans more than one fact type. In these cases, a join condition must be specified (or inferred) to define the manner by which the populations of the involved fact types are related. We call those roles which are involved in the join condition of such a role sequence, the join roles for that role sequence, and define them through the function:  $\text{JoinRoles} : \mathcal{RO}^+ \rightarrow \wp(\mathcal{RO})$ , such that:

$$\text{JoinRoles}(v) \triangleq \{p \in \text{OtherRoles}(v) \mid \exists_{q,r \in v, s \in \text{CoRoles}(r)} [\text{Rel}(q) \neq \text{Rel}(r) \wedge \text{Rel}(p) = \text{Rel}(q) \wedge p \sim s]\}$$

AutoWeight Rule 10 anchors those unanchored fact types, in which only one role is the join role for some set constraint role sequence:

$$\text{JoinSingleSetCon}(p) \triangleq \exists_{v:\text{SetCon}(v)} [p \in \text{JoinRoles}(v)]$$

This join role becomes the anchor:

$$\text{AutoWeight}_{10}(p) \triangleq \text{if } \text{JoinSingleSetCon}(p) \wedge \forall q \in \text{OtherRoles}(p) [\neg \text{JoinSingleSetCon}(q)] \text{ then } 2 \text{ else } 0$$

Figure 21 shows an example to which this rule is applicable. The ‘works for’ predicate is first anchored to ‘Employee’ when Rule 1 (mandatory roles) is fired. The ‘involved in’ predicate is then anchored to ‘Employee’ by the activation of Rule 8 (multi-set constraint with single anchor). Lastly, rule 10 causes the ‘sponsors’ predicate to be anchored to ‘Department’, since the role played by ‘Department’ is the one which is used to join together the *target* role sequence of the subset constraint.

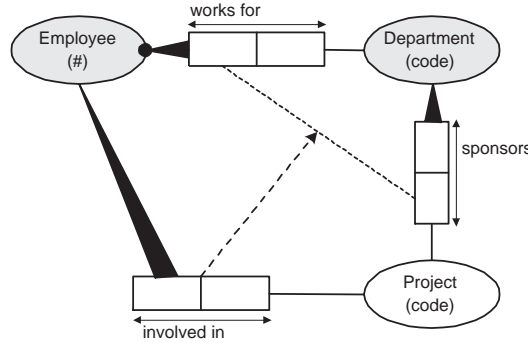


Figure 21: Roles involved in the join condition of a multi-predicate set constraint role sequence are anchored.

### Rule 11 – First Role of Set Constraints

If there is a multi-role, non-implied set constraint (i.e. subset, equality or exclusion constraint) and one of the involved roles has the lowest sequence position within one of the constraint’s role sequence (for its predicate), then this role should become an anchor.

$$\text{AutoWeight}_{11}(p) \triangleq$$

$$\text{if } \exists v: \text{SetCon}(v) \forall s \in \text{CoRoles}(p) \cap \text{Set}(v) [\text{AnchorPoint}_i(\text{Player}(s)) \Rightarrow \text{Pos}(p, v) \leq \text{Pos}(s, v)]$$

$$\text{then } 1 \text{ else } 0$$

This choice is based on the semantics which is derivable from the order in which the modeler chose to initially verbalize the fact type. An example of such a situation is depicted in figure 22.

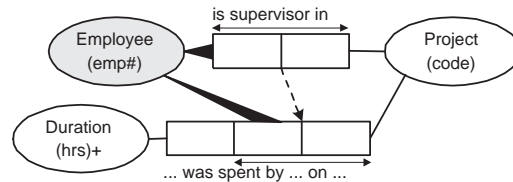


Figure 22: Roles which appear first in set-constraint role sequences are anchored.

Similarly to rules 7 and 8, it is important to consider the case in figure 23, in which the multi-role set constraints contradict each other. In this case, rule 11 would not produce a determinant anchorage for the predicate ‘works on’. We therefore only use this rule on non-anchored fact types which are only involved in one multi-role set constraint. Note, that in figure 23 we denote the role sequence order by specifying the sequence of role numbers that are involved at the ‘source’ of the set constraint.

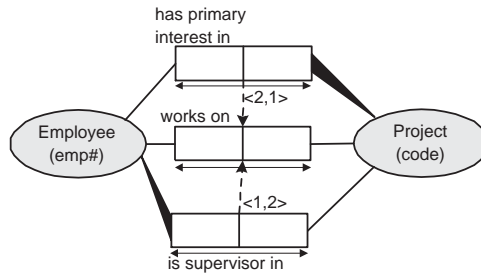


Figure 23: Rule 11 does not consider cases in which a multi-role set constraint contradicts another.

### Rule 12 – First Role of Internal Uniqueness Constraint

Any fact type, which is not already anchored, should be anchored, by default, on the first role that is involved in an internal uniqueness constraint. This choice is based on the semantics which is derived from the order in which the modeler chose to initially verbalize the fact type.

$$\text{AutoWeight}_{1,2}(p) \triangleq \mathbf{if} \text{PosN}(p) = \min \{ \text{PosN}(q) \mid \exists v [\text{Unique}(v) \in \mathcal{U} \wedge q \in v] \} \mathbf{then} 1 \mathbf{else} 0$$

Figure 24 shows an example in which rule 12 is triggered.

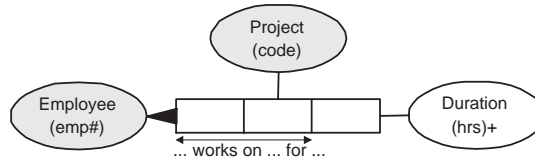


Figure 24: Roles which are positioned in the first 'keyed' position are anchored.

## 4 Deriving Abstraction Levels

When a conceptual schema is abstracted, each progressively higher level of abstraction includes all the most conceptually important components from the previous level. To define an abstraction level, we must therefore first select the major object types and major fact types. We refer to those major types which form the foundation for an abstraction view at level  $i$  as  $\mathcal{KER}_i$  (which will be defined formally below).

Once the kernel of an abstraction level has been calculated, there are still a number of steps which must be performed before the abstraction is complete. Firstly, any fact type predicate in which every role is played by a major object type is included in the abstraction. We do not include these predicates in the kernel itself, because we do not want these fact types to effect the outcome of future abstraction levels. Secondly, we include the identification scheme of all object types which appear in the abstraction. Finally, we restore the connectivity of our abstracted conceptual schema. This involves retaining both the connectivity of subtyping hierarchies, and the connectivity of non-type related object types.

The following subsections formally describe these steps in the abstraction process.

### 4.1 Major Types

In a conceptual schema at a particular level of abstraction ( $i$ ), the set of object types which do not have the lowest conceptual significance are referred to as the major object types ( $\text{MajorOT}_i$ ). Those object types which are of least conceptual significance are referred to as minor. We identify the major object types as the set of object types which have a higher total object type weight ( $\text{OTWeight}$ ) than the minimum total object type weight for the current schema level.

