

Inferring the Principal Type and the Schema Requirements of an OQL Query

A. Trigoni and G.M. Bierman

University of Cambridge Computer Laboratory, UK

Abstract. In this paper, we present an inference algorithm for OQL which both identifies the most general type of a query in the absence of schema type information, and derives the minimum type requirements a schema should satisfy to be compatible with this query. Our algorithm is useful in any database application where heterogeneity is encountered, for example, schema evolution, queries addressed against multiple schemata, inter-operation or reconciliation of heterogeneous schemata. Our inference algorithm is technically interesting as it concerns an object functional language with a rich semantics and complex type system. More precisely, we have devised a set of constraints and an algorithm to resolve them. Our resulting type inference system for OQL should be useful in any open distributed, or even semi-structured, database environment.

1 Introduction

The ODMG Standard [6] (hereafter referred to as simply the Standard) presents, rather informally, some details of a type system for checking OQL queries using type information about the classes, extents, named objects and query definitions from a given database schema. Recently there have been some efforts to formalise this type system [2, 3]. This paper builds on our earlier work [3] and considers the problem of inferring the most general type of an OQL query in the absence of any schema information.

For example, consider the following OQL definition and query:

```
define Dept_Managers(dept) as
  select e
  from   Employees as e
  where  e.position="manager" and e.department=dept;

select d
from   Departments as d
where  count(Dept_Managers(d))>5
```

This query yields those departments that have more than five managers. It is interesting to notice that this information could be drawn by running the query against databases with significantly different schemata. For instance, consider schema A, which has two classes, *Employee* and *Department*, defined as follows.

```

class Employee (extent Employees)      class Department (extent Departments)
{ attribute string    name;            { attribute string id;}
  attribute string    position;
  attribute int       year_of_birth;
  attribute float     salary;
  attribute Department department;}

```

On the other hand, consider a second schema B, which has a class `Employee` and a named collection object `Departments` of type `List(int)`.

```

class Employee (extent Employees)
{ attribute string name;
  attribute string position;
  attribute int    department;}

```

The query could potentially run against both A and B without causing any type errors. In the case of schema A, the result of the query would be a bag of `Department` objects. In a database with schema B, the result of the query would be a bag of integers. Two vital questions arise at this point. First, how we can draw limits, or put restrictions, on the properties of a schema, so that a certain query is well-typed with respect to it? Second, what information we can derive about the type of the result of the query, supposing that we have no specific schema in mind? In this paper, we study these two questions in detail, but first let us consider the setting where this could be important.

For example, this information could be exploited in distributed database applications. Suppose we have time critical queries addressed against multiple schemata. If frequent updates on parts of these schemata are likely to occur, then many of the queries will inevitably fail to be executed. In order to avoid this situation, we should register interest in specific updates of each schema—at least in those that would affect the critical queries—and resolve the type incompatibility in due course and not at the time the queries get executed.

Our work is equally useful in contexts where we need to achieve inter-operation between heterogeneous sources. There has been a lot of research on reconciling schemata with semantic heterogeneity [4, 7]. One approach to this problem identifies the semantic inconsistencies of the ontologies in different domains and creates a global ontology that combines all of them. Another approach identifies the intersection of domains where the inconsistencies occur and tries to resolve them by introducing matching rules between them. In both cases, queries that are initially written to be executed on one domain need to be rephrased to fit the needs of more domains. Knowing the schema requirements of a query and the schema mappings to a (global or just different) ontology, the task of rephrasing queries becomes a trivial automatic process. Suppose that a group of airline companies cooperate to create a single uniform system for booking tickets. In order to do that they define a global ontology that is very close to each of the distinct ontologies. Each query is initially phrased to conform to the global ontology and is then transformed to appropriate queries addressed to the individual schemata. The transformation is much easier to perform if besides the schema mappings (from one ontology to the other), we are aware of the query

schema requirements. The latter effectively point out the exact mappings we need to use.

This paper is organised as follows. In section 2 we recall our earlier [3] definition of a core OQL—a fragment of the language defined in the Standard, but which has the same expressive power. We give a brief overview of the type system of OQL, including the notion of subtyping. In section 3, we study the type system of our inference model introducing a new relation between types, called *more specific*. In section 4, we describe the kinds of constraints generated by our type inference algorithm, and in section 5, we present an algorithm for resolving these constraints. The core of our inference system, the inference rules, are given in section 6. Finally, in section 7, we present the inference algorithm which yields the most general type of a query along with its type schema requirements.

2 Core OQL

In this section we fix the syntax and type system for OQL. This is explained in greater detail in an earlier paper [3]; space restrictions mean that here we simply give the syntax for queries and definitions¹ in Figure 1. An OQL **program** consists of a number (maybe zero) of named definitions followed by a query.

The syntax for OQL types is also given in Figure 1. In what follows we will write $\text{Col}(\sigma)$, to denote an arbitrary collection type (set, bag, list or array), with elements of type σ .

Implicit in the ODMG model is a notion of subtyping; the underlying idea is that σ is said to be a subtype of τ , if a value of type σ can be used in any context in which a value of type τ is expected. This we shall write $\sigma \leq \tau$ and define as the least relation closed under the rules given in Figure 1.

We use the \sqsubseteq symbol to denote single inheritance between two classes, referred to in the Standard as the “*derives from*” relation. To simplify our presentation we do not consider interfaces.

An interesting feature of our subtype relation is the treatment of structures. A type $\sigma = \mathbf{struct}(l_1:\sigma_1, \dots, l_m:\sigma_m)$ is considered to be a subtype of $\tau = \mathbf{struct}(l_1:\tau_1, \dots, l_n:\tau_n)$ if τ is obtained from σ by dropping some labels. (In fact, we generalise this a little and also allow subtyping between the label types). This so-called width-subtyping is an extension to the Standard, but we feel it offers considerable flexibility.

The type system and the subtype relation are given in detail in an earlier paper [3]. In that work, we aimed at deriving the type of an OQL query given specific schema information. In order to do that, we defined typing judgements of the form:

$$\mathcal{S}; \mathcal{D}; \mathcal{N}; \mathcal{Q} \vdash q: \sigma$$

where \mathcal{S} are the class definitions, \mathcal{D} are the persistent query definitions and \mathcal{N} are the named objects of a specific schema. \mathcal{Q} represents the query typing

¹ Naturally as we are interested in *inferring* types we drop the requirement that definition parameters be explicitly typed.

