

# A Formal Framework for Reasoning on UML Class Diagrams

Andrea Calì, Diego Calvanese, Giuseppe De Giacomo, Maurizio Lenzerini

Dipartimento di Informatica e Sistemistica  
Università di Roma “La Sapienza”  
Via Salaria 113, I-00198 Roma, Italy  
*lastname@dis.uniroma1.it*

**Abstract.** In this paper we formalize UML class diagrams in terms of a logic belonging to Description Logics, which are subsets of First-Order Logic that have been thoroughly investigated in Knowledge Representation. The logic we have devised is specifically tailored towards the high expressiveness of UML information structuring mechanisms, and allows one to formally model important properties which typically can only be specified by means of qualifiers. The logic is equipped with decidable reasoning procedures which can be profitably exploited in reasoning on UML class diagrams. This makes it possible to provide computer aided support during the application design phase in order to automatically detect relevant properties, such as inconsistencies and redundancies.

## 1 Introduction

There is a vast consensus on the need for a precise semantics for UML [10, 13], in particular for UML class diagrams. Indeed, several types of formalization of UML class diagrams have been proposed in the literature [9–11, 7]. Many of them have been proved very useful with respect to the task of establishing a common understanding of the formal meaning of UML constructs. However, to the best of our knowledge, none of them has the explicit goal of building a solid basis for allowing automated reasoning techniques, based on algorithms that are sound and complete wrt the semantics, to be applicable to UML class diagrams.

In this paper, we propose a new formalization of UML class diagrams in terms of a particular formal logic of the family of Description Logics (DL). DLs<sup>1</sup> have been proposed as successors of semantic network systems like KL-ONE, with an explicit model-theoretic semantics. The research on these logics has resulted in a number of automated reasoning systems [14, 15, 12], which have been successfully tested in various application domains (see e.g., [18, 19, 17]). Our long term goal is to exploit the deductive capabilities of DL systems, and show that effective reasoning can be carried out on UML class diagrams, so as to provide support during the specification phase of software development.

In DLs, the domain of interest is modeled by means of *concepts* and *relationships*, which denote classes of objects and relations, respectively. Generally speaking, a DL is formed by three basic components:

---

<sup>1</sup> See <http://dl.kr.org> for the home page of Description Logics.

- A *description language*, which specifies how to construct complex concept and relationship expressions (also called simply concepts and relationships), by starting from a set of atomic symbols and by applying suitable constructors,
- a *knowledge specification mechanism*, which specifies how to construct a DL knowledge base, in which properties of concepts and relationships are asserted, and
- a set of *automatic reasoning procedures* provided by the DL.

The set of allowed constructors characterizes the expressive power of the description language. Various languages have been considered by the DL community, and numerous papers investigate the relationship between expressive power and computational complexity of reasoning (see [8] for a survey).

Several works point out that DLs can be profitably used to provide both formal semantics and reasoning support to formalisms in areas such as Natural Language, Configuration Management, Database Management, Software Engineering. For example, [5] illustrates the use of DLs for database modeling. However, to the best of our knowledge, DLs have not been applied to the Unified Modeling Language (UML) (with the exception of [3]). The goal of this work is to present a formalization of UML class diagrams in terms of DLs. In particular, we show how to map the constructs of a class diagram onto those of Description Logics. The mapping provides us with a rigorous logical framework for representing and automatically reasoning on UML class specifications. The logic we have devised is specifically tailored towards the high expressiveness of UML information structuring mechanisms, and allows one to formally model important properties which typically can only be specified by means of constraints. The logic is equipped with decidable reasoning procedures which can be profitably exploited in reasoning on UML class diagrams. This makes it possible to provide computer aided support during the application design phase, in order to automatically detect relevant properties, such as inconsistencies and redundancies.

The paper is organized as follows: in Section 2 we give an overview of the Description Logic we use, called  $\mathcal{DLR}$ . In Sections 3, 4, 5 and 6, we illustrate the formalization of UML class diagrams in terms of  $\mathcal{DLR}$ , focusing on classes, associations, generalization, and constraints, respectively. In Section 7 we discuss the use of the reasoning procedures associated to  $\mathcal{DLR}$  in order to support the specification of UML class diagrams. Section 8 concludes the paper.

## 2 The Description Logic $\mathcal{DLR}$

The goal of this section is to give an overview on the Description Logic  $\mathcal{DLR}$  introduced in [4], which is able to capture a great variety of data models with many forms of constraints [2, 6].

### 2.1 Syntax

The basic elements of  $\mathcal{DLR}$  are *concepts* (unary relations), and *n-ary relations*. We assume to deal with a finite set of atomic relations and atomic concepts, denoted by  $P$  and  $A$ , respectively. In this section, we use  $R$  to denote arbitrary relations (of given arity

between 2 and  $n_{max}$ ), and  $C$  to denote arbitrary concepts, respectively built according to the following syntax:

$$\begin{aligned} R &::= \top_n \mid P \mid (i/n : C) \mid \neg R \mid R_1 \sqcap R_2 \\ C &::= \top_1 \mid A \mid \neg C \mid C_1 \sqcap C_2 \mid (\leq k [i] R) \end{aligned}$$

where  $i$  denotes a component of a relation, i.e., an integer between 1 and  $n_{max}$ ,  $n$  denotes the *arity* of a relation, i.e., an integer between 2 and  $n_{max}$ , and  $k$  denotes a non-negative integer. We consider only concepts and relations that are *well-typed*, which means that (i) only relations of the same arity  $n$  are combined to form expressions of type  $R_1 \sqcap R_2$  (which inherit the arity  $n$ ), and (ii)  $i \leq n$  whenever  $i$  denotes a component of a relation of arity  $n$ .