Each object type also has an object type weight (OTWeight) associated with it at a particular abstraction level. The object type weight represents the sum of the weights of those fact type roles which are anchored to it; i.e.:

$$\text{OTWeight}_i(x) \triangleq \sum_{r: \exists y [\text{AnchoredRole}(r, y) \wedge y \sim x]} \text{Weight}_i(r)$$

where  $\text{AnchoredRole} \subseteq \mathcal{FR} \times \mathcal{OB}$  is true when the given fact type role is in  $\mathcal{KER}_i$  and is anchored to an object type which is type related to the given object type; i.e.:

$$\text{AnchoredRole}_i(r, x) \triangleq \text{Player}(r) \sim x \wedge \text{Anchor}_i(r) \wedge \text{Rel}(r) \in \mathcal{KER}_i$$

An object type is considered to be *major* at a particular level of abstraction ( $\text{MajorOT} \subseteq \mathcal{OB}$ ) if its OTWeight is greater than the minimum weight for object types in the kernel at that level. Formally, MajorOT is defined as:

$$\text{MajorOT}_i(x) \triangleq \text{OTWeight}_i(x) > \min \{ \text{OTWeight}_i(y) \mid y \in \mathcal{KER}_i \}$$

The *major fact types* at a particular abstraction level ( $\text{MajorFT}_i \subseteq \mathcal{FT}$ ) are defined as those fact types which bridge between more than one subtyping hierarchy and in which every participant is a major object type at that level:

$$\text{MajorFT}_i(x) \triangleq \forall_{s \in \text{Roles}(x)} [\text{MajorOT}_i(\text{Player}(s))] \wedge \exists_{s, t \in \text{Roles}(x)} [s \not\sim t]$$

## 4.2 Algorithm for Determining Next Abstraction Level

We refer to the set of component types and constraints included in the level  $i$  abstraction view of conceptual data schema  $\mathcal{CS}$  as  $\mathcal{CS}_i$ . In the level 1 conceptual schema,  $\mathcal{CS}_1$  (often abbreviated to  $\mathcal{CS}$ ), all component elements are present. Increasing the level of abstraction will never increase the number of populatable types visible in the conceptual schema:

$$\mathcal{TP}_i \supseteq \mathcal{TP}_{i+1}$$

Of even greater importance, though, increasing the level of abstraction will necessary strictly decrease the number of populatable types within the abstraction kernel:

$$\mathcal{KER}_i \supset \mathcal{KER}_{i+1}$$

As described previously, each progressively higher level of abstraction includes all the most conceptually important components from the kernel of the previous level. For this reason, at each level of abstraction, we include all the major fact types from the previous level, plus all the major object types which participate in at least one of these major fact-types. Formally, we define the abstraction kernel at level  $i + 1$  ( $\mathcal{KER}_{i+1}$ ) as:

$$\mathcal{KER}_{i+1} \triangleq \{x \mid \text{MajorFT}_i(x)\} \cup \{\text{Player}(p) \mid \text{MajorFT}_i(\text{Rel}(p))\}$$

Notice that this definition does not necessarily include every major object type of one abstraction level in the kernel of the next level. The kernel at a particular level of abstraction will only include those major object types which participate in some fact type role contained in the kernel. This conforms to the standard rules of conceptual schema design, as defined in [17].

The definition of  $\mathcal{KER}$  removes all objectified relationship types which neither participate in a major fact-type, nor are major fact types themselves (as shown in figure 25, below).

It is interesting to compare the differences in the way objectified fact types and co-referenced object types are treated. Consider the examples shown in figure 26.

In figure 26(a), the objectified fact type (identified by the participating A and B) is included in the next higher level of abstraction, because it is considered to be a *major fact type*. In figure 26(b), however, the co-referenced object type (AB) is not included in the next higher level of abstraction, because it does not participate in any fact types at this level. AB may, however, be added to  $\mathcal{CS}_{i+1}$  if it is required for connectivity. We justify the difference in treatment of objectified fact types and co-referenced object types by the observation that an objectified fact type can, itself be thought of as a type of abstraction on a co-referenced fact type, which must, itself be ‘unwrapped’ [9].

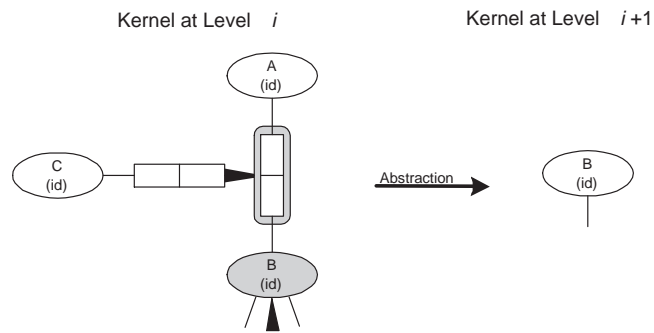


Figure 25: Removal of objectified relationship types

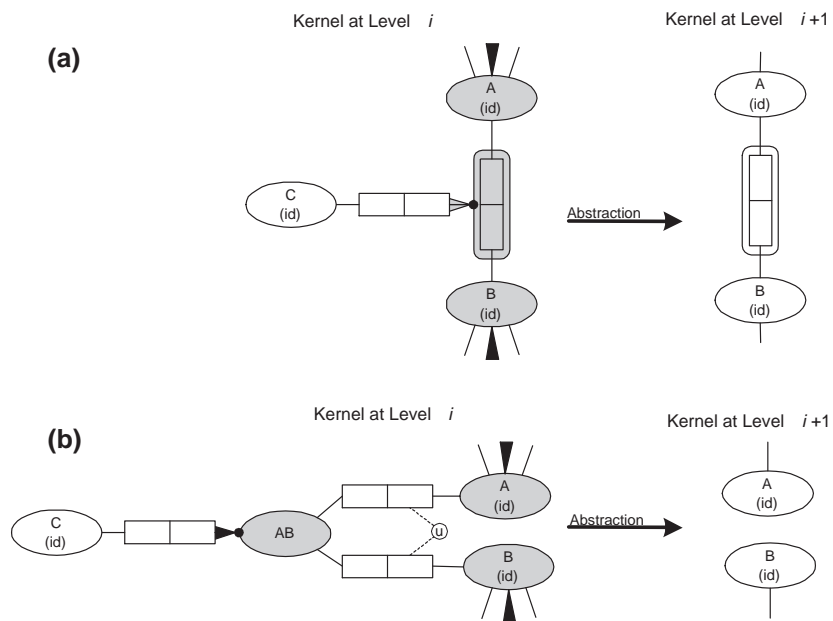


Figure 26: (a) Objectified Fact Type (b) Co-Referenced Object Type



### 4.3 Ring Fact Types

At this stage, the kernel only contains fact types which bridge between subtyping hierarchies. The kernel does not retain those (ring) fact types for which every participant is a type-related object type because we do not want these fact type to perpetually cause their player to be an anchor point. The user, however, is probably interested in viewing all fact types which are played entirely by major object types at the previous level. This includes the ring fact types. To this end, we therefore apply the procedure `AddRingFTs` to the types in the kernel.

```
AddRingFTs :  $\wp(\mathcal{OB}) \rightarrow \wp(\mathcal{OB})$ 
AddRingFTs(Types)
BEGIN
  RETURN Types  $\cup \{r \in \mathcal{FT} \mid \forall_{s,t \in \text{Player}(r)} [s, t \in \text{Types} \wedge s \sim t]\}$ ;
END AddRingFTs;
```

### 4.4 Object Type Identification Schemes

Since the identification scheme of an object type is often important for its understanding, we ensure that the identification scheme of every object type is included in each abstraction level. To this end, we define that relation,  $\text{IsIdRel} \subseteq \mathcal{RL} \times \wp(\mathcal{TP})$ , which is true when the given relationship type is involved the primary identification of some object type in the given set of types.

$$\text{IsIdRel}(r, z) \triangleq \exists_{x \in z} [r \in \text{PIdRels}(x) - z]$$

The function `IdentifiedSchema` :  $\wp(\mathcal{TP}) \rightarrow \wp(\mathcal{TP})$  takes the types in the current abstraction level and adds to them, those types that are required to identify the input types.

```
IdentifiedSchema(Types)
VAR
  r :  $\mathcal{FT}$ ;
BEGIN
  WHILE  $\exists_r [r \notin \text{Types} \wedge \text{IsIdRel}(r, \text{Types})]$  DO
    LET r BE SUCH THAT  $r \notin \text{Types} \wedge \text{IsIdRel}(r, \text{Types})$ ;
    Types +=  $\{r\} \cup \{\text{Player}(s) \mid s \in \text{Roles}(r)\}$ ;
  END WHILE;
  RETURN Types;
END IdentifiedSchema.
```

It is interesting to consider the effect that this algorithm has on the subtyping hierarchies of the abstracted schema. In figure 23, the only major object types in  $\mathcal{KER}_i$ , which participate in a major fact type, are C and F. C and F are therefore the only object types to appear in  $\mathcal{KER}_{i+1}$ . Since C inherits its identification from its (indirect) supertype A, however, A is included in `IdentifiedSchema`( $\mathcal{KER}_{i+1}$ ) as the player of C's identifying relationship type.