Queries $q ::= b \mid f \mid i \mid c \mid s$

- | x
- | $\text{bag}(q, \dots, q) \mid \text{set}(q, \dots, q) \mid \text{list}(q, \dots, q) \mid \text{array}(q, \dots, q)$
- | $\text{struct}(l: q, \dots, l: q)$
- | $C(l: q, \dots, l: q) \mid q.l \mid (C)q$
- | $q[q] \mid q \text{ in } q \mid q() \mid q(q, \dots, q)$
- | $\text{forall } x \text{ in } q: q \mid \text{exists } x \text{ in } q: q$
- | $q \text{ binop } q \mid \text{unop}(q)$
- | $\text{select} [\text{distinct}] q$
- | $\text{from } (q \text{ as } x, \dots, q \text{ as } x)$
- | $\text{where } q$
- | $[\text{group by } (l: q, \dots, l: q)]$
- | $[\text{having } q]$
- | $[\text{order by } (q \text{ asc} \mid \text{desc}, \dots, q \text{ asc} \mid \text{desc})]$

Definitions $d ::= \text{define } x \text{ as } q$

- | $\text{define } x(x, \dots, x) \text{ as } q$

Here b, f, i, c, s range over booleans, floats, integers, characters and strings respectively, x is taken from a countable set of identifiers, l is taken from a countable set of labels, and C ranges over a countable set of class names. We assume sets of unary and binary operators, ranged over by unop and binop respectively.

Types $\sigma ::= \text{int} \mid \text{float} \mid \text{bool} \mid \text{char} \mid \text{string} \mid \text{void}$

- | $\sigma \times \dots \times \sigma \rightarrow \sigma$
- | $\text{bag}(\sigma) \mid \text{set}(\sigma) \mid \text{list}(\sigma) \mid \text{array}(\sigma)$
- | $\text{struct}(l: \sigma, \dots, l: \sigma)$
- | C

We assume a distinguished class name `Object`.

Sub-typing

$$\begin{array}{c}
\frac{}{C \leq \text{Object}} \text{Top} \quad \frac{C \sqsubseteq C'}{C \leq C'} \text{Sub-Class} \\
\\
\frac{\sigma'_1 \leq \sigma_1 \dots \sigma'_k \leq \sigma_k \quad \tau \leq \tau'}{\sigma_1 \times \dots \times \sigma_k \rightarrow \tau \leq \sigma'_1 \times \dots \times \sigma'_k \rightarrow \tau'} \text{Sub-Fun} \quad \frac{\sigma \leq \tau}{\text{Col}(\sigma) \leq \text{Col}(\tau)} \text{Sub-Coll} \\
\\
\frac{\sigma_1 \leq \tau_1 \quad \dots \quad \sigma_k \leq \tau_k}{\text{struct}(l_1: \sigma_1, \dots, l_k: \sigma_k, \dots, l_{k+n}: \sigma_{k+n}) \leq \text{struct}(l_1: \tau_1, \dots, l_k: \tau_k)} \text{Sub-Struct} \\
\\
\frac{}{\sigma \leq \sigma} \text{Sub-Refl} \quad \frac{\sigma \leq \sigma' \quad \sigma' \leq \sigma''}{\sigma \leq \sigma''} \text{Sub-Trans}
\end{array}$$

Fig. 1. Syntax, Types and Subtyping for Core OQL

environment, i.e. it contains the types of any free identifiers in q . A simple example of the typing rules used to derive the type of a query is the following:

$$\frac{\mathcal{S}; \mathcal{D}; \mathcal{N}; \mathcal{Q} \vdash q: \text{list}(\sigma)}{\mathcal{S}; \mathcal{D}; \mathcal{N}; \mathcal{Q} \vdash \text{first}(q): \sigma} \text{First-list}$$

In the current context, we have no information about the classes, the query definitions, and the named objects, as we have no schema information. The problem we address can thus be written $? \vdash q: ?$, i.e. given an arbitrary query q , can we infer its type and the type of any supporting schemata?

3 Type and Schema Inference

In this section, we present the extended type system behind our inference algorithm and a new relation between types called *Generalisation-Specialisation* relation. We also discuss the notions of *Least Upper Bound* and *Greatest Lower Bound* of two types that occur frequently in our inference algorithm.

3.1 Extended Type System

It turns out to be convenient to extend the notion of type given in Figure 1; an example should make this clear. Consider the following:

```
define q1 as select x from Students as x;
define q2 as set(first(Students));
q1 union q2
```

Considering the query `q1 union q2` first, we can infer immediately that `q1` and `q2` should be either sets or bags of elements. To represent this we introduce a new type constructor `set/bag(-)`. Moreover, the elements of this collection cannot be of a *function* type, since the Standard does not allow functions to be members of a collection; thus we introduce the types `nonfunctional` and `function`. From the definition `q1` we infer that `Students` is some collection (set, bag, list or array) of elements of any (non-functional) type. Considering the definition `q2` we can infer further information about `Students`. As it is the argument of a `first` operation it must be an *ordered* collection, i.e. a list or an array. Again we introduce a new type constructor `list/array(-)`. In summary, the algorithm should infer that `Students` is a list or an array of a nonfunctional member type τ , and that the query `q1 union q2` is of type `bag(τ)`.²

The above example motivates our need to extend the initial *specific* types (given in Figure 1) with the following so-called *general* types.

```
{any, nonfunctional, atomic, orderable, int/float} ∪
{collection( $\tau$ ), set/bag( $\tau$ ), list/array( $\tau$ ), constructor( $l_i: \tau_i$ )} ∪
{all types from the core type system with at
least one component type being a general type, e.g. set(any)}
```

² The Standard [§4.10.11] states that merging a set and a bag results in a bag.

where τ, τ_i are *specific* or *general* types.

Given these general types the resulting type system is then as follows:

$$\begin{aligned}
 \sigma ::= & \text{int} \mid \text{float} \mid \text{bool} \mid \text{char} \mid \text{string} \mid \text{void} \\
 & \mid \sigma \times \dots \times \sigma \rightarrow \sigma \\
 & \mid \text{bag}(\sigma) \mid \text{set}(\sigma) \mid \text{list}(\sigma) \mid \text{array}(\sigma) \mid \text{struct}(l_1: \sigma, \dots, l_n: \sigma) \\
 & \mid \mathbb{C} \\
 & \mid \text{any} \mid \text{nonfunctional} \\
 & \mid \text{atomic} \mid \text{orderable} \mid \text{int/float} \\
 & \mid \text{constructor}(l_1: \sigma, \dots, l_n: \sigma) \\
 & \mid \text{collection}(\sigma) \mid \text{set/bag}(\sigma) \mid \text{list/array}(\sigma)
 \end{aligned} \tag{1}$$

This extended type system (1) is coupled with the type hierarchy illustrated in Figure 2. It is worth noting that the *general* types, which are the internal nodes of the tree, are not types that can be found in a database schema, but rather abstractions or families of types that encapsulate the common features of their children.