We also make use of the following abbreviations:  $C_1 \sqcup C_2$  for  $\neg(\neg C_1 \sqcap \neg C_2)$ ;  $C_1 \Rightarrow C_2$  for  $\neg C_1 \sqcup C_2$ ;  $(\geq k [i] R)$  for  $\neg(\leq k-1 [i] R)$ ;  $\exists [i] R$  for  $(\geq 1 [i] R)$ ;  $\forall [i] R$  for  $\neg \exists [i] \neg R$ . Moreover, we abbreviate  $(i/n : C)$  with  $(i : C)$  when  $n$  is clear from the context.

A  $\mathcal{DLR}$  Knowledge Base (KB) is constituted by a finite set of *inclusion assertions*, where each assertion has one of the forms:

$$R_1 \sqsubseteq R_2 \qquad C_1 \sqsubseteq C_2$$

with  $R_1$  and  $R_2$  of the same arity.

Besides inclusion assertions,  $\mathcal{DLR}$  KBs allow for assertions expressing identification constraints and functional dependencies.

An *identification assertion* on a concept has the form:

$$(\text{id } C [i_1] R_1, \dots, [i_h] R_h)$$

where  $C$  is a concept, each  $R_j$  is a relation, and each  $i_j$  denotes one component of  $R_j$ . Intuitively, such an assertion states that no two different instances of  $C$  agree on the participation to  $R_1, \dots, R_h$ . In other words, if  $a$  is an instance of  $C$  that is the  $i_j$ -th component of a tuple  $t_j$  of  $R_j$ , for  $j \in \{1, \dots, h\}$ , and  $b$  is an instance of  $C$  that is the  $i_j$ -th component of a tuple  $s_j$  of  $R_j$ , for  $j \in \{1, \dots, h\}$ , and for each  $j$ ,  $t_j$  agrees with  $s_j$  in all components different from  $i_j$ , then  $a$  and  $b$  coincide.

A *functional dependency assertion* on a relation has the form:

$$(\text{fd } R i_1, \dots, i_h \rightarrow j)$$

where  $R$  is a relation,  $h \geq 2$ , and  $i_1, \dots, i_h, j$  denote components of  $R$ . The assertion imposes that two tuples of  $R$  that agree on the components  $i_1, \dots, i_h$ , agree also on the component  $j$ .

Note that unary functional dependencies (i.e., functional dependencies with  $h = 1$ ) are ruled out in  $\mathcal{DLR}$ , since these lead to undecidability of reasoning [4]. Note also that the right hand side of a functional dependency contains a single element. However, this is not a limitation, because any functional dependency with more than one element in the right hand side can always be split into several dependencies of the above form.

$\top_n^{\mathcal{I}} \subseteq (\Delta^{\mathcal{I}})^n$	$\top_1^{\mathcal{I}} = \Delta^{\mathcal{I}}$
$P^{\mathcal{I}} \subseteq \top_n^{\mathcal{I}}$	$A^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}}$
$(i/n : C)^{\mathcal{I}} = \{t \in \top_n^{\mathcal{I}} \mid t[i] \in C^{\mathcal{I}}\}$	$(\neg C)^{\mathcal{I}} = \Delta^{\mathcal{I}} \setminus C^{\mathcal{I}}$
$(\neg R)^{\mathcal{I}} = \top_n^{\mathcal{I}} \setminus R^{\mathcal{I}}$	$(C_1 \sqcap C_2)^{\mathcal{I}} = C_1^{\mathcal{I}} \cap C_2^{\mathcal{I}}$
$(R_1 \sqcap R_2)^{\mathcal{I}} = R_1^{\mathcal{I}} \cap R_2^{\mathcal{I}}$	$(\leq k [i] R)^{\mathcal{I}} = \{a \in \Delta^{\mathcal{I}} \mid \#\{t \in R_1^{\mathcal{I}} \mid t[i] = a\} \leq k\}$

**Fig. 1.** Semantic rules for  $\mathcal{DLR}$  ( $P$ ,  $R$ ,  $R_1$ , and  $R_2$  have arity  $n$ )

## 2.2 Semantics

The semantics of  $\mathcal{DLR}$  is specified through the notion of interpretation. An *interpretation*  $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$  of a  $\mathcal{DLR}$  KB  $\mathcal{K}$  is constituted by an *interpretation domain*  $\Delta^{\mathcal{I}}$  and an *interpretation function*  $\cdot^{\mathcal{I}}$  that assigns to each concept  $C$  a subset  $C^{\mathcal{I}}$  of  $\Delta^{\mathcal{I}}$  and to each relation  $R$  of arity  $n$  a subset  $R^{\mathcal{I}}$  of  $(\Delta^{\mathcal{I}})^n$ , such that the conditions in Figure 1 are satisfied. (In the figure,  $t[i]$  denotes the  $i$ -th component of tuple  $t$ .) We observe that  $\top_1$  denotes the interpretation domain, while  $\top_n$ , for  $n > 1$ , does *not* denote the  $n$ -Cartesian product of the domain, but only a subset of it, that covers all relations of arity  $n$ . It follows, from this property, that the “ $\neg$ ” constructor on relations is used to express difference of relations, rather than complement.