Notice that the subtyping arrows in the various schema fragments adapt automatically to the set of object types included in the diagram. This is possible because of the fact that subtyping relationships are inherently transitive, with only the non-implied arrows being displayed on the diagram.

### 4.5 Connectivity

Since we wish to retain the connectivity of our conceptual schema throughout each level of abstraction, we must define the concept of connectivity. We begin by defining a connected path ( $\mathcal{PA}$ ) through a conceptual schema. A path is a sequence of types in which each element (except the first) is either a relationship type of which the previous conceptual type is a player, or is one of the object types which plays the previous conceptual type. Note that a path does not necessarily define a unique traversal through a relationship type, since an object type may play more than one role in the same relationship type.

$$\mathcal{PA} \triangleq \{x \in \mathcal{TP}^+ \mid \forall_{1 \leq i \leq |x|} [\text{Connected}(x[i], x[i+1])]\}$$

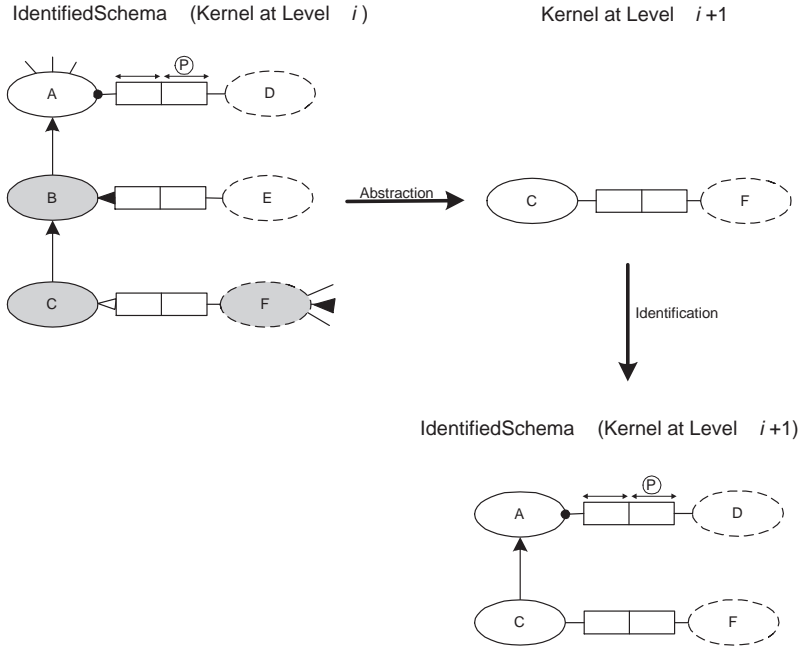


Figure 27: The primary identification of subtype C is inherited from supertype A

where `Connected` identifies whether or not the given conceptual types are connected (in that one participates in the other):

$$\text{Connected}(y, z) \triangleq \exists_{p \in \text{Roles}(y)} [\text{Player}(p) \sim z] \vee \exists_{p \in \text{Roles}(z)} [\text{Player}(p) \sim y]$$

The predicate `PathBetween`  $\subseteq \mathcal{PA} \times \mathcal{TP} \times \mathcal{TP}$ , holds exactly when the given path exists, starting at conceptual type  $x$  and ending at conceptual type  $y$ .

$$\text{PathBetween}(v, x, y) \triangleq v[1] = x \wedge v[|v|] = y$$

`ShortestPaths`( $x, y$ ) returns the set of paths which start at type  $x$  and end at type  $y$ , which contain the least number of conceptual types in between (`ShortestPaths` :  $\mathcal{TP} \times \mathcal{TP} \rightarrow \wp(\mathcal{PA})$ ).

$$\text{ShortestPaths}(x, y) \triangleq \{v \in \mathcal{PA} \mid \text{PathBetween}(v, x, y) \wedge |v| = \min \{|w| \mid \text{PathBetween}(w, x, y)\}\}$$

Another type of connectivity which is useful to maintain is the connectivity of subtyping hierarchies. To this end, each set of unconnected, type-related types in `IdentifiedSchema`( $\mathcal{KER}_{i+1}$ ) are reconnected via the lowest common supertype which has its own identification scheme. The notion of a *lowest common identified supertype* is therefore introduced. A set of object types can actually have more than one *lowest common identified supertype*, in the case in which the subtyping hierarchy forms a lattice. The function `LowestCIdObjs` :  $\wp(\mathcal{TP}) \rightarrow \wp(\mathcal{TP})$  is defined on a set of type-related object types as:

$$\text{LowestCIdObjs}(w) \triangleq \{x \in \text{CommonIDSup}(w) \mid \neg \exists_{y \in \text{CommonIDSup}(w)} [y \text{ SubOf } x]\}$$

where `CommonIDSup`( $w$ )  $\triangleq \text{CommonSup}(w) \cap \text{Players}(\text{OtherRoles}(\text{PIdRoles}(w)))$ . Notice that the only case in which the *lowest common identified supertype* for a set of object types will be the same as the *lowest common supertype*, defined previously, is when the lowest common supertype has an identification scheme directly attached to it.

The function `ConnectSchema` is used to add those types which are required to connect the given set of types by means of subtyping hierarchies and shortest paths.

```

ConnectSchema :  $\wp(\mathcal{TP}) \rightarrow \wp(\mathcal{TP})$ 
ConnectSchema(Types)
  # PostCondition:  $\forall x, y \in \text{Types} [\text{PathBetween}(x, y)]$ 
BEGIN
  1. WHILE  $\exists x, y \in \text{Types} [x \sim y \wedge \neg \exists z \in \text{Types} [z \in \text{CommonSup}(\{x, y\})]]$  DO
    Types += LowestCldObjs( $\{s \mid s \sim x\}$ );
  2. WHILE  $\exists x, y \in \text{Types} [\neg \text{PathBetween}(x, y)]$  DO
    Types +=  $\bigcup_{c \in \text{ShortestPath}(x, y)} \text{Set}(c)$ ;
    Types := IdentifiedSchema(Types);
  END 2;
END 1;
RETURN Types;
END ConnectSchema;

```

The first loop (marked 1.) ensures that all type related object types are connected via a subtyping hierarchy in the abstraction schema. The second loop (marked 2.) ensures that the non-type related object types are connected via relationship types in the resulting set of types. Notice that IdentifiedSchema is reapplied during ConnectSchema to ensure that any newly added types are also identifiable in the abstraction schema.

Figure 28 shows an example in which IdentifiedSchema( $\mathcal{KER}_{i+1}$ ) is disconnected, despite the fact that its components contain type related object types. In order to connect the two schema fragments in IdentifiedSchema( $\mathcal{KER}_{i+1}$ ), ConnectSchema adds A to the set of included types. Object type A represents the lowest common supertype of G and C which has its own identification scheme. In this case, C and G have their own identification scheme. However, in order to retain the notion that the instances of C and G come from a common domain, we include in the abstraction the lowest common supertype which is directly attached to a unique identification scheme for this domain.