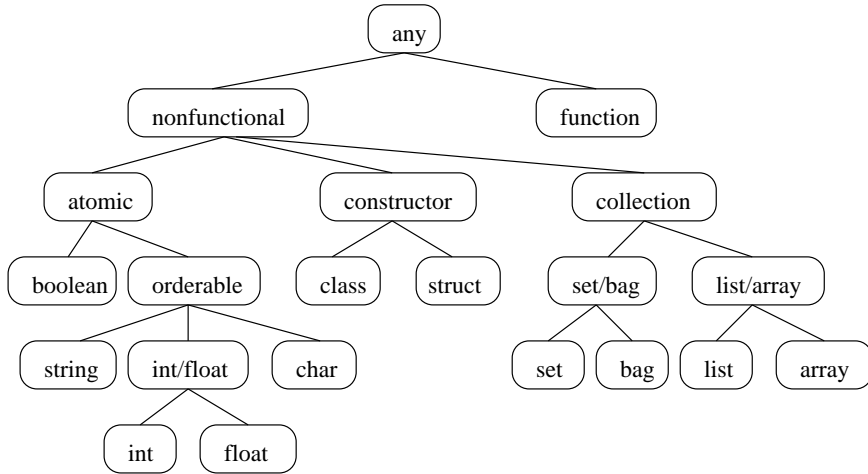


Fig. 2. Type Hierarchy

A type is said to be *specific*, if it can be derived by the type system given in Figure 1. Otherwise, it is said to be *general*. All the leaves of the hierarchy tree (in Figure 2) are specific types, if they are nullary (non parametric) types (**int**, **float**, **char**, **string**, **bool**) or if they are parametric types with all the parameter types being specific types (e.g. **set(int)**).

Given this more general type system, we need to extend our notion of subtyping given earlier. We define a new relation, GSR (Generalisation-Specialisation Relationship). Given types σ, τ , we write $\sigma \subseteq \tau$ to express that σ is *more specific*

than τ .

$$\frac{\sigma \leq \tau}{\sigma \subseteq \tau} \text{GSR - Type} \quad \frac{\sigma'_1 \subseteq \sigma_1 \cdots \sigma'_k \subseteq \sigma_k \quad \tau \subseteq \tau'}{\sigma_1 \times \cdots \times \sigma_k \rightarrow \tau \subseteq \sigma'_1 \times \cdots \times \sigma'_k \rightarrow \tau'} \text{GSR - Fun}$$

$$\frac{\text{coll}_1 \text{ is child of coll}_2 \quad \sigma \subseteq \tau}{\text{coll}_1(\sigma) \subseteq \text{coll}_2(\tau)} \text{GSR - Coll}$$

where $\text{coll}_1, \text{coll}_2$ are nodes in the sub-tree (figure 2) with root **collection** and **is child of** signifies that coll_2 is a direct or indirect parent of coll_1 or that coll_1 and coll_2 are the same.

$$\frac{\text{constr}_1 \text{ is child of constr}_2 \quad \sigma_1 \subseteq \tau_1 \cdots \sigma_k \subseteq \tau_k}{\text{constr}_1(\mathbf{l}_1: \sigma_1, \dots, \mathbf{l}_k: \sigma_k, \dots, \mathbf{l}_{k+n}: \sigma_{k+n}) \subseteq \text{constr}_2(\mathbf{l}_1: \tau_1, \dots, \mathbf{l}_k: \tau_k)} \text{GSR - Constr}$$

where $\text{constr}_1, \text{constr}_2$ are nodes in the sub-tree (figure 2) with root **constructor**, and **is child of** signifies that constr_2 is a direct parent of constr_1 or constr_1 and constr_2 are the same.

$$\frac{\text{atom}_1 \text{ is child of atom}_2}{\sigma \subseteq \text{atom}_2} \text{GSR - Atomic}$$

where atom_1 and atom_2 are nodes in the sub-tree (figure 2) with root **atomic** and **is child of** signifies that atom_2 is direct or indirect parent of atom_1 or atom_1 and atom_2 are the same.

$$\frac{}{\sigma \subseteq \sigma} \text{GSR - Refl} \quad \frac{\sigma_1 \subseteq \sigma_2 \quad \sigma_2 \subseteq \sigma_3}{\sigma_1 \subseteq \sigma_3} \text{GSR - Trans}$$

$$\frac{\sigma \neq \tau_1 \rightarrow \tau_2}{\sigma \subseteq \text{nonfunctional}} \text{GSR - NonFun} \quad \frac{}{\sigma \subseteq \text{any}} \text{GSR - All}$$

Given this definition, we can define the *Greatest Lower Bound* (GLB) and the *Least Upper Bound* (LUB) of two types τ_1 and τ_2 . Insight into these concepts can be gained through the following simple example. Consider the types $\tau_1 = \text{set}(\text{atomic})$ and $\tau_2 = \text{set}/\text{bag}(\text{int})$. The GLB of the two types is derived by taking the most specific of the collection constructors, **set**, and the most specific of the parameter types, **int**. Thus $\text{GLB}(\tau_1, \tau_2) = \text{set}(\text{int})$. Likewise, for the LUB, we take the most general of the two collection constructors, **set/bag**, and the most general of the two parameter types, **atomic**. Thus, $\text{LUB}(\tau_1, \tau_2) = \text{set}/\text{bag}(\text{atomic})$.

We may now formally present GLB and LUB. In the following definitions, we assume that $\text{constr}_1, \text{constr}_2$ are nodes in the sub-tree (figure 2) with root **constructor** and that constr_1 is a child of constr_2 or $\text{constr}_1 = \text{constr}_2$. Moreover, coll_1 and coll_2 are nodes in the sub-tree (figure 2) with root **collection** and coll_1 is a child of coll_2 or $\text{coll}_1 = \text{coll}_2$.