To specify the semantics of a KB we first define when an interpretation satisfies an assertion as follows:

- An interpretation  $\mathcal{I}$  *satisfies* an inclusion assertion  $R_1 \sqsubseteq R_2$  (resp.  $C_1 \sqsubseteq C_2$ ) if  $R_1^{\mathcal{I}} \subseteq R_2^{\mathcal{I}}$  (resp.  $C_1^{\mathcal{I}} \subseteq C_2^{\mathcal{I}}$ ).
- An interpretation  $\mathcal{I}$  *satisfies* the assertion  $(\mathbf{id} C [i_1]R_1, \dots, [i_h]R_h)$  if for all  $a, b \in C^{\mathcal{I}}$  and for all  $t_1, s_1 \in R_1^{\mathcal{I}}, \dots, t_h, s_h \in R_h^{\mathcal{I}}$  we have that:

$$\left. \begin{array}{l} a = t_1[i_1] = \dots = t_h[i_h], \\ b = s_1[i_1] = \dots = s_h[i_h], \\ t_j[i] = s_j[i], \text{ for } j \in \{1, \dots, h\}, \text{ and for } i \neq i_j \end{array} \right\} \text{ implies } a = b$$

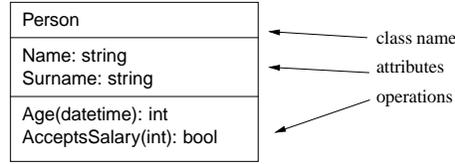
- An interpretation  $\mathcal{I}$  *satisfies* the assertion  $(\mathbf{fd} R i_1, \dots, i_h \rightarrow j)$  if for all  $t, s \in R^{\mathcal{I}}$ , we have that:

$$t[i_1] = s[i_1], \dots, t[i_h] = s[i_h] \quad \text{implies} \quad t[j] = s[j]$$

An interpretation that satisfies all assertions in a KB  $\mathcal{K}$  is called a *model* of  $\mathcal{K}$ .

## 2.3 Reasoning tasks

Several reasoning services are applicable to  $\mathcal{DLR}$  KBs. The most important ones are KB satisfiability and logical implication. A KB  $\mathcal{K}$  is *satisfiable* if there exists a model of  $\mathcal{K}$ . A concept  $C$  is *satisfiable* in a KB  $\mathcal{K}$  if there is a model  $\mathcal{I}$  of  $\mathcal{K}$  such that  $C^{\mathcal{I}}$  is nonempty. A concept  $C_1$  is *subsumed* by a concept  $C_2$  in a KB  $\mathcal{K}$  if  $C_1^{\mathcal{I}} \subseteq C_2^{\mathcal{I}}$  for every model  $\mathcal{I}$  of  $\mathcal{K}$ . An assertion  $\alpha$  is *logically implied* by  $\mathcal{K}$  if all models of  $\mathcal{K}$  satisfy  $\alpha$ . One can easily verify that logical implication and KB unsatisfiability are mutually reducible.



**Fig. 2.** Representation of a class in UML

One of the distinguishing features of  $\mathcal{DLR}$  is that it is equipped with reasoning algorithms that are sound and complete wrt to the semantics. Such algorithms allow one to decide all the above reasoning tasks in deterministic exponential time [4]. Indeed, the proposed algorithms are computationally optimal, since reasoning in  $\mathcal{DLR}$  is EXPTIME-complete [2].

### 3 Classes

In this paper we concentrate on class diagrams for the conceptual perspective. Hence, we do not deal with those features that are relevant for the implementation perspective, such as public, protected, and private qualifiers for methods and attributes.

A *class* in an UML class diagram denotes a *sets of objects* with common features. A class is graphically rendered as a rectangle divided into parts, as shown in Figure 2.

The first part contains the *name* of the class, which has to be unique in the whole diagram. The second part contains the *attributes* of the class, each denoted by a name (possibly followed by the *multiplicity*, between square brackets) and with an associated *class*, which indicates the domain of the attribute values. For example, an attribute `phoneNumber[1..*]: String` means that each instance of the class has at least one phone number, and possibly more, and that each phone numbers is an instance of `String`. If not otherwise specified, attributes are *single-valued*. The third part contains the *operations* of the class, i.e., the operations associated to the objects of the class. An operation definition has the form:

$$\textit{operation-name}(\textit{parameter-list}): (\textit{return-list})$$

Observe that an operation may return a *tuple* of objects as result.

An UML class is represented by a  $\mathcal{DLR}$  concept. This follows naturally from the fact that both UML classes and  $\mathcal{DLR}$  concepts denote *sets of objects*.

An UML *attribute*  $a$  of type  $C'$  for a class  $C$  associates to each instance of  $C$ , zero, one, or more instances of a class  $C'$ . An optional *multiplicity*  $[i..j]$  for  $a$  specifies that  $a$  associates to each instance of  $C$ , at least  $i$  and most  $j$  instances of  $C'$ . When the multiplicity is missing,  $[1..1]$  is assumed, i.e., the attribute is *mandatory* and *single-valued*.

To formalize attributes we have to think of an attribute  $a$  of type  $C'$  for a class  $C$  as a binary relation between instances of  $C$  and instances of  $C'$ . We capture such a binary relation by means of a binary relation  $a$  of  $\mathcal{DLR}$ . To specify the type of the attribute we use the assertion:

$$C \sqsubseteq \forall[1](a \Rightarrow (2 : C'))$$

Such an assertion specifies precisely that, for each instance  $c$  of the concept  $C$ , all objects related to  $c$  by  $a$ , are instances of  $C'$ . Note that an attribute name is not necessarily unique in the whole diagram, and hence two different classes could have the same attribute, possibly of different types. This situation is correctly captured by the formalization in  $\mathcal{DLR}$ .

To specify the multiplicity  $[i..j]$  associated to the attribute we add the assertion:

$$C \sqsubseteq (\geq i [1]a) \sqcap (\leq j [1]a)$$

Such an assertion specifies that each instance of  $C$  participates at least  $i$  times and at most  $j$  times to relation  $a$  via component 1. If  $i = 0$ , i.e., the attribute is *optional*, we omit the first conjunct, and if  $j = *$  we omit the second one. Observe that for attributes with multiplicity  $[0..*]$  we omit the whole assertion, and that, when the multiplicity is missing the above assertion becomes:

$$C \sqsubseteq \exists[1]a \sqcap (\leq 1 [1]a)$$

An operation of a class is a function from the objects of the class to which the operation is associated, and possibly additional parameters, to tuples of objects. In class diagrams, the code associated to the operation is not considered and typically, what is represented is only the signature of the operation.