## 4.6 Complete Abstraction of Conceptual Schema

The complete set of abstraction levels for a given conceptual schema,  $\mathcal{CS}$ , can now be defined. The level  $i$  abstraction view,  $\mathcal{CS}_i$ , of a given conceptual schema,  $\mathcal{CS}$ , includes the types in  $\mathcal{KER}_i$ , the types required for connectivity and identification, and all constraints from the original schema in which these included types are involved. We, therefore, formally define  $\mathcal{CS}_i$  as follows:

$$\begin{aligned} \mathcal{CS}_1 &\triangleq \langle \mathcal{IS}, \text{Mand}, \text{Unique}, \text{PUUnique}, \text{Frequency}, \text{SetCon}, \text{Weight} \rangle \\ \mathcal{CS}_i &\triangleq \langle \mathcal{IS}_i, \text{Mand}_i, \text{Unique}_i, \text{PUUnique}_i, \text{Frequency}_i, \text{SetCon}_i, \text{Weight}_i \rangle \text{ for } i > 1 \end{aligned}$$

where

$$\begin{aligned} \mathcal{IS}_i &\triangleq \langle \mathcal{TP}_i, \mathcal{RO}_i, \text{Roles}_i, \text{PosN}, \text{SubOf}_i, \text{Player} \rangle \\ \mathcal{TP}_i &\triangleq \text{ConnectSchema}(\text{IdentifiedSchema}(\text{AddRingFTs}(\mathcal{KER}_i))) \\ \mathcal{RO}_i &\triangleq \{p \in \mathcal{RO} \mid \text{Rel}(p) \in \mathcal{TP}_i\} \\ \text{SubOf}_i &\triangleq \{\langle x, y \rangle \in \mathcal{TP}_i \times \mathcal{TP}_i \mid x \text{ SubOf}_{i-1} y\} \\ \text{Mand}_i &\triangleq \{x \in \wp(\mathcal{RO}_i) \mid \text{Mand}(x) \wedge \exists p \in \mathcal{RO}_i \forall q \in x [p \sim q \wedge p \notin x]\} \\ \text{Unique}_i &\triangleq \{x \in \wp(\mathcal{RO}_i) \mid \text{Unique}(x)\} \\ \text{Frequency}_i &\triangleq \{x \in \wp(\mathcal{RO}_i) \mid \text{Frequency}(x)\} \\ \text{SetCon}_i &\triangleq \{\langle v, w \rangle \in \mathcal{RO}_i^+ \times \mathcal{RO}_i^+ \mid \text{SetCon}(v, w)\} \end{aligned}$$

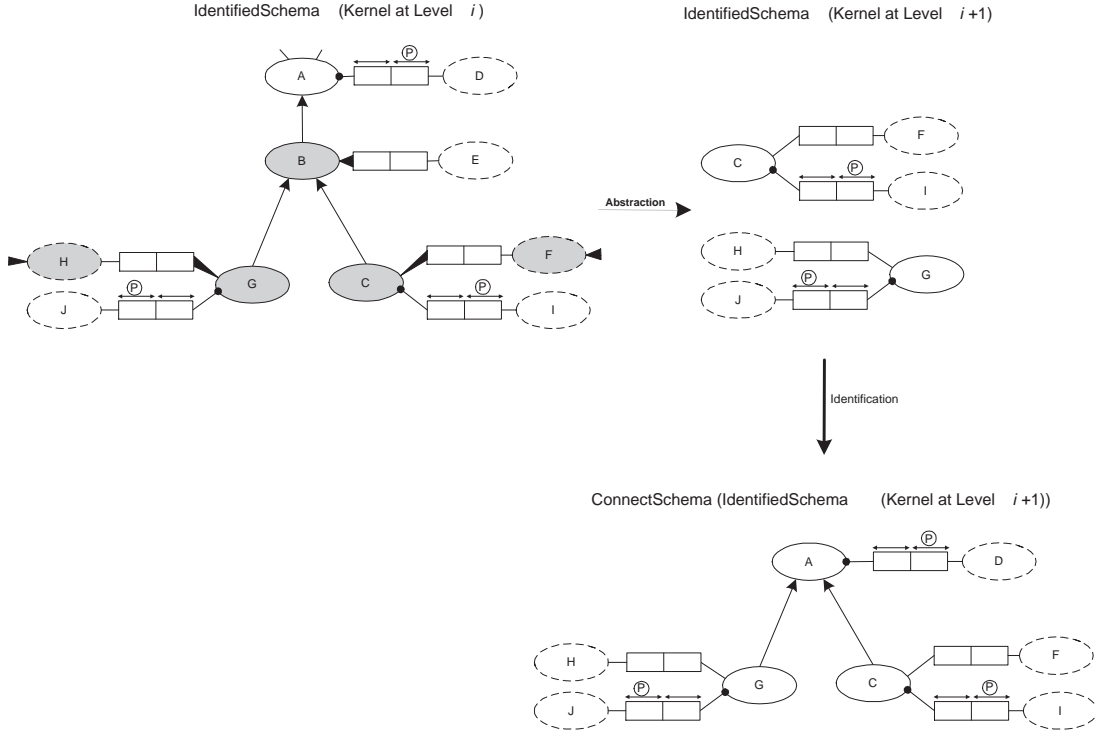


Figure 28: Connecting the kernel by means of the subtype hierarchy

## 4.7 Relation to Clusters

The definitions so-far allow us to take a flat ORM conceptual schema and derive a number of abstraction layers for this flat ORM schema. Each of these layers is still essentially a flat subsection of the original ORM schema. Therefore, we now introduce the glue that actually holds these layers of abstraction together.

The idea is to view each major object type as becoming the centre for a clustering of surrounding minor object types. As a result, the object types in each  $\mathcal{KER}_i$  are clusterings of types from  $\mathcal{KER}_{i-1}$ . This idea of using clustering as a binding mechanism for abstraction layers for ORM schemas was proposed previously in [6], [7] and [5]. In [9] a possible formalization of the clustering mechanism is presented.

In this subsection, we show how to derive a clustering of minor object types for each major object type. The presented style of clustering conforms to the requirements given in [9]. This means that when applying the abstraction algorithm discussed in this article, together with the clustering mechanism presented below, a three dimensional ORM schema results that is in line with the 3-Dimensional Conceptual Modeling Kernel as proposed in [9].

The clustering mechanism is defined as a set of derivation rules. An actual clustering is given as a function  $\text{Cluster} : \mathbb{N} \times \mathcal{OB} \rightarrow \wp(\mathcal{TP})$ . The intuition is that if  $x \in \text{Cluster}(i, y)$ , then at abstraction level  $i$  type  $x$  has been grouped into the cluster surrounding object type  $y$ .

The first derivation rule clusters all fact types which have disappeared since the last level of abstraction towards the object type to which they were anchored.

$$\text{[CL1]} \quad \text{Anchor}_{i-1}(r) \wedge \text{Rel}(r) \in \mathcal{KER}'_i \vdash \text{Rel}(r) \in \text{Cluster}(i, \text{Player}(r))$$

where  $\mathcal{KER}'_i \triangleq \mathcal{KER}_{i-1} - \mathcal{KER}_i$  for  $i > 1$ .

Object types which participate in any relationship type included in a cluster should, obviously, be included in the same cluster:

$$\text{[CL2]} \quad x \in \text{Cluster}(i, c) \wedge y \in \text{Players}(x) \vdash y \in \text{Cluster}(i, c)$$

Please note that an object type could occur in more than one clustering if it is involved in relationship types anchored to different major object types. As a result, the clustering is not a partition of the types.

The following two rules are concerned with subtyping. If an object type in a type hierarchy is removed from the kernel (i.e. it is in  $\mathcal{KER}'_i$ ), we must still cluster those fact types that were anchored to it. We anchor such fact types towards the lowest supertype which remains in the kernel (LowestKernelSup).

**[CL3]**  $x \in \mathcal{KER}'_i \wedge y \text{ LowestKernelSup } x \wedge r \text{ AnchoredTo}_{i-1} x \vdash \text{Rel}(r) \in \text{Cluster}(i, y)$

where  $y \text{ LowestKernelSup } x$  indicates that  $y$  is the lowest supertype of  $x$  which remains in the abstraction kernel. That is:

$$y \text{ LowestKernelSup } x \triangleq x \text{ SubOf } y \wedge y \in \mathcal{KER}_i \wedge \neg \exists_{z \in \mathcal{KER}_i} [x \text{ SubOf } z \text{ SubOf } y]$$

If no supertype remains in the kernel, however, the relationship types anchored towards a disappearing subtype can be clustered towards the nearest subtype which remains in the kernel.