$$\begin{aligned}
\text{GLB}(\tau_1, \tau_2) &\stackrel{\text{def}}{=} \tau_1 \text{ if } \tau_1 \subseteq \tau_2 \wedge \tau_2 \subseteq \text{atomic} \\
\text{LUB}(\tau_1, \tau_2) &\stackrel{\text{def}}{=} \tau \text{ if } \tau_1 \subseteq \text{atomic} \wedge \tau_2 \subseteq \text{atomic} \wedge \\
&\text{(there exists no } \tau' \text{ s.t. } (\tau' \neq \tau \wedge \tau_1 \subseteq \tau' \wedge \tau_2 \subseteq \tau' \wedge \tau \subseteq \tau')) \\
\text{GLB}(\text{constr}_1(\mathbf{l}_1: \sigma_1, \dots, \mathbf{l}_k: \sigma_k), \text{constr}_2(\mathbf{l}_1: \sigma'_1, \dots, \mathbf{l}_k: \sigma'_k, \dots, \mathbf{l}_{k+n}: \sigma'_{k+n})) &\stackrel{\text{def}}{=} \\
&\text{constr}_1(\mathbf{l}_1: \text{GLB}(\sigma_1, \sigma'_1), \dots, \mathbf{l}_k: \text{GLB}(\sigma_k, \sigma'_k), \dots, \mathbf{l}_{k+n}: \sigma_{k+n}) \\
\text{LUB}(\text{constr}_1(\mathbf{l}_1: \sigma_1, \dots, \mathbf{l}_k: \sigma_k), \text{constr}_2(\mathbf{l}_1: \sigma'_1, \dots, \mathbf{l}_k: \sigma'_k, \dots, \mathbf{l}_{k+n}: \sigma'_{k+n})) &\stackrel{\text{def}}{=} \\
&\text{constr}_2(\mathbf{l}_1: \text{LUB}(\sigma_1, \sigma'_1), \dots, \mathbf{l}_k: \text{LUB}(\sigma_k, \sigma'_k)) \\
\text{GLB}(\text{coll}_1(\sigma_1), \text{coll}_2(\sigma_2)) &\stackrel{\text{def}}{=} \text{coll}_1(\text{GLB}(\sigma_1, \sigma_2)) \\
\text{LUB}(\text{coll}_1(\sigma_1), \text{coll}_2(\sigma_2)) &\stackrel{\text{def}}{=} \text{coll}_2(\text{LUB}(\sigma_1, \sigma_2)) \\
\text{GLB}(\tau_1 \times \dots \times \tau_k \rightarrow \sigma, \tau'_1 \times \dots \times \tau'_k \rightarrow \sigma') &\stackrel{\text{def}}{=} \text{LUB}(\tau_1, \tau'_1) \times \dots \times \text{LUB}(\tau_k, \tau'_k) \rightarrow \text{GLB}(\sigma, \sigma') \\
\text{LUB}(\tau_1 \times \dots \times \tau_k \rightarrow \sigma, \tau'_1 \times \dots \times \tau'_k \rightarrow \sigma') &\stackrel{\text{def}}{=} \text{GLB}(\tau_1, \tau'_1) \times \dots \times \text{LUB}(\tau_k, \tau'_k) \rightarrow \text{LUB}(\sigma, \sigma') \\
\text{LUB}(\sigma, \tau) &\stackrel{\text{def}}{=} \text{any, if } \sigma \subseteq \text{function} \wedge \tau \subseteq \text{nonfunctional} \\
\text{LUB}(\sigma_1, \sigma_2) &\stackrel{\text{def}}{=} \text{nonfunctional, if } \forall \tau_1, \tau_2. \sigma_1 \subseteq \tau_1 \wedge \sigma_2 \subseteq \tau_2 \wedge \tau_1 \neq \tau_2 \wedge \\
&\tau_1, \tau_2 \in \{\text{atomic}, \text{constructor}(), \text{collection}(\text{any})\}
\end{aligned}$$

4 Type compatibility - Constraints

The inference algorithm we present in section 7 analyses an OQL construct and infers the most general type of the query and the schema requirements that should be satisfied so that the query is well-typed. Before being able to present the inference algorithm, we first discuss an important mechanism that the algorithm is based upon—the generation of constraints. When the inference algorithm analyses a certain query construct, it often infers several relations or associations amongst the types of the query and its subqueries. These associations are given in the form of constraints.

For example, the analysis of a query $q_1 \text{ union } q_2$ would generate the constraint that the type τ of the query is the merge result of the types τ_1 and τ_2 of the two subqueries, i.e. $\tau = \text{Merge_Result}(\tau_1, \tau_2)$. In section 5, we show how this constraint is simplified and is *assimilated* in the set of the existing constraints.

Analysing the query `exists x in Customers: x.income > 40,000` our algorithm generates the following constraints. First, it introduces the constraint $\tau_1 = \text{Member_Type}(\tau_2)$, where τ_1 is the type of `x` and τ_2 is the type of `Customers`. The type of `x` is expected to be the same as that of the members of the collection `Customers`. Second, the constraint $\tau_3 = \text{Constructor_Member_Type}(\tau_1, \text{income})$ is generated, where τ_3 is the type of `x.income`. This signifies that `x` is a constructor type (a structure or a class) with at least one member `income` of type τ_3 . Another interesting constraint is that the types of `x.income` (τ_3) and of the literal `40,000` (`int`) should be compatible in the sense that they than can be compared

for inequality. This is expressed by the constraint `Greater_Less_Than_Compatible(τ_3 , int)`. Later, we will show how this constraint is simplified to the constraint $\tau_3 \subseteq \text{int/float}$.

In section 6, we give a set of inference rules, one for each query construct. Each rule starts with an existing set of constraints and generates a number (possibly zero) of new constraints. The different kinds of constraints generated by the rules in our inference algorithm are given below:

1. <code>Equality_Compatible(τ_1, \dots, τ_n)</code>	6. <code>$\tau_0 = \text{Merge_Result}(\tau_1, \tau_2)$</code>
2. <code>Greater_Less_Than_Compatible(τ_1, \dots, τ_n)</code>	7. <code>$\tau_0 = \text{Distinct_Result}(\tau_1)$</code>
3. <code>$\tau_1 = \tau_2$</code>	8. <code>$\tau = \text{Member_Type}(\sigma)$</code>
4. <code>$\tau_1 \subseteq \tau_2$</code>	9. <code>$\tau = \text{Constructor_Member_Type}(\sigma, 1)$</code>
5. <code>$\tau_0 = \text{Arith_Result}(\tau_1, \tau_2)$</code>	

We briefly explain the constraints used in our inference model. The constraint `Equality_Compatible(q_1, \dots, q_n)` is analysed in section 4.1. The constraint `Greater_Less_Than_Compatible(q_1, q_2)` is useful for ensuring typability for queries like $q_1 < q_2$. Type equality, and GSR (Generalisation-Specialisation Relationship) are handled by constraints 3 and 4. `$\tau_0 = \text{Arith_Result}(\tau_1, \tau_2)$` is needed for the type inference of queries of the form $q_1 \text{ op } q_2$ where $\text{op} \in \{+, -, /, *\}$. Likewise, `$\tau_0 = \text{Merge_Result}(\tau_1, \tau_2)$` arises as a constraint from inferencing the type of `union`, `intersect` or `except` query expressions. The constraint `$\tau_0 = \text{Distinct_Result}(\tau_1)$` implies that both τ_0 and τ_1 are collection types and the collection constructor of τ_0 is the distinct equivalent of the collection constructor of τ_1 . Moreover, `$\tau = \text{Member_Type}(\sigma)$` implies that σ is a collection type and that τ is the type of its members. Finally, the constraint `$\tau = \text{Constructor_Member_Type}(\sigma, 1)$` is used to denote that σ is a `class` or `struct` type with at least one member 1 of type τ . If σ is a `class` type then 1 can be any of its properties, relationships or methods.