In  $\mathcal{DLR}$ , we model operations by means of  $\mathcal{DLR}$  relations. Let

$$f(P_1, \dots, P_m) : (R_1, \dots, R_n)$$

be an operation of a class  $C$  that has  $m$  parameters belonging to the classes  $P_1, \dots, P_m$  respectively and  $n$  return values belonging to  $R_1, \dots, R_n$  respectively. We formalize such an operation as a  $\mathcal{DLR}$  relation, named  $\text{op}_{f(P_1, \dots, P_m) : (R_1, \dots, R_n)}$ , of arity  $m+n+1$  among instances of the  $\mathcal{DLR}$  concepts  $C, P_1, \dots, P_m, R_1, \dots, R_n$ . On such a relation we enforce the following assertions:

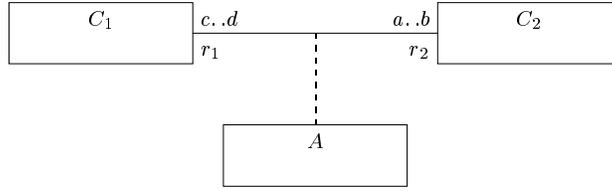
- An assertion imposing the correct types to parameters and return values:

$$C \sqsubseteq \forall[1](\text{op}_{f(P_1, \dots, P_m) : (R_1, \dots, R_n)} \Rightarrow ((2 : P_1) \sqcap \dots \sqcap (m+1 : P_m) \sqcap (m+2 : R_1) \sqcap \dots \sqcap (m+n+1 : R_n)))$$

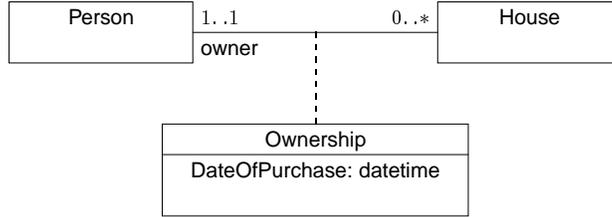
- Assertions imposing that invoking the operation on a given object with given parameters determines in a unique way each return value (i.e., the relation corresponding to the operation is in fact a function from the invocation object and the parameters to the returned values):

$$\begin{aligned} & (\text{fd}_{\text{op}_{f(P_1, \dots, P_m) : (R_1, \dots, R_n)}} 1, \dots, m+1 \rightarrow m+2) \\ & \dots \\ & (\text{fd}_{\text{op}_{f(P_1, \dots, P_m) : (R_1, \dots, R_n)}} 1, \dots, m+1 \rightarrow m+n+1) \end{aligned}$$

These functional dependencies are determined only by the number of parameters and the number of result values, and not by the specific class for which the operation is defined, nor by the types of parameters and result values.



**Fig. 3.** Association in UML



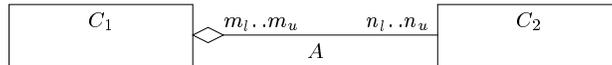
**Fig. 4.** Example of association in UML

The *overloading* of operations does not pose any difficulty in the formalization since an operation is represented in  $\mathcal{DLR}$  by a relation having as name the whole signature of the operation, which consists not only the name of the operation but also the parameter and return value types. Observe that the formalization of operations in  $\mathcal{DLR}$  allows one to have operations with the same name or even with the same signature in two different classes.

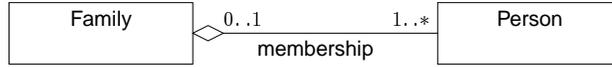
#### 4 Associations and Aggregations

An *association* in UML is a relation between the instances of two or more classes. An association often has a related *association class* that describes properties of the association such as attributes, operations, etc. An association  $A$  between the instances of two classes  $C_1$  and  $C_2$  is graphically rendered as in Figure 3, where the class  $A$  is the association class related to the association,  $r_1$  and  $r_2$  are the *role names* of  $C_1$  and  $C_2$  respectively, i.e., they specify the role that each class plays within the relation  $R$ , and where the *multiplicity*  $a..b$  specifies that each instance of class  $C_1$  can participate at least  $a$  times and at most  $b$  times to relation  $A$ ;  $c..d$  has an analogous meaning for class  $C_2$ .

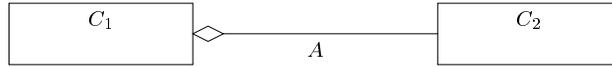
*Example 1.* The association in Figure 4 models ownership of houses; a person can own 0 or more houses (the symbol  $*$  means  $\infty$ ), while a house has to be owned by exactly one person, which is denoted by the role *owner*.



**Fig. 5.** Aggregation in UML



**Fig. 6.** Example of aggregation in UML



**Fig. 7.** Aggregation in UML

An *aggregation* in UML is a binary association between the instances of two classes, denoting a part-whole relationship, i.e., a relationship that specifies that each instance of a class is made up of a set of instances of another class. An aggregation is graphically rendered as shown in Figure 5, where the diamond indicates the *containing class*, opposed to the *contained class*. The multiplicity has the same meaning as in associations. As for associations, also for aggregation it is possible to define role names which denote the role each class plays in the aggregation.

*Example 2.* In figure 6 we have persons belonging to families; a family can have one or more persons, while a person belongs to at most one family.

In UML class diagrams, both associations and aggregations denote relationships between classes. However aggregations are simpler than general associations, since they are necessarily binary, while association can be of any arity. Moreover aggregations do not have an associated class, while associations typically have a related association class. Observe that names of associations and names of aggregations (as names of classes) are *unique*. In other words there cannot be two associations/aggregations with the same name.

