

Number 547



**UNIVERSITY OF
CAMBRIDGE**

Computer Laboratory

Semantic optimization of OQL queries

Agathoniki Trigoni

October 2002

15 JJ Thomson Avenue
Cambridge CB3 0FD
United Kingdom
phone +44 1223 763500
<http://www.cl.cam.ac.uk/>

© 2002 Agathoniki Trigoni

This technical report is based on a dissertation submitted October 2001 by the author for the degree of Doctor of Philosophy to the University of Cambridge, Pembroke College.

Technical reports published by the University of Cambridge Computer Laboratory are freely available via the Internet:

<http://www.cl.cam.ac.uk/TechReports/>

Series editor: Markus Kuhn

ISSN 1476-2986

Abstract

This work explores all the phases of developing a query processor for OQL, the Object Query Language proposed by the Object Data Management Group (ODMG 3.0). There has been a lot of research on the execution of relational queries and their optimization using syntactic or semantic transformations. However, there is no context that has integrated and tested all the phases of processing an object query language, including the use of semantic optimization heuristics. This research is motivated by the need for query execution tools that combine two valuable properties: i) the expressive power to encompass all the features of the object-oriented paradigm and ii) the flexibility to benefit from the experience gained with relational systems, such as the use of semantic knowledge to speed up query execution.

The contribution of this work is twofold. First, it establishes a rigorous basis for OQL by defining a type inference model for OQL queries and proposing a complete framework for their translation into calculus and algebraic representations. Second, in order to enhance query execution it provides algorithms for applying two semantic optimization heuristics: constraint introduction and constraint elimination techniques. By taking into consideration a set of association rules with exceptions, it is possible to add or remove predicates from an OQL query, thus transforming it to a more efficient form.

We have implemented this framework, which enables us to measure the benefits and the cost of exploiting semantic knowledge during query execution. The experiments showed significant benefits, especially in the application of the constraint introduction technique. In contexts where queries are optimized once and are then executed repeatedly, we can ignore the cost of optimization, and it is always worth carrying out the proposed transformation. In the context of adhoc queries the cost of the optimization becomes an important consideration. We have developed heuristics to estimate the cost as well as the benefits of optimization. The optimizer will carry out a semantic transformation only when the overhead is less than the expected benefit. Thus transformations are performed safely even with adhoc queries. The framework can often speed up the execution of an OQL query to a considerable extent.

To my family

Acknowledgements

I would like to thank my supervisor, Dr Ken Moody, for his invaluable guidance, sustained interest and encouragement during the three years of this research. In spite of his busy schedule, he always made himself readily available to discuss a variety of issues and to assist me in the shaping of my thesis. I am also grateful to Dr Gavin Bierman, my second supervisor, for teaching me type theory and formalization and for patiently helping me with a variety of difficult theoretical problems. The contribution of Dr Jean Bacon must also be acknowledged: working in a slightly different research area, she often provided me with a much needed external viewpoint. I also owe much to Robin Fairbarns and Graham Titmus who provided me with extensive software and hardware support.

This research would never have been possible without the kind and joint financial support from the Greek State Scholarships Foundation (IKY) and my employer, the National Bank of Greece (NBG). In particular, I wish to thank my IKY advisor, Dr Alexiou, for his insightful advice and encouraging progress reports at the end of each year. I am also much indebted to my superiors in the IT Division of NBG, K. Economopoulos, K. Marinakis and E. Synodinou, for promoting my application for graduate studies, and to the Human Resource Development Division for arranging my leave of absence.

I also feel compelled to thank my many friends and colleagues in the OPERA Group. I am grateful to Oliver and Chaoying for introducing me to the OPERA Group research and to Alan, Alexis, Andrei, Andras, Aras, Brian, Peter, Walt and Wei for the enjoyable experiences that we shared together.

On a more personal note, I would like to thank Constantin for making Cambridge seem like a wonderful place even at the worst of times. I greatly benefited from both his practical and emotional support. Back home, I am indebted to my friend Spyros for acting on my behalf on a variety of issues and for making Greece seem a little closer than it really was. I am also grateful to my flatmates and other friends in Cambridge for those precious relaxing moments away from my computer.

Finally, I owe a great deal to my family and in particular to my parents for their love and unceasing support. I may never have undertaken this thesis had they not always encouraged me to pursue my interests and inspired me with their idealistic view of research.

Contents

1	Introduction	15
1.1	Relevant Research Areas	16
1.2	Research Context	16
1.3	Existing Object-Oriented Database Systems	17
1.4	Type Checking vs Type Inference	19
1.5	Calculus and Algebraic Representations	19
1.6	Semantic optimization using association rules	20
1.7	Dissertation outline	21
2	Applications and Related Work	23
2.1	Applications that motivate this