4.1 Collections - Membership Type Compatibility

The first two constraints refer to type compatibility w.r.t. equality or non-equality comparison. These constraints arise in OQL constructs that involve a merge of two or more elements, or a membership test. For instance, the first constraint results from considering a query of the form `set(q_1, \dots, q_n)`, which includes the merge of n query results.

First of all, we should stress the fact that in order for two values (objects or literals) to be eligible as members of the same collection, they should be eligible for *equality comparison*. If two types are compatible (membership-wise), two values of these types may be members of a set. In order to insert an element into a set, we need to test if its value is equal to any existing value. Thus, we need to ensure that these values have types which are compatible (equality-wise). Inversely, if two types are compatible equality-wise, then their values may be inserted into any collection, therefore these types are also compatible membership-wise.

The Standard [§4.10] defines recursively when two types are compatible, and thus when elements of these types can be put in the same collection. The Stan-

dard then defines the notion of least upper bound (LUB) of two types to derive the type of the collection elements. In a context where we need to check and derive the type of a query based on *specific* type information (from a schema), this approach is sufficient and straightforward. However, in our context, where we aim to infer the type of a query without any schema information, the compatibility issue becomes more complicated. The use of LUB to infer the type of a query like `set(q1, ..., qn)` does not yield the appropriate result, for example consider the following.

```
define q1 as struct(x:12,y:30);
define q2 as element(select z from People as z where z.x=14);
set(q1,q2)
```

If we call the inference algorithm on `q1` and `q2` the inferred types (IT) would be `struct(x : int, y : int)` and `constructor(x : int)` respectively. The least upper bound of these two types is `constructor(x : int)`. Thus, the inferred type of the query `set(q1,q2)` would be `set(constructor(x : int))`.

However, the correct inferred type should be `set(struct(x : int))`, since we may not merge objects and structures in the same collection, and therefore we know that the `constructor` should be a `struct` and not a `class`.

To overcome this problem we define another relation between types, namely CUB (*Compatible Upper Bound*). Intuitively, CUB combines the behaviour of both LUB and GLB (Greatest Lower Bound). In the previous example, the CUB of the two types `IT(q1)` and `IT(q2)` would be derived by taking the most specific of the two constructor types (`constructor` and `struct`), but the least general of the element types (`(x : int)` and `(x : int, y : int)`). Before we define CUB, we define *compatibility* (Membership- or Equality- wise) for our typing system.

Compatibility is recursively defined as follows:

- τ is compatible with τ
- if σ is compatible with τ and $\text{coll}_1, \text{coll}_2 \in \{\text{collection}, \text{set/bag}, \text{list/array}, \text{set}, \text{bag}, \text{list}, \text{array}\}$ and either the collection constructors are the same or one is child of the other in the hierarchy tree then $\text{coll}_1(\sigma)$ is compatible with $\text{coll}_2(\tau)$.
- Any two class types `class_name1` and `class_name2` are compatible.
- If σ_i is compatible with $\tau_i, \forall i = 1, \dots, n$ and `constr1, constr2` $\in \{\text{constructor}, \text{struct}\}$ and no labels other than l_1, \dots, l_n are common in both constructor types then `constr1(l1: $\sigma_1, \dots, l_n: \sigma_n, l_{11}: \sigma_{11}, \dots, l_{1k}: \sigma_{1k}$)` and `constr2(l1: $\tau_1, \dots, l_n: \tau_n, l_{21}: \tau_{21}, \dots, l_{2m}: \tau_{2m}$)` are compatible.
- If σ_i is compatible with $\tau_i, \forall i = 1, \dots, n$ and `class_name` is a class type such that `Constructor_Member_Type(class_name, li) = $\tau_i, \forall i = 1, \dots, n$` , then the types `class_name` and `constructor(l1: $\sigma_1, \dots, l_n: \sigma_n, l_{11}: \sigma_{11}, \dots, l_{1k}: \sigma_{1k}$)` are compatible, provided that no labels other than l_1, \dots, l_n are common members in the two types.
- If $\sigma, \tau \in \{\text{atomic}, \text{orderable}, \text{int/float}, \text{int}, \text{float}, \text{char}, \text{string}, \text{bool}\}$ and either they are the same, or one is a child of the other in the hierarchy

- tree, or one is `int` and the other is `float` then σ is compatible with τ .
- If either of the types is `nonfunctional` and the other type is not a function type then these types are compatible.
- If either of two types is `any`, then these types are compatible.

Note that we do not define compatibility for function types, as no two function values may be members of the same collection or may be compared for equality. Only the results of function application may be considered for compatibility.

Given that two types are compatible (based on the recursive definition above), their CUB is defined recursively and in accordance with the compatibility category that they fall into.

- $\text{CUB}(\tau, \tau) = \tau$.
- If $\text{coll}_1(\sigma)$ is compatible with $\text{coll}_2(\tau)$ and coll_1 is a child of (or the same as) coll_2 , then $\text{CUB}(\text{coll}_1(\sigma), \text{coll}_2(\tau)) = \text{coll}_1(\text{CUB}(\sigma, \tau))$.
- If the types τ_1 and τ_2 are class types, then $\text{CUB}(\tau_1, \tau_2)$ is the least common superclass of the two classes.
- If the types $\sigma = \text{constr}_1(l_1: \sigma_1, \dots, l_n: \sigma_n, l_{11}: \sigma_{11}, \dots, l_{1k}: \sigma_{1k})$ and $\tau = \text{constr}_2(l_1: \tau_1, \dots, l_n: \tau_n, l_{21}: \tau_{21}, \dots, l_{2m}: \tau_{2m})$ are compatible, where $\text{constr}_1, \text{constr}_2 \in \{\text{constructor}, \text{struct}\}$ then $\text{CUB}(\sigma, \tau) = \text{constr}_1(l_1: \text{CUB}(\sigma_1, \tau_1), \dots, l_n: \text{CUB}(\sigma_n, \tau_n))$.
- If $\sigma = \text{class_name}_1$ and $\text{Constructor_Member_Type}(\text{class_name}_1, l_i) = \sigma_i$, $\forall i = 1, \dots, n$ and $\tau = \text{constructor}(l_1: \tau_1, \dots, l_n: \tau_n, l_{21}: \tau_{21}, \dots, l_{2m}: \tau_{2m})$ then $\text{CUB}(\sigma, \tau)$ is the least superclass of class_name_1 , say class_name_2 , satisfying the following condition: For all l'_j , l'_j is a property or a relationship of class_name_2 , if $\text{Constructor_Member_Type}(\text{class_name}_2, l'_j) = \phi_j$ then there exists $k, 1 \leq k \leq n$, s.t. $l'_j = l_k \wedge \text{CUP}(\tau_k, \sigma_k) \subseteq \phi_j, \forall j = 1, \dots, m, m \leq n$.
- If σ is compatible with τ , then
 1. if either of them is `int/float` or one is `int` and the other is `float` then $\text{CUB}(\sigma, \tau) = \text{int/float}$
 2. else $\text{CUB}(\sigma, \tau) = \text{GLB}(\sigma, \tau)$.
- If $\sigma = \text{nonfunctional}$ and $\tau \subseteq \text{nonfunctional}$ then $\text{CUB}(\sigma, \tau) = \tau$.
- If $\sigma = \text{any}$ then, for any type τ , $\text{CUB}(\sigma, \tau) = \tau$.