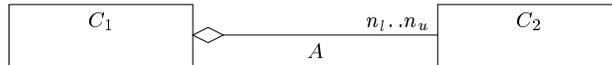
Next we turn to the formalization in  $\mathcal{DLR}$ . We first concentrate on aggregations, which are simpler to model than associations. An aggregation  $A$ , saying that instances of the class  $C_1$  have 0 or more components that are instances of the class  $C_2$ , depicted in Figure 7, is formalized in  $\mathcal{DLR}$  by means of a binary relation  $A$  on which the following assertion is enforced:

$$A \sqsubseteq (1 : C_1) \sqcap (2 : C_2).$$

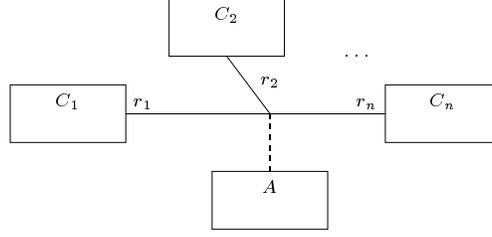
Note that the distinction between the contained class and the containing class is not lost. Indeed, we simply use the following convention: *the first argument of the relation is the containing class*.

As we have seen for class attributes, the multiplicity of an aggregation can be easily expressed in  $\mathcal{DLR}$ . For example, if we have a multiplicity on the participation of instances of  $C_2$  for each given instance of  $C_1$ , as in Figure 8, we simply state the assertion

$$C_1 \sqsubseteq (\geq n_l [1]A) \sqcap (\leq n_u [1]A)$$



**Fig. 8.** Multiplicity in aggregation



**Fig. 9.** Association in UML

We can use a similar assertion for a multiplicity on the participation of instances of  $C_1$  for each given instance of  $C_2$ .

Observe that, in the formalization in  $\mathcal{DLR}$  of aggregation, role names do not play any role. If we want to keep track of them in the formalization, it suffices to consider them as convenient abbreviations for the components of the  $\mathcal{DLR}$  relation modeling the aggregation.

Next we focus on *associations*. Since associations have often a related association class, we formalize associations in  $\mathcal{DLR}$  by reifying each association  $A$  into a  $\mathcal{DLR}$  concept  $A$  with suitable properties. Let us consider the association shown in Figure 9. We represent it in  $\mathcal{DLR}$  by introducing a concept  $A$  and  $n$  binary relations  $r_1, \dots, r_n$ , one for each component of the association  $A$ <sup>2</sup>. Then we enforce the following assertion:

$$\begin{aligned}
C \sqsubseteq & \exists[1]r_1 \sqcap (\leq 1 [1]r_1) \sqcap \forall[1](r_1 \Rightarrow (2 : C_1)) \sqcap \\
& \exists[1]r_2 \sqcap (\leq 1 [1]r_2) \sqcap \forall[1](r_2 \Rightarrow (2 : C_2)) \sqcap \\
& \vdots \\
& \exists[1]r_n \sqcap (\leq 1 [1]r_n) \sqcap \forall[1](r_n \Rightarrow (2 : C_n))
\end{aligned}$$

where  $\exists[1]r_i$  (with  $i \in \{1, \dots, n\}$ ) specifies that the concept  $A$  must have all components  $r_1, \dots, r_n$  of the association  $A$ ,  $(\leq 1 [1]r_i)$  (with  $i \in \{1, \dots, n\}$ ) specifies that each such component is single-valued, and  $\forall[1](r_i \Rightarrow (2 : C_i))$  (with  $i \in \{1, \dots, n\}$ ) specifies the class each component has to belong to. Finally, we use the assertion

$$(\mathbf{id} A [1]r_1, \dots, [1]r_n)$$

to specify that each instance of  $A$  represents a *distinct* tuple in  $C_1 \times \dots \times C_n$ .

We can easily represent in  $\mathcal{DLR}$  a multiplicity on a binary association, by imposing a number restriction on the relations modeling the components of the association. Differently from aggregation, however, the names of such relations (which correspond to roles) are unique wrt to the association only, not the entire diagram. Hence we have to state such constraints in  $\mathcal{DLR}$  in a slightly more involved way.

Suppose we have a situation like that in Figure 10. Consider the association  $A_1$  and the constraint saying that for each instance of  $C$  there can be at least  $n_{l_1}$  and at most  $n_{u_1}$  instances of  $D_1$  related by  $A_1$  to it. We capture this constraint as follows:

<sup>2</sup> These relations may have the name of the roles of the association if available in the UML diagram, or an arbitrary name if role names are not available. In any case, we preserve the possibility of using the same role name in different associations.

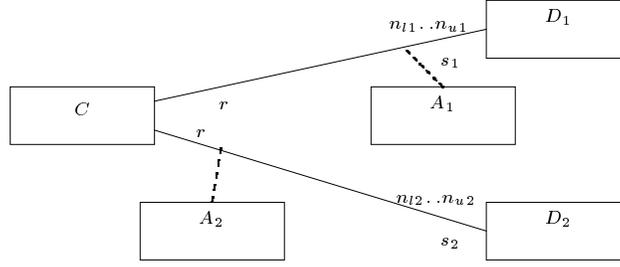


Fig. 10. Multiplicity in aggregation

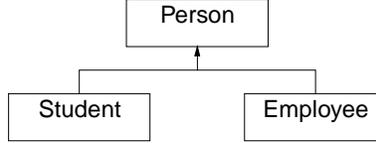


Fig. 11. A class hierarchy in UML

$$C \sqsubseteq (\geq n_{l1} [2](r \sqcap (1 : A_1))) \sqcap (\leq n_{u1} [2](r \sqcap (1 : A_1)))$$

Observe that nothing prevents  $C$  to participate to a different association  $A_2$  with the same role  $r$  but with different multiplicity  $n_{l2}..n_{u2}$ . Observe that this is modeled by the totally unrelated assertion:

$$C \sqsubseteq (\geq n_{l2} [2](r \sqcap (1 : A_2))) \sqcap (\leq n_{u2} [2](r \sqcap (1 : A_2)))$$