**[CL4]**  $x \in \mathcal{KER}'_i \wedge \neg \exists_z [z \text{ LowestKernelSup } x] \wedge y \text{ HighestKernelSub } x \wedge r \text{ AnchoredTo}_{i-1} x \vdash \text{Rel}(r) \in \text{Cluster}(i, y)$

where  $y \text{ HighestKernelSub } x$  indicates that  $y$  is the highest subtype of  $x$  which remains in the abstraction kernel. That is:

$$y \text{ HighestKernelSub } x \triangleq y \text{ SubOf } x \wedge y \in \mathcal{KER}_i \wedge \neg \exists_{z \in \mathcal{KER}_i} [y \text{ SubOf } z \text{ SubOf } x]$$

The remaining derivation rules are completeness rules on clusters. Clustered types are inherited between layers of abstraction. So we have:

**[CL5]**  $x \in \text{Cluster}(i, c) \vdash x \in \text{Cluster}(i+1, c)$

The reference types needed to identify any of the types in a cluster are also included:

**[CL6]**  $x \in \text{Cluster}(i, c) \wedge y \in \text{PldRels}(x) \vdash y \in \text{Cluster}(i, c)$

The above definition of cluster is a ‘maximally complete’ one. However, when displaying clusters to a user, for example, one may chose to only show those clustered types which are part of  $\mathcal{KER}'_i$ . That is:

$$\text{Cluster}(i, x) - \text{Cluster}(i-1, x)$$

It may also be decided to only show the clusters for those types which appear in  $\mathcal{KER}_i$ , and ignore the clusters for those object types which were major at the previous level, but do not participate in a fact types in  $\mathcal{KER}_i$ . Choices like this are up to the designer of the actual abstraction tool and often depend solely upon the purpose for which the abstraction and clustering was created.

## 5 Case Study

Now that we have developed a theory for the creation of abstractions for a conceptual schema, it is time to study the effect that such a mechanism has on an application example. Applying WeightSchema to the conceptual schema shown in figure 6, we achieve the anchored schema shown below in figure 29. As in previous examples, we have shaded the major object types and indicated the anchors by an arrowed role connector line. For the purposes of our example application, we have also included the Weight assigned to each anchor (e.g. “.9.”). This will hopefully help the reader to retrace the AutoWeight rules that have been fired to achieve this result.

Applying WeightSchema to our application example of figure 6 helps us identify the major object types in the Universe of Discourse. In this case, they are ‘Motel’, ‘Committee’, ‘Institution’, ‘Country’, ‘Request’, ‘Person’, ‘Paper’, ‘Accepted Paper’, ‘Paper Slot’, ‘Room’, ‘Lab or Lecture Room’, ‘Lecture Room’ and ‘Laboratory’. This follows our own intuition of the most ‘conceptually important object types’.

It is important to notice that the relationship types ‘Room is in Building’ and ‘Room has Room#’ are not anchored. This is because they are part of the primary identification scheme for ‘Room’, and are therefore reference types. Only fact types are anchored.

Among the anchored fact types are ‘*Person chairs Committee*’, which is anchored towards ‘*Committee*’. The maximum frequency of 2 on the role played by ‘*Committee*’ causes *AutoWeight* rule 4 to assign a *Weight* of 7 to this role. *AutoWeight* rule 8 is responsible for anchoring the fact type ‘*Person presents Accepted Paper*’ towards ‘*Accepted Paper*’, due to the fact that it is associated via a set constraint to a fact type already anchored towards ‘*Paper*’.

Figure 30 shows the kernel types which form the foundation of the second level of abstraction for our example conceptual schema. Notice that no ring fact types or identification schemes are included in the kernel, and that the kernel is actually disconnected.

When we apply *AddRingFTs*, *IdentifiedSchema* and *ConnectSchema* to the kernel types in figure 30, add the constraints that are still relevant and re-apply *WeightSchema*, we achieve the complete weighted, second level abstraction schema shown in figure 31.

Notice that the major object types of  $\mathcal{CS}_2$  are ‘*Request*’, ‘*Person*’, ‘*Institution*’, ‘*Paper*’, ‘*Accepted Paper*’ and ‘*Paper Slot*’. Because these object types are major at both the first and second level of abstraction, we consider them to be more ‘conceptually important’ than those object types which are only major at the first level of abstraction. In fact, we gauge an object type’s degree of majorness ( $\text{DegreeMajor} : \mathcal{OB} \rightarrow \mathbb{N}$ ) by calculating the highest level of abstraction at which that object type is major. We define:

$$\text{DegreeMajor}(x) \triangleq \max(\{i \mid \text{MajorOT}_i(x)\} \cup \{0\})$$

and we know that:

$$\text{MajorOT}_i(x) \Rightarrow \text{DegreeMajor}(x) \geq i$$

For example, so far we know that:

$$\text{DegreeMajor}(\textit{Rating}') = 0; \text{DegreeMajor}(\textit{Motel}') = 1 \text{ and } \text{DegreeMajor}(\textit{Person}') \geq 2$$

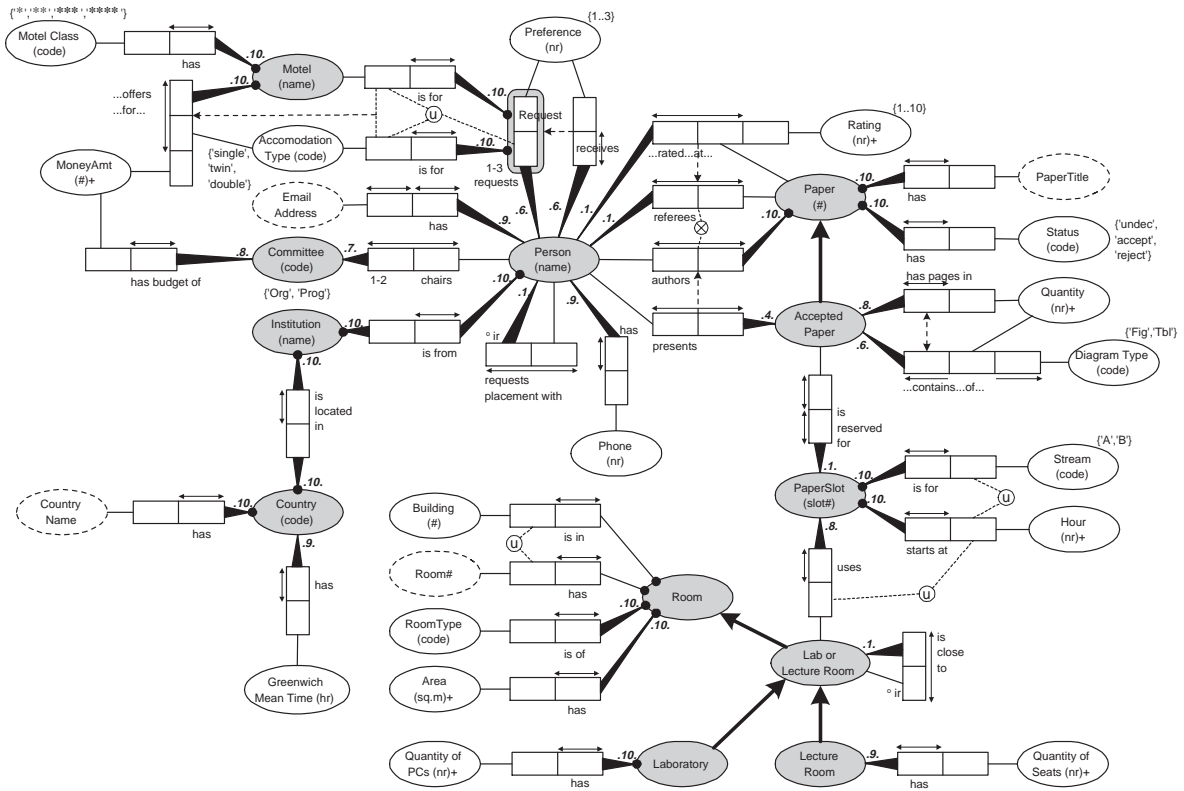
Conceptual importance, or conceptual relevance (as indicated by *DegreeMajor*) plays a key role in a number of areas. For example, in computer supported query formulation conceptual importance is used to help select between alternative interpretations of queries ([4], [21], [22]).