As discussed earlier, the notion of CUB is used for the inference of the types of queries like `set(q1, ..., qn)`. However, the OQL construct `q1 in q2` raises another issue of a slightly different nature. The Standard [§4.10.8.3] states that if the type of `q2` is $\text{coll}(\tau)$ then the type of `q1` should be τ . This is not the case in our context. Suppose that the type of `q2` is inferred to be `bag(struct(x: int, y: string))`; then according to the Standard `q1` should have the type `struct(x: int, y: string)`. Since a value of type `struct(x: int)` could potentially be added in the collection `q2` (that is, since `struct(x: int)` and `struct(x: int, y: float)` are compatible types), there is no reason why `q1` could not be of type `struct(x: int)` or even `struct()`.

The same situation occurs when dealing with a collection of objects of different classes. Suppose $\text{lub_class}(l_1:\sigma_1, \dots, l_n:\sigma_n)$ is the LUB of all classes of the objects in the collection and `Object` is the most general class that all other classes derive from (the top of the class hierarchy). We should be able to check whether an object of type $\text{object_class}(l'_1:\sigma'_1, \dots, l'_m:\sigma'_m)$ is a member of the collection, even if its class is not a subclass of $\text{lub_class}(l_1:\sigma_1, \dots, l_n:\sigma_n)$. This allows more queries to (safely) type-check, for example:

```
select x
from   People as x
where  x.father in School_Teachers
```

does not type-check according to the Standard. In order to be type-correct, `x.father` requires an explicit type cast, i.e. `(School_Teacher)x.father`. The problem is that this query, despite being well-typed, can generate a run-time error (if the cast does not succeed). We choose not to enforce that `x.father` has a type which is more specific than the member type or the collection `School_Teachers`. Rather, we simply ensure that `x.father` could potentially be a member of `School_Teachers`. To do this we add the constraint `Equality-Compatible(σ, τ)`, where σ is the type of `x.father` and τ is the member type of `School_Teachers`.

5 Resolving constraints

Now that we have studied the kinds of constraints that are generated by our inference algorithm, we can discuss how these constraints are resolved. When a constraint is generated by an inference rule, it is added to the set of existing constraints. If this was a simple insertion procedure, we would end up having a huge set of constraints, that would include redundant and often incomprehensible type information; there is obviously a need to resolve the inserted constraints. Due to the complexity of the type system and the expressiveness of the language, we have a wide variety of constraints, that cannot be solved using a standard unification mechanism alone [10]. In our system, the insertion of a new constraint in a set of existing constraints may have one of the following effects:

- A constraint is deleted, if it is always satisfied, e.g. the constraint `Equality-Compatible(set(int), set(float))` is always true, so it does not need to be maintained.
- A constraint raises a type error or exception, if it is never satisfied, e.g. `set(Employee(name: string)) \subseteq list/array(Employee(name: string))`.
- A constraint might be maintained as it is. This usually occurs when some of the types involved are *general*; it may be that when refined, these types no longer satisfy the constraint. Therefore, they must be preserved as required schema information. For instance, if $\tau_i \subseteq \text{set/bag}(\tau) \forall i = 0, 1, 2$ are in the set of already produced constraints, the constraint $\tau_0 = \text{Merge_Result}(\tau_1, \tau_2)$ needs to be preserved.
- A constraint is often simplified, i.e. replaced by one or more simpler constraints. For instance, the constraint `set(σ) = set(τ)` is replaced by the simpler one `$\sigma = \tau$` .

- A constraint occasionally implies one or more constraints. The latter need to be added to the set of constraints already produced. For example, the constraint `Greater_Less_Than_Compatible(τ_1, τ_2)` is inserted as a new constraint along with the implied constraints $\tau_1 \subseteq \text{orderable}$ and $\tau_2 \subseteq \text{orderable}$.

The effect varies depending on the constraint kind, the types involved in the constraint and the already existing constraints on these types. The details of the constraint resolution algorithm will appear in [11].

It is worth pointing out that the resolution of constraints could take place either at the time each constraint is generated (*gradual resolution*) or at the time all the constraints have been produced (*accumulative resolution*).

The *gradual resolution* is very simple, since it usually concerns the insertion of a few constraints whose simplification (*unification*) is straightforward. If their simplification produces new constraints then these are simplified as well, until no more constraints are produced.

The *accumulative resolution* starts from the constraints of the form $\tau_1 = \tau_2$. It simplifies them to constraints of the form `type_var = τ` and replaces `type_var` by τ in all other constraints that involve `type_var`. Then it proceeds to simplify all other kinds of constraints. If a simplification leads to more constraints, the latter are added to the set of unprocessed constraints and are simplified in due course.

6 Type Inference Rules

Having explained the type system underlying our inference model and the various constraints generated and resolved by our algorithm, we are now in a position to present the backbone of our work, the *inference rules*. Note that there is a single rule for each OQL construct, and, therefore, the use of the rules by the inference algorithm is syntax driven. In the remainder of this section, we present a substantial part of the inference rules; the complete set will be given in [11].

In the following rules, \mathcal{H} signifies the *type environment*, that is $\mathcal{H} = \{var_i : \tau_i\}$, and \mathcal{C} denotes the *constraints* added so far. The inference rules for the literal and the identifier queries are given first:

$$\frac{}{\mathcal{H}; \mathcal{C} \vdash b: \text{bool} \Rightarrow \mathcal{C}} \quad \frac{}{\mathcal{H}; \mathcal{C} \vdash i: \text{int} \Rightarrow \mathcal{C}} \quad \frac{}{\mathcal{H}; \mathcal{C} \vdash f: \text{float} \Rightarrow \mathcal{C}}$$

$$\frac{}{\mathcal{H}; \mathcal{C} \vdash c: \text{char} \Rightarrow \mathcal{C}} \quad \frac{}{\mathcal{H}; \mathcal{C} \vdash s: \text{string} \Rightarrow \mathcal{C}} \quad \frac{}{\mathcal{H} \cup \{x: \sigma\}; \mathcal{C} \vdash x: \sigma \Rightarrow \mathcal{C}}$$

There are several rules to deal with various collections (sets, bags, lists, arrays). We just give one representative rule, that concerns the query construct `set(q_1, \dots, q_n)`. As expected, the rule generates a constraint that ensures that the types of the queries are compatible equality-wise (or membership-wise). We also give the rules for accessing the first, last or i -th member of an ordered collection, as well as checking whether an element belongs in a certain collection.