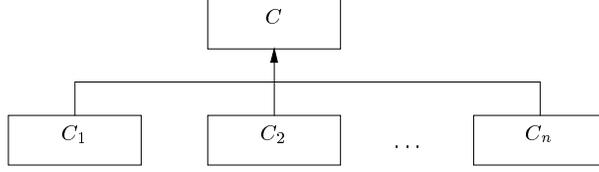
## 5 Generalization and Inheritance

In UML one can use *generalization* between a parent class and a child class to specify that each instance of the child class is also an instance of the parent class. Hence, the instances of the child class inherit the properties of the parent class, but typically they satisfy additional properties that in general do not hold for the parent class. For example, in Figure 11, an employee may have a salary and be associated with an employer, while this does not hold for a person in general.

Generalization is naturally supported in  $\mathcal{DLR}$ . If an UML class  $C_2$  generalizes a class  $C_1$ , we can express this by the  $\mathcal{DLR}$  assertion:

$$C_1 \sqsubseteq C_2$$

Inheritance between  $\mathcal{DLR}$  concepts works exactly as inheritance between UML classes. This is an obvious consequence of the semantics of  $\sqsubseteq$  which is based on subsetting. Indeed, given an assertion  $C_1 \sqsubseteq C_2$ , every tuple in a  $\mathcal{DLR}$  relation having  $C_2$  as  $i$ -th argument type may have as  $i$ -th component an instance of  $C_1$ , which is in fact also an instance of  $C_2$ . As a consequence, in the formalization, each attribute or operation of



**Fig. 12.** A class hierarchy in UML

$C_2$ , and each aggregation and association involving  $C_2$  is correctly inherited by  $C_1$ . Observe that the formalization in  $\mathcal{DLR}$  also captures directly inheritance among association classes, which are treated exactly as all other classes, and multiple inheritance between classes (including association classes).

Moreover in UML, one can group several generalizations into a class hierarchy, as shown in Figure 12. Such a hierarchy is captured in  $\mathcal{DLR}$  by a set of inclusion assertions, one between each child class and the parent class:

$$C_i \sqsubseteq C \quad \text{for each } i \in \{1, \dots, n\}$$

We discuss in Section 6 how to formalize in  $\mathcal{DLR}$  additional properties of a class hierarchy, such as mutual disjointness between the child classes, or covering of the parent class.

In UML it is possible to override attributes or operations of a superclass. That is, it is possible to specialize an attribute or an operation for the subclass. From the conceptual point of view such a specialization needs to remain compatible with the original definition of the attribute/operation, i.e., the attribute/operation of the subclass can only be a *restriction* of the corresponding attribute/operation belonging to the superclass. For attributes, this means that one can restrict the type of the attribute to be a subclass of the original type, or restrict the multiplicity wrt to the one specified for the superclass. For operations, while keeping the same signature, one may restrict (by means of constraints) the return types and possibly also the argument types to be subclasses of the original ones<sup>3</sup>.

We illustrate by means of an example how one can correctly model such forms of overriding in  $\mathcal{DLR}$ .

*Example 3.* If  $C'$  is a subclass of  $C$ , and  $C$  has an operation  $f(C_1, C_2) : C_3$ , then  $C'$  inherits the operation. In  $\mathcal{DLR}$  we have that  $C' \sqsubseteq C$  and as a consequence  $C'$  inherits the properties that hold for  $C$  wrt the participation in relation  $\text{op}_{f(C_1, C_2):C_3}$ , which models the operation. In  $C'$  we can further restrict such properties by requiring for example that the result value belongs to a class  $C'_3$  that is a specialization of  $C_3$ . We can do so by adding the assertion

$$C' \sqsubseteq \forall[1](\text{op}_{f(C_1, C_2):C_3} \Rightarrow (4 : C'_3))$$

where we assume that we have also the assertion  $C'_3 \sqsubseteq C_3$ .

<sup>3</sup> Observe that restricting the argument types corresponds, in the implementation of the operation, to restrict the preconditions for the applicability of the operation.

## 6 Constraints

In UML it is possible to add information to a class diagram by using *constraints*. In general, constraints are used to express in an informal way information which cannot be expressed by other constructs of UML class diagrams. We discuss here common types of constraints that occur in UML class diagrams and how they can be taken into account when formalizing class diagrams in  $\mathcal{DLR}$ .

Often, when defining generalizations between classes, we need to add additional constraints among the involved classes. For example, for the class hierarchy in Figure 12, a constraint may express that  $C_1, \dots, C_n$  are *mutually disjoint*. In  $\mathcal{DLR}$ , such a relationship can be expressed by the following assertions:

$$C_i \sqsubseteq \neg C_j \quad \text{for each } i, j \in \{1, \dots, n\} \text{ with } i \neq j$$

In general, in UML, if not otherwise specified by a constraint, two classes may have common instances, i.e., they are *not disjoint*. If a constraint imposes the disjointness of two classes, say  $C$  and  $C'$ , this can be formalized in  $\mathcal{DLR}$  by means of the assertion:

$$C \sqsubseteq \neg C'$$

Disjointness of classes is just one example of *negative information*. Again, by exploiting the expressive power of  $\mathcal{DLR}$ , we can express additional forms of negative information, usually not considered in UML, by introducing suitable assertions. For example, we can enforce that no instance of a class  $C$  has an attribute  $a$  by means of the assertion:

$$C \sqsubseteq \neg \exists[1]a$$

Analogously, one can assert that no instance of a class is involved in a given association or aggregation.