It is important to understand how the schema abstraction in figure 30 was obtained. The object types ‘*Room*’, ‘*Building*’, ‘*Room#*’ and ‘*Preference*’ were all added by the *IdentifiedSchema* procedure, because they are used in the identification of some kernel object type. The fact types ‘*Lab or Lecture Room is close to Lab or Lecture Room*’ and ‘*Person requests placement with Person*’ were added during *AddRingFTs*, and are not actually part of  $\mathcal{KER}_2$ . This explains why we do not consider these fact types to be anchored.

There are a few interesting things to observe with respect to the anchorage of  $\mathcal{CS}_2$ . Firstly, notice that ‘*Institution is located in Country*’, in contrast to  $\mathcal{CS}_1$ , is not anchored towards ‘*Country*’. This is because at abstraction level 2, this role becomes implied mandatory. Secondly, the anchor on ‘*Person chairs Committee*’ was on the role played by ‘*Committee*’ in  $\mathcal{CS}_1$ , but has now moved to the role played by ‘*Person*’. This is because ‘*Committee*’ has become a leaf object type, causing the role played by ‘*Person*’ to gain a new weight of .9.. The fact type ‘*Person referees Paper*’ also has a change in anchorage. In the previous level of abstraction, it was anchored by Rule 12. Since it is now only associated with a single set constraint, however, Rule 8 now triggers a weight of .4. on the role played by *Paper*. Lastly, notice that the weight of the anchor on ‘*PaperSlot uses Lab or Lecture Room*’ has increased from .8. to .9. because ‘*Lab or Lecture Room*’ is now a leaf object type, triggering rule 3.

Figure 32 depicts a different view of our second level abstraction. In figure 32, we explicitly show the clusterings that have occurred during the abstraction process. We have chosen to represent those object types which are repeated in more than one cluster by surrounding them with a second ellipse; and have shown only those constraints which are completely within or completely external to a clustering. It is particularly interesting to observe the subtype clustering that has occurred around the object type ‘*Lab or Lecture Room*’.

Taking things one step further, we can easily extend our results from figure 31 and 32 to show a corresponding Entity Relationship (ER) representation of the application. Figure 33 shows this ER view. There are many notations used for ER modeling. The one presented here uses rounded rectangles to represent entities, named lines to represent relationship types, crow feet to indicate that the opposite entity can play that role many times, a double rectangle to represent a ‘weak entity type’, and the letters ‘ID’ placed on its identifying relationship. Attributes are not shown in this diagram.



**SUBTYPE DEFINITIONS**  
 each Accepted Paper is a Paper that has Status 'accept'  
 each LabOrLecRoom is a Room that is of RoomType in {'lab','lec'}  
 each Laboratory is a Room that is of RoomType 'lab'  
 each Lecture Room is a Room that is of RoomType 'lec'

Figure 29: WeightSchema( $\mathcal{CS}_1$ ): Abstraction Level 1

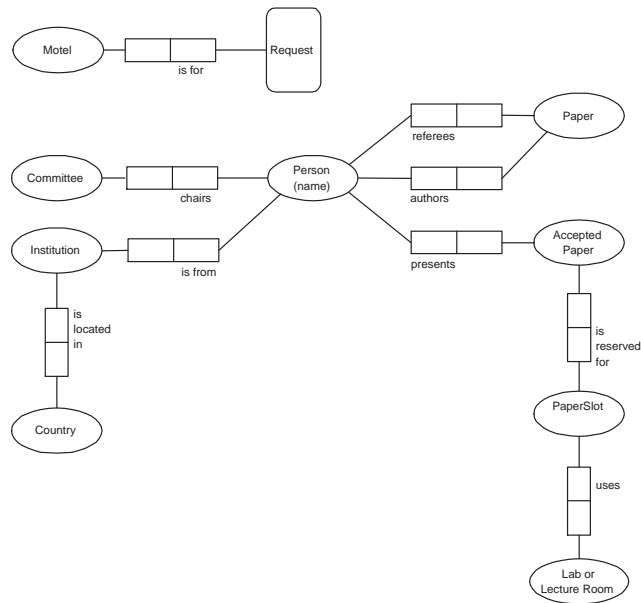


Figure 30:  $\mathcal{KER}_2$ : The kernel for Abstraction Level 2

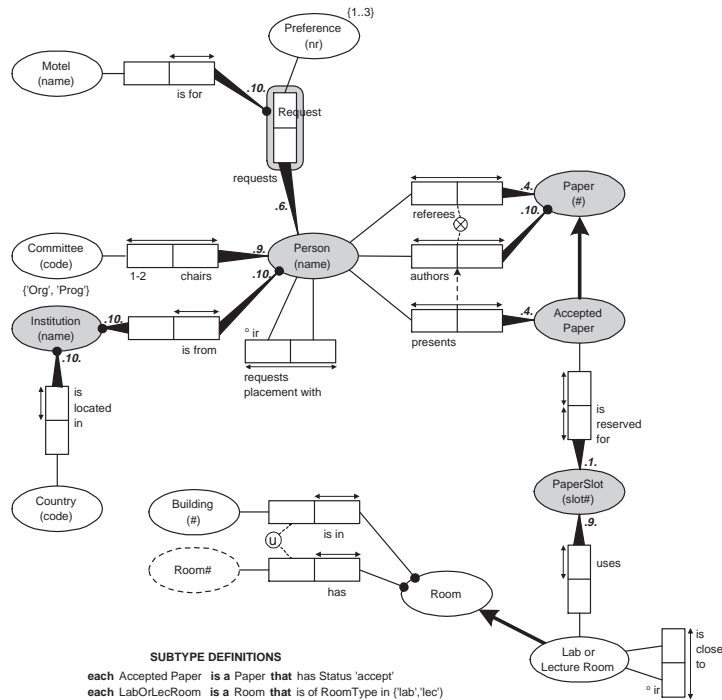


Figure 31: WeightSchema( $\mathcal{CS}_2$ ): Abstraction Level 2

It is important to realize that the version of ER used above allows multi-valued and composite attributes. Making this assumption allows us to achieve an intuitive overview of the original ORM diagram using ER notation. We also allow relationships to have attributes. For example, the relationship ‘*Person rated Paper*’ has the attribute ‘*Rating*’. Notice that relationships, such as this one, which have attributes do not appear in the ORM abstraction since some of their participants are minor. Also note that ‘*Request*’ is represented as a weak entity because its identification scheme involves both an attribute (‘*Preference*’) and a relationship to an entity (‘*Person*’).

We now take the second level abstraction shown in figure 31 and abstract again. Figure 34 illustrates the third level of abstraction,  $\mathcal{CS}_3$ . When we apply WeightSchema to  $\mathcal{CS}_3$ , the only changes in anchorage that can be seen from the previous abstraction are in ‘*Paper Slot is reserved for Accepted Paper*’ (as ‘*Paper Slot*’ becomes a leaf) and ‘*Person is from Institution*’ (as the role played by ‘*Institution*’ becomes implied mandatory). We can now determine that the degree of majorness for ‘*Institution*’ and ‘*Request*’ is 2, and for ‘*Person*’, ‘*Paper*’, ‘*Accepted Paper*’ and ‘*Paper Slot*’ is greater than, or equal to three.