The constraint $\text{Member_Type}(\sigma) = \tau$ denotes that σ is a collection (set, bag, list or array) with members of type τ .

$$\begin{array}{c}
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_1: \sigma_1 \Rightarrow \mathcal{C}_1 \dots \mathcal{H}; \mathcal{C}_{n-1} \vdash \mathbf{q}_n: \sigma_n \Rightarrow \mathcal{C}_n \\
\hline
\mathcal{H}; \mathcal{C} \vdash \text{set}(\mathbf{q}_1, \dots, \mathbf{q}_n): \text{set}(\text{CUB}(\sigma_1, \dots, \sigma_n)) \Rightarrow \mathcal{C}_n \wedge \{\text{Equality_Compatible}(\sigma_1, \dots, \sigma_n)\} \\
\hline
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_1: \sigma \Rightarrow \mathcal{C}_1 \\
\hline
\mathcal{H}; \mathcal{C} \vdash \text{first}(\mathbf{q}_1): \phi \Rightarrow \mathcal{C}_1 \wedge \{\text{Member_Type}(\sigma) = \phi\} \wedge \{\sigma \subseteq \text{list/array}(\phi)\} \\
\hline
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_1: \sigma \Rightarrow \mathcal{C}_1 \quad \mathcal{H}; \mathcal{C}_1 \vdash \mathbf{q}_2: \tau \Rightarrow \mathcal{C}_2 \\
\hline
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_1[\mathbf{q}_2]: \phi \Rightarrow \mathcal{C}_2 \wedge \{\tau = \text{int}\} \wedge \{\text{Member_Type}(\sigma) = \phi\} \wedge \{\sigma \subseteq \text{list/array}(\phi)\} \\
\hline
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_1: \sigma \Rightarrow \mathcal{C}_1 \quad \mathcal{H}; \mathcal{C}_1 \vdash \mathbf{q}_2: \tau \Rightarrow \mathcal{C}_2 \\
\hline
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_1 \text{ in } \mathbf{q}_2: \text{bool} \Rightarrow \mathcal{C}_2 \wedge \{\text{Equality_Compatible}(\sigma, \text{Member_Type}(\tau))\}
\end{array}$$

The rules for constructing a structure or an object, as well as for accessing a member of a structure or an object are given below. The constraint $\text{Constructor_Member_Type}(\sigma, \ell) = \tau$ denotes that type σ is a class or a structure with a member called ℓ of type τ .

$$\begin{array}{c}
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_1: \sigma_1 \Rightarrow \mathcal{C}_1 \dots \mathcal{H}; \mathcal{C}_{n-1} \vdash \mathbf{q}_n: \sigma_n \Rightarrow \mathcal{C}_n \\
\hline
\mathcal{H}; \mathcal{C} \vdash \text{class_name}(\mathbf{l}_1: \mathbf{q}_1, \dots, \mathbf{l}_n: \mathbf{q}_n): \text{class_name} \Rightarrow \mathcal{C}_n \wedge \\
\{\sigma_1 \subseteq \text{Constructor_Member_Type}(\text{class_name}, \mathbf{l}_1)\} \wedge \dots \wedge \\
\{\sigma_n \subseteq \text{Constructor_Member_Type}(\text{class_name}, \mathbf{l}_n)\} \\
\hline
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_1: \sigma_1 \Rightarrow \mathcal{C}_1 \dots \mathcal{H}; \mathcal{C}_{n-1} \vdash \mathbf{q}_n: \sigma_n \Rightarrow \mathcal{C}_n \\
\hline
\mathcal{H}; \mathcal{C} \vdash \text{struct}(\mathbf{l}_1: \mathbf{q}_1, \dots, \mathbf{l}_n: \mathbf{q}_n): \text{struct}(\mathbf{l}_1: \sigma_1, \dots, \mathbf{l}_n: \sigma_n) \Rightarrow \mathcal{C}_n \\
\hline
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_1: \tau \Rightarrow \mathcal{C}_1 \\
\hline
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_1.\mathbf{l}: \sigma \Rightarrow \mathcal{C}_1 \wedge \{\text{Constructor_Member_Type}(\tau, \mathbf{l}): \sigma\}
\end{array}$$

The inference rules for the existential and the universal quantification follow. It is worth noting that the variable x is bound in query \mathbf{q}_2 to the member type of the collection \mathbf{q}_1 .

$$\begin{array}{c}
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_1: \sigma_1 \Rightarrow \mathcal{C}_1 \quad \mathcal{H} \cup \{x: \phi\}; \mathcal{C}_1 \vdash \mathbf{q}_2: \sigma_2 \Rightarrow \mathcal{C}_2 \\
\hline
\mathcal{H}; \mathcal{C} \vdash \text{exists } x \text{ in } \mathbf{q}_1: \mathbf{q}_2: \text{bool} \Rightarrow \mathcal{C}_2 \wedge \{\sigma_2 = \text{bool}\} \wedge \{\text{Member_Type}(\sigma_1) = \phi\} \\
\hline
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_1: \sigma_1 \Rightarrow \mathcal{C}_1 \quad \mathcal{H} \cup \{x: \phi\}; \mathcal{C}_1 \vdash \mathbf{q}_2: \sigma_2 \Rightarrow \mathcal{C}_2 \\
\hline
\mathcal{H}; \mathcal{C} \vdash \text{forall } x \text{ in } \mathbf{q}_1: \mathbf{q}_2: \text{bool} \Rightarrow \mathcal{C}_2 \wedge \{\sigma_2 = \text{bool}\} \wedge \{\text{Member_Type}(\sigma_1) = \phi\}
\end{array}$$

An interesting set of rules concerns the application of methods with or without parameters. The inferred types of the queries used as arguments are not constrained to be the same as the types of the parameters of the method involved. They only need to be their subtypes.

$$\begin{array}{c}
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_1: \sigma \Rightarrow \mathcal{C}_1 \\
\hline
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_1(): \phi \Rightarrow \mathcal{C}_1 \wedge \{\sigma = \text{unit} \rightarrow \phi\} \\
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_0: \sigma_0 \Rightarrow \mathcal{C}_0 \quad \dots \quad \mathcal{H}; \mathcal{C}_{n-1} \vdash \mathbf{q}_n: \sigma_n \Rightarrow \mathcal{C}_n \\
\hline
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_0(\mathbf{q}_1, \dots, \mathbf{q}_n): \phi \Rightarrow \mathcal{C}_n \wedge \{\sigma_0 = \tau_1 \times \dots \times \tau_n \rightarrow \phi\} \wedge \{\sigma_1 \subseteq \tau_1\} \wedge \dots \wedge \{\sigma_n \subseteq \tau_n\}
\end{array}$$

The rules concerning the query constructs `q binop q` or `unop q` are omitted for space reasons. We finish by giving the rule for a simple `select` query; the judgements dealing with a `group by` or an `order by` clause can be found in [11].