Turning again the attention to generalization hierarchies, by default, in UML a generalization hierarchy is open, in the sense that there may be instances of the superclass that are not instances of any of the subclasses. This allows for extending the diagram more easily, in the sense that the introduction of a new subclass does not change the semantics of the superclass. However, in specific situations, it may happen that in a generalization hierarchy, the superclass  $C$  is a covering of the subclasses  $C_1, \dots, C_n$ . We can represent such a situation in  $\mathcal{DLR}$  by simply including the additional assertion:

$$C \sqsubseteq C_1 \sqcup \dots \sqcup C_n$$

The above assertion models a form of *disjunctive information*: each instance of  $C$  is either an instance of  $C_1$ , or an instance of  $C_2$ , ... or an instance of  $C_n$ .

Other forms of disjunctive information can be modeled by exploiting the expressive power of  $\mathcal{DLR}$ . For example, that an attribute  $a$  is present only for a specified set  $C_1, \dots, C_n$  of classes can be modeled by suitably using union of classes as follows:

$$\exists[1]a \sqsubseteq C_1 \sqcup \dots \sqcup C_n$$

*Keys* are a modeling notion that is very common in databases, and they are used to express that certain attributes uniquely identify the instances of a class. We can exploit

the expressive power of  $\mathcal{DLR}$  in order to associate keys to classes. If an attribute  $a$  is a key for a class  $C$  this means that there is no pair of instances of  $C$  that have the same value for  $a$ . We can capture this in  $\mathcal{DLR}$  by means of the following assertion:

$$(\text{id } C [1]a)$$

More generally, we are able to specify that a *set* of attributes  $\{a_1, \dots, a_n\}$  is a key for  $C$ ; in this case we use the assertion:

$$(\text{id } C [1]a_1, \dots, [1]a_n)$$

As already discussed in Section 5, constraints that correspond to the specialization of the type of an attribute or its multiplicity can be represented in  $\mathcal{DLR}$ . Similarly, consider the case of a class  $C$  participating in an aggregation  $A$  with a class  $D$ , and where  $C$  and  $D$  have subclasses  $C'$  and  $D'$  respectively, related via an aggregation  $A'$ . A *subset constraint* from  $A'$  to  $A$  can be modeled correctly in  $\mathcal{DLR}$  by means of the assertion

$$A \sqsubseteq A'$$

involving the two binary relations  $A$  and  $A'$  that represent the aggregations.

More generally, one can exploit the expressive power of  $\mathcal{DLR}$  to formalize several types of constraints that allow one to better represent the application semantics and that are typically not dealt with in a formal way. Observe that this allows one to take such constraints fully into account when reasoning on the class diagram.

## 7 Reasoning on Class Diagrams

Traditional CASE tools support the designer with a user friendly graphical environment and provide powerful means to access different kinds of repositories that store information associated to the elements of the developed project. However, no support for higher level activities related to managing the complexity of the design is provided. In particular, the burden of checking relevant properties of class diagrams, such as consistency or redundancy (see below), is left to the responsibility of the designer.

Thus, the formalization in  $\mathcal{DLR}$  of UML class diagrams, and the fact that properties of inheritance and relevant types of constraints are perfectly captured by the formalization in  $\mathcal{DLR}$  and the associated reasoning tasks, provides the ability to reason on class diagrams. This represents a significant improvement and is a first step towards the development of modeling tools that offer an automated reasoning support to the designer in his modeling activity.

We briefly discuss the tasks that can be performed by exploiting the reasoning capabilities of a  $\mathcal{DLR}$  reasoner [15, 16], and that allow a modeling tool to take over tasks traditionally left to the responsibility of the designer. Such a tool may construct from a class diagram a  $\mathcal{DLR}$  knowledge base, and manage it in a way completely transparent to the designer. By exploiting the  $\mathcal{DLR}$  reasoning services various kinds of checks can be performed on the class diagram.<sup>4</sup>

<sup>4</sup> A prototype design tool with such a kind of automated reasoning support is available at <http://www.cs.man.ac.uk/~franconi/icom/>.

*Consistency of the class diagram* A class diagram is *consistent*, if its classes can be populated without violating any of the constraints in the diagram. Observe that the interaction of various types of constraints may make it very difficult to detect inconsistencies. By exploiting the formalization in  $\mathcal{DLR}$ , the consistency of a class diagram can be checked by checking the satisfiability of the corresponding  $\mathcal{DLR}$  knowledge base.

*Class Consistency* A class is *consistent*, if it can be populated without violating any of the constraints in the class diagram. The inconsistency of a class may be due to a design error or due to over-constraining. In any case, the designer can be forced to remove the inconsistency, either by correcting the error, or by relaxing some constraints, or by deleting the class, thus removing redundancy from the diagram. Exploiting the formalization in  $\mathcal{DLR}$ , class consistency can be checked by checking satisfiability of the corresponding concept in the  $\mathcal{DLR}$  knowledge base representing the class diagram.

*Class Equivalence* Two classes are *equivalent* if they denote the same set of instances whenever the constraints imposed by the class diagram are satisfied. Determining equivalence of two classes allows for their merging, thus reducing the complexity of the diagram. Again, checking class equivalence amounts to check the equivalence in  $\mathcal{DLR}$  of the corresponding concepts.

*Class Subsumption* A class  $C_1$  is *subsumed* by a class  $C_2$  if, whenever the constraints imposed by the class diagram are satisfied, the extension of  $C_1$  is a subset of the extension of  $C_2$ . Such a subsumption allows one to deduce that properties for  $C_1$  hold also for  $C_2$ . It is also the basis for a *classification* of all the classes in a diagram. Such a classification, as in any object-oriented approach, can be exploited in several ways within the modeling process [1]. Subsumption, and hence classification, can be checked by verifying subsumption in  $\mathcal{DLR}$ .