The highest level of abstraction that can be reached for our example application is four. Figure 35 shows  $\mathcal{CS}_4$ . No higher level of abstraction can be reached, because there are no major fact types in  $\mathcal{CS}_4$ , and therefore no conceptual types would be present at a higher level of abstraction. We can now conclude that the most conceptually important object type in our Universe of Discourse is ‘*Paper*’ (and ‘*Accepted Paper*’), with a degree of majorness equal to four.

## 6 Conclusions

In this article, we have presented an algorithm to derive layers of abstraction for a given flat conceptual schema. The cornerstone of this abstraction algorithm is the notion of a major object type. We have defined a prioritized set of derivation rules to assist in the selection of these major object types at each level of abstraction. In comparison to other approaches which determine the major object types of a conceptual schema, our approach considers more of the semantics that are hidden in the constraints and verbalizations. Alternative approaches have instead relied more heavily on user input.



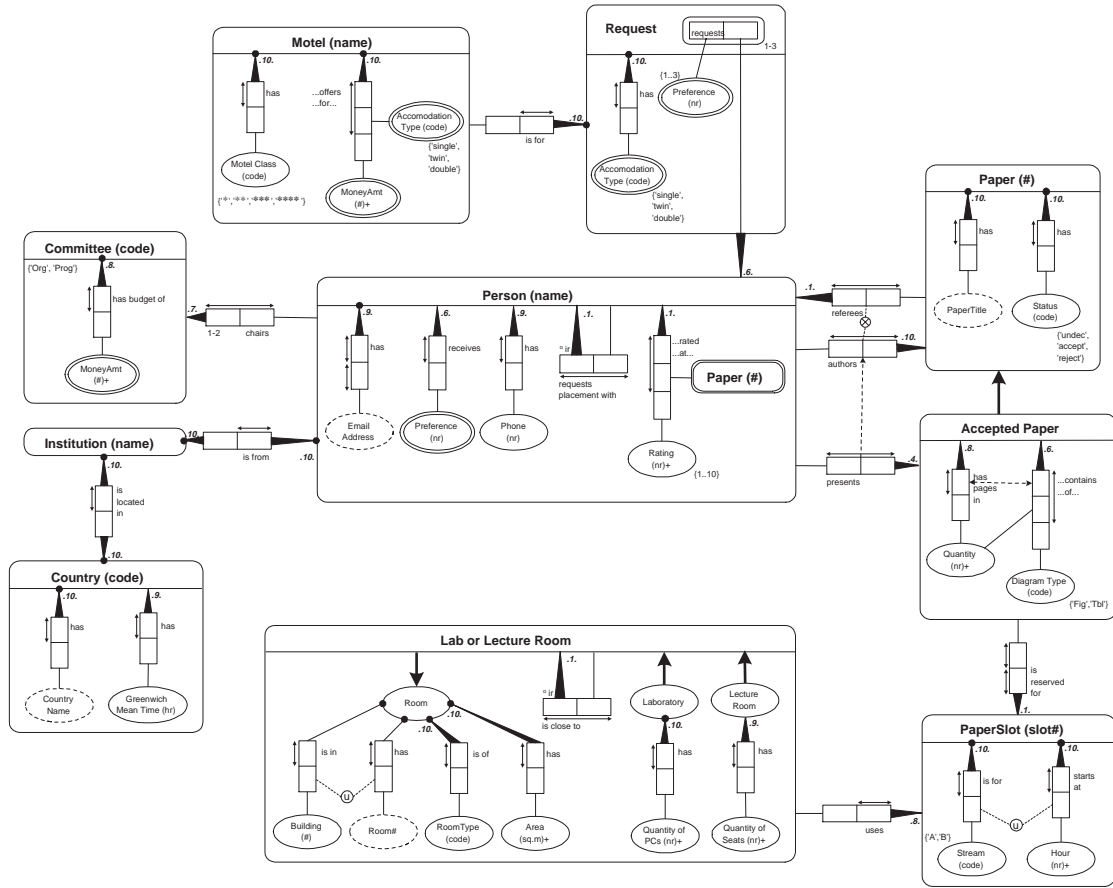


Figure 32: Abstraction Level 2 showing the clusterings that have occurred

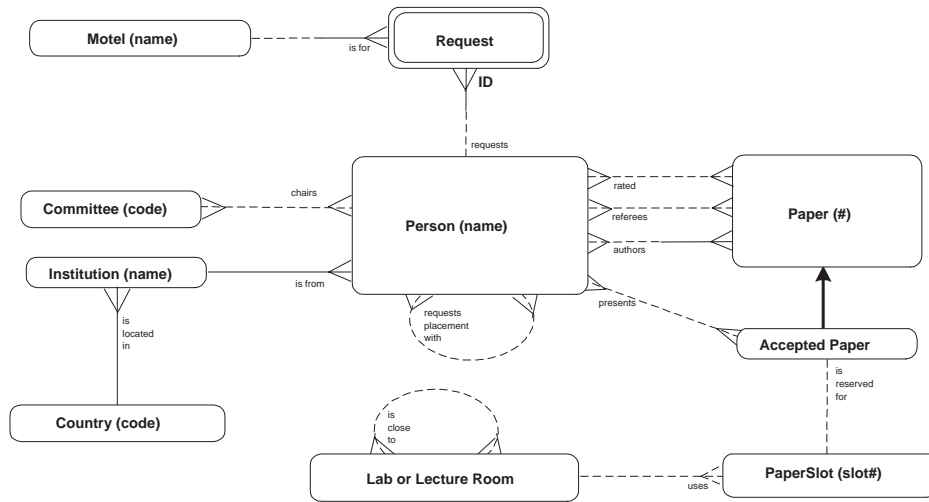


Figure 33: An ER View of the example application

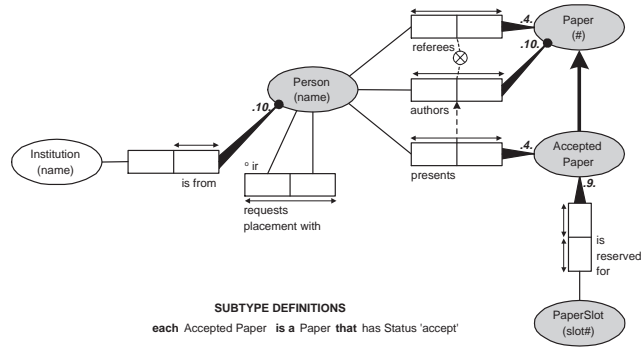


Figure 34: WeightSchema( $\mathcal{CS}_3$ ): Abstraction Level 3

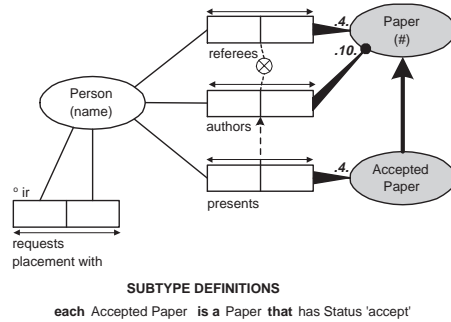


Figure 35: WeightSchema( $\mathcal{CS}_4$ ): Abstraction Level 4

The paper presents an iterative method for using the major object types to determine the kernel types in each subsequent abstraction level. This kernel is then embellished with additional conceptual types to enhance its comprehensibility, by incorporating identification schemes, ring predicates and connectivity. The result of these processes, together with the constraints applicable to the included types, form the resulting abstraction schema, which is demonstrated on a concrete case study in section 5.

In addition to this, we have also shown how the resulting layers of abstraction provide a three dimensional view of the underlying flat conceptual schema, where object types at each level of abstraction can be regarded as clusterings of object types from a lower level of abstraction.

Future plans include the tuning of priorities and weights in the derivation rules, which currently reflect the intuition of the authors, to be based on empirical evidence gained through concrete testing. Because the derivation rule approach taken in this paper allows a very modular approach to be taken in its implementation, the tuning of the priority and weighting values can be performed locally.