$$\begin{array}{c}
\mathcal{H}; \mathcal{C} \vdash \mathbf{q}_1: \sigma_1 \Rightarrow \mathcal{C}_1 \\
\mathcal{H} \cup \{\mathbf{x}_1: \tau_1\}; \mathcal{C}_1 \vdash \mathbf{q}_2: \sigma_2 \Rightarrow \mathcal{C}_2 \\
\vdots \\
\mathcal{H} \cup \{\mathbf{x}_1: \tau_1, \dots, \mathbf{x}_{n-1}: \tau_{n-1}\}; \mathcal{C}_{n-1} \vdash \mathbf{q}_n: \sigma_n \Rightarrow \mathcal{C}_n \\
\mathcal{H} \cup \{\mathbf{x}_1: \tau_1, \dots, \mathbf{x}_n: \tau_n\}; \mathcal{C}_n \vdash \mathbf{q}_{01}: \sigma_{01} \Rightarrow \mathcal{C}_{01} \\
\mathcal{H} \cup \{\mathbf{x}_1: \tau_1, \dots, \mathbf{x}_n: \tau_n\}; \mathcal{C}_{01} \vdash \mathbf{q}_{00}: \sigma_{00} \Rightarrow \mathcal{C}_{00} \\
\hline
\mathcal{H}; \mathcal{C} \vdash \text{select } \mathbf{q}_{00} \text{ from } \mathbf{q}_1 \text{ as } \mathbf{x}_1, \dots, \mathbf{q}_n \text{ as } \mathbf{x}_n \text{ where } \mathbf{q}_{01}: \text{bag}(\sigma_{00}) \Rightarrow \mathcal{C}_{00} \wedge \\
\{\text{Member_Type}(\sigma_1) = \tau_1\} \wedge \dots \wedge \{\text{Member_Type}(\sigma_n) = \tau_n\} \wedge \{\sigma_{01} = \text{bool}\}
\end{array}$$

7 Inference Algorithm

Having given an overview of the type system, the constraints and the rules involved in our inference model, we may now present the core of our work, which is the inference algorithm. The algorithm takes as input a query `q` and returns its inferred type, as well as a pair $(\mathcal{H}, \mathcal{C})$ of a type environment and its constraints. This pair is a synopsis of the requirements a schema should satisfy so that the query `q` can be executed against it without any type-errors.

1. For each free variable `var` in the query `q`, $\mathcal{H} = \mathcal{H} \cup \{\text{var}: \text{new_type_var}\}$. Initially $\mathcal{C} = \{\}$.
2. Based on the construct of the query `q` recursively apply the appropriate inference rule.
3. Depending on the unification strategy, either simplify the constraints as soon as they are produced (gradual unification) or simplify them all in the end after having applied all the inference rules. If the unification process produces a type-error then the query is not typable and the algorithm is interrupted.
4. The final \mathcal{H}, \mathcal{C} include the requirements a schema should have to be compatible with the query `q`. The type of `q`, which is the type inferred by the outer inference rule, also satisfies the constraints \mathcal{C} .

8 Related Work and Conclusions

Fundamental to the work described in this report is the type system for ODMG OQL described in an earlier paper [3]. Alagić [2] independently gave a number

of typing rules for OQL; see our earlier paper for a comparison. The canonical reference for work on type systems for database programming languages is the work of Buneman and Ohori [5]. The goals of their type inference algorithm are identical to ours; both approaches infer the most general type of an expression (if one exists) without accessing any schema information, and in this sense determine the constraints placed on the schema by the query. However, the underlying languages, the type systems, and some parts of the inference algorithms differ considerably. Buneman and Ohori introduce kinded types to infer the type of a record based on selections of fields on this record. Instead, we use the notion of *general* types; in this way we are able to express general type information not only for records, but also for parametric collection types, structures and classes. Moreover, due to the syntax of OQL, we define a wider variety of constraints than those introduced in their framework and therefore a different algorithm to resolve them.

The work most related to ours arises from studying type systems for object oriented programming languages, see for example [8, 9, 1]. However none of these studies consider the various issues arising from studying *database* type systems; for example, the complications arising from combining parametric collection types with subtyping.

Much research on schema evolution, schema inter-operation, distributed or semi-structured database applications has pointed out that there is a need to run queries in the presence of changing or heterogeneous schemata, or even in the absence of specific schema information. Our work addresses this problem, by proposing an inference algorithm for the ODMG query language OQL. This algorithm infers the most general type of an OQL query and derives the schema information required so that the query can be executed against it without any type errors. In contrast to other work, we deal with a rather complex type system, which includes atomic types, structures, classes, various (parameterised) collection types (set, bag, list, array) and function types. This, in connection with the rich semantics of OQL, results in the generation of a wide variety of constraints by the inference rules. We discuss the semantics of these constraints and provide a mechanism for their solution. Finally, we present a set of inference rules for OQL, which is the core of our type inference algorithm. Based on our experience, this algorithm, as well as all the formalisms prior to it, are easy to implement, and hence, we believe that they could prove to be useful in many applications.

Acknowledgements

Trigoni is funded by the State Scholarships Foundation of Greece and the National Bank of Greece.

References

1. O. Agesen and U. Holzle. Type feedback vs. concrete type inference: a comparison of optimization techniques for object-oriented languages. In *OOPSLA*, pages 91–

- 107, 1995.
2. S. Alagić. Type checking OQL queries in the ODMG type systems. *ACM Transactions on Database Systems*, 24(3):319–360, September 1999.
 3. G.M. Bierman and A. Trigoni. Towards a formal type system for ODMG OQL. Technical Report 497, University of Cambridge, Computer Laboratory, October 2000.
 4. M.W. Bright, A.R. Hurson, and S. Pakzad. Automated resolution of semantic heterogeneity in multidatabases. *ACM Transactions on Database Systems*, 19(2), 1994.
 5. P. Buneman and A. Ohori. Polymorphism and type inference in database programming. *ACM Transactions on Database Systems*, 21(1):30–76, March 1996.
 6. R.G.G. Cattell et al. *The Object Data Standard: ODMG 3.0*. Morgan Kaufmann, 2000.
 7. R. Hull. Managing semantic heterogeneity in databases: A theoretical perspective. In *PODS*, 1997.
 8. J. Palsberg and M. I. Schwartzbach. Object-oriented type inference. In *OOPSLA*, pages 146–161, 1991.
 9. J. Plevyak and A.A. Chien. Precise concrete type inference for object-oriented languages. In *OOPSLA*, pages 324–340, 1994.
 10. J. A. Robinson. A machine-oriented logic based on the resolution principle. *Journal of the ACM*, 12(1):23–41, January 1965.
 11. A. Trigoni. Phd thesis, to appear. 2001.