*Logical Consequence* A property is a *logical consequence* of a class diagram if it holds whenever all constraints specified in the diagram are satisfied. As an example, consider the generalization hierarchy depicted in Figure 12 and assume that a constraint specifies that it is complete. If an attribute  $a$  is defined as mandatory for all classes  $C_1, \dots, C_n$  then it follows logically that the same attribute is mandatory also for class  $C$ , even if not explicitly present in the diagram. Determining logical consequence is useful on the one hand to reduce the complexity of the diagram by removing those constraints that logically follow from other ones, and on the other hand it can be used to explicit properties that are implicit in the diagram, thus enhancing its readability.

Logical consequence can be captured by logical implication in  $\mathcal{DLR}$ , and determining logical implication is at the basis of all types of reasoning that a  $\mathcal{DLR}$  reasoning system can provide. In particular, observe that all reasoning tasks we have considered above can be rephrased in terms of logical consequence.

## 8 Conclusions

We have proposed a new formalization of UML class diagrams in terms of a particular formal logic of the family of Description Logics. Our long term goal is to exploit the

deductive capabilities of DL systems, thus showing that effective reasoning can be carried out on UML class diagrams, so as to provide support during the specification phase of software development. As a first step, we have shown in this paper how to map the constructs of a class diagram onto those of Description Logics. The mapping provides us with a rigorous logical framework for representing and automatically reasoning on UML class specifications.

We have already started experimenting our approach. In particular, we have used FACT for representing and reasoning on class diagrams. Although FACT does not yet incorporate all features required by our formalization (e.g., keys), the first results are encouraging. In particular, we have been able to draw interesting, non-trivial inferences on class diagrams containing about 50 classes. More experiments are under way, and we plan to report on them in the near future.

In the future, we aim at extending our formalization in order to capture further aspects of the UML. Our first step in this direction will be to add to our formal framework the possibility of modeling and reasoning on objects and links (i.e., instances of classes and associations).

## References

1. S. Bergamaschi and B. Nebel. Acquisition and validation of complex object database schemata supporting multiple inheritance. *Applied Intelligence*, 4(2):185–203, 1994.
2. D. Calvanese, G. De Giacomo, and M. Lenzerini. On the decidability of query containment under constraints. In *Proc. of PODS'98*, pages 149–158, 1998.
3. D. Calvanese, G. De Giacomo, and M. Lenzerini. Reasoning in expressive description logics with fixpoints based on automata on infinite trees. In *Proc. of IJCAI'99*, pages 84–89, 1999.
4. D. Calvanese, G. De Giacomo, and M. Lenzerini. Identification constraints and functional dependencies in description logics. In *Proc. of IJCAI 2001*, pages 155–160, 2001.
5. D. Calvanese, M. Lenzerini, and D. Nardi. Description logics for conceptual data modeling. In J. Chomicki and G. Saake, editors, *Logics for Databases and Information Systems*, pages 229–264. Kluwer Academic Publisher, 1998.
6. D. Calvanese, M. Lenzerini, and D. Nardi. Unifying class-based representation formalisms. *J. of Artificial Intelligence Research*, 11:199–240, 1999.
7. T. Clark and A. S. Evans. Foundations of the Unified Modeling Language. In D. Duke and A. Evans, editors, *Proc. of the 2nd Northern Formal Methods Workshop*. Springer-Verlag, 1997.
8. F. M. Donini, M. Lenzerini, D. Nardi, and A. Schaerf. Reasoning in description logics. In G. Brewka, editor, *Principles of Knowledge Representation*, Studies in Logic, Language and Information, pages 193–238. CSLI Publications, 1996.
9. A. Evans, R. France, K. Lano, and B. Rumpe. The UML as a formal modeling notation. In H. Kilov, B. Rumpe, and I. Simmonds, editors, *Proc. of the OOPSLA'97 Workshop on Object-oriented Behavioral Semantics*, pages 75–81. Technische Universität München, TUM-I9737, 1997.
10. A. Evans, R. France, K. Lano, and B. Rumpe. Meta-modelling semantics of UML. In H. Kilov, editor, *Behavioural Specifications for Businesses and Systems*, chapter 2. Kluwer Academic Publisher, 1999.
11. A. S. Evans. Reasoning with UML class diagrams. In *Second IEEE Workshop on Industrial Strength Formal Specification Techniques (WIFT'98)*. IEEE Computer Society Press, 1998.
12. V. Haarslev and R. Möller. Expressive ABox reasoning with number restrictions, role hierarchies, and transitively closed roles. In *Proc. of KR 2000*, pages 273–284, 2000.

13. D. Harel and B. Rumpe. Modeling languages: Syntax, semantics and all that stuff. Technical Report MCS00-16, The Weizmann Institute of Science, Rehovot, Israel, 2000.
14. I. Horrocks. Using an expressive description logic: FaCT or fiction? In *Proc. of KR'98*, pages 636–647, 1998.
15. I. Horrocks and P. F. Patel-Schneider. Optimizing description logic subsumption. *J. of Log. and Comp.*, 9(3):267–293, 1999.
16. I. Horrocks, U. Sattler, and S. Tobies. Practical reasoning for expressive description logics. In H. Ganzinger, D. McAllester, and A. Voronkov, editors, *Proc. of LPAR'99*, number 1705 in LNAI, pages 161–180. Springer-Verlag, 1999.
17. T. Kirk, A. Y. Levy, Y. Sagiv, and D. Srivastava. The Information Manifold. In *Proceedings of the AAAI 1995 Spring Symp. on Information Gathering from Heterogeneous, Distributed Enviroments*, pages 85–91, 1995.
18. D. L. McGuinness and J. R. Wright. An industrial strength description logic-based configuration platform. *IEEE Intelligent Systems*, pages 69–77, 1998.
19. U. Sattler. *Terminological Knowledge Representation Systems in a Process Engineering Application*. PhD thesis, LuFG Theoretical Computer Science, RWTH-Aachen, Germany, 1998.