Other topics for future research include investigating how bottom up abstraction (as described in this paper) and top down abstraction mechanisms can complement each other in a single conceptual schema design procedure; and investigating how the algorithms presented in this paper impact on other types of conceptual schema abstraction.

## Acknowledgments

We would like to thank the anonymous referees for their comments and suggestions, which have lead to improvements of the original article. Furthermore, we would particularly like to thank A.H.M. ter Hofstede for his comments the first drafts of this article

## References

- [1] C. Batini, S. Ceri, and S.B. Navathe. *Conceptual Database Design - An Entity-Relationship Approach*. Benjamin Cummings, Redwood City, California, 1992.
- [2] 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):471–495, October 1991.
- [3] 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):213–235, 1995.
- [4] C.A.J. Burgers, H.A. Proper, and Th.P. van der Weide. An Information System organized as Stratified Hypermedia. In N. Prakash, editor, *CISM94, International Conference on Information Systems and Management of Data*, pages 159–183, Madras, India, October 1994.
- [5] L.J. Campbell. Adding a New Dimension to Flat Conceptual Modelling. In T.A. Halpin and R. Meersman, editors, *Proceedings of the First International Conference on Object-Role Modelling (ORM-1)*, pages 294–309, Magnetic Island, Australia, July 1994.
- [6] L.J. Campbell and T.A. Halpin. Automated Support for Conceptual to External Mapping. In S. Brinkkemper and F. Harmsen, editors, *Proceedings of the Fourth Workshop on the Next Generation of CASE Tools*, pages 35–51, Paris, France, June 1993.
- [7] L.J. Campbell and T.A. Halpin. Abstraction Techniques for Conceptual Schemas. In R. Sacks-Davis, editor, *Proceedings of the 5th Australasian Database Conference*, volume 16, pages 374–388, Christchurch, New Zealand, January 1994. Global Publications Services.
- [8] C.R. Carlson, W. Ji, and A.K. Arora. The Nested Entity-Relationship Model. In F.H. Lochovsky, editor, *Proceedings of the Eight International Conference on Entity-Relationship Approach, Entity-Relationship Approach to Database Design and Querying*, pages 43–57, Toronto, Canada, 1990. Elsevier Science Publishers.
- [9] P.N. Creasy and H.A. Proper. A Generic Model for 3-Dimensional Conceptual Modelling. *Data & Knowledge Engineering*, 20(2):119–162, 1996.
- [10] B. Czejdo and D.W. Embley. View Specification and Manipulation for a Semantic Data Model. *Information Systems*, 16(4):28–44, 1991.
- [11] O.M.F. De Troyer. A Logical Formalization of the Binary Object-Role Model. In T.A. Halpin and R. Meersman, editors, *Proceedings of the First International Conference on Object-Role Modelling (ORM-1)*, pages 28–44, Magnetic Island, Australia, July 1994.
- [12] R. Elmasri and S.B. Navathe. *Fundamentals of Database Systems*. Benjamin Cummings, Redwood City, California, 1994. Second Edition.
- [13] P. Feldman and D. Miller. Entity Model Clustering: Structuring a Data Model by Abstractions. *The Computer Journal*, 29(4):348–360, 1986.
- [14] C. Francalanci and B. Pernici. Abstraction Levels for Entity-Relationship Schemas. In P. Loucopoulos, editor, *Proceedings of the Fourth International Conference CAiSE'92 on Advanced Information Systems Engineering*, volume 593 of *Lecture Notes in Computer Science*, pages 456–473, Manchester, United Kingdom, 1992. Springer-Verlag.
- [15] M. Ghandi, E.L. Robertson, and D.V. Gucht. Leveled Entity Relationship Model. In P. Loucopoulos, editor, *Proceedings of the Fourth International Conference CAiSE'92 on Advanced Information Systems Engineering*, volume 593 of *Lecture Notes in Computer Science*, pages 420–436, Manchester, United Kingdom, May 1992. Springer-Verlag.
- [16] T.A. Halpin. *A logical analysis of information systems: static aspects of the data-oriented perspective*. PhD thesis, University of Queensland, Brisbane, Australia, 1989.

- [17] T.A. Halpin. *Conceptual Schema and Relational Database Design*. Prentice-Hall, Sydney, Australia, 2nd edition, 1995.
- [18] T.A. Halpin and H.A. Proper. Subtyping and Polymorphism in Object-Role Modelling. *Data & Knowledge Engineering*, 15:251–281, 1995.
- [19] A.H.M. ter Hofstede. *Information Modelling in Data Intensive Domains*. PhD thesis, University of Nijmegen, Nijmegen, The Netherlands, 1993.
- [20] 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):489–523, October 1993.
- [21] A.H.M. ter Hofstede, H.A. Proper, and Th.P. van der Weide. Computer Supported Query Formulation in an Evolving Context. In R. Sacks-Davis and J. Zobel, editors, *Proceedings of the Sixth Australasian Database Conference, ADC'95*, volume 17(2) of *Australian Computer Science Communications*, pages 188–202, Adelaide, Australia, January 1995.
- [22] A.H.M. ter Hofstede, H.A. Proper, and Th.P. van der Weide. Query formulation as an information retrieval problem. *The Computer Journal*, 39(4):255–274, September 1996.
- [23] A.H.M. ter Hofstede and Th.P. van der Weide. Expressiveness in conceptual data modelling. *Data & Knowledge Engineering*, 10(1):65–100, February 1993.
- [24] A.H.M. ter Hofstede and Th.P. van der Weide. Fact Orientation in Complex Object Role Modelling Techniques. In T.A. Halpin and R. Meersman, editors, *Proceedings of the First International Conference on Object-Role Modelling (ORM-1)*, pages 45–59, Townsville, Australia, July 1994.
- [25] A.H.M. ter Hofstede and Th.P. van der Weide. Deriving Identity from Extensionality. *International Journal of Software Engineering and Knowledge Engineering*, 8(2):189–221, June 1997.
- [26] S. Huffman and R.V. Zoeller. A Rule-Based System Tool for Automated ER Model Clustering. In F.H. Lochovsky, editor, *Proceedings of the Eight International Conference on Entity-Relationship Approach, Entity-Relationship Approach to Database Design and Querying*, pages 221–236, Toronto, Canada, 1990. Elsevier Science Publishers.
- [27] J. Martin. *Strategic Data Planning Methodologies*. Prentice-Hall, Englewood Cliffs, New Jersey, 1982.
- [28] D. Moody. A Practical Methodology for the Representation of Enterprise Data Models. In *Proceedings of the 2nd Annual conference on Information Systems and Database Special Interest Group*, Sydney, Australia, 1991.
- [29] A.H. Seltviet. An Abstraction-Based Approach to Large-Scale Information System Development. In C. Rolland, F. Bodart, and C. Cauvet, editors, *Proceedings of the Fifth International Conference CAiSE'93 on Advanced Information Systems Engineering*, volume 685 of *Lecture Notes in Computer Science*, Paris, France, 1993. Springer-Verlag.
- [30] G.C. Simson. A Structured Approach to Data Modelling. *The Australian Computer Journal*, 21(3), August 1989.
- [31] 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):975–987, August 1989.
- [32] D. Vermeir. Semantic Hierarchies and Abstractions in Conceptual Schemata. *Information Systems*, 8(2):117–124, 1983.

- [33] L.A. Walko. Caves: Visualization and Abstraction Mechanism for Object-Oriented Databases. In *Proceedings of the 3rd Australian Database Conference*, pages 10–35, Melbourne, Australia, 1992. World Scientific.
- [34] Th.P. van der Weide, A.H.M. ter Hofstede, and P. van Bommel. Uniquet: Determining the Semantics of Complex Uniqueness Constraints. *The Computer Journal*, 35(2):148–156, April 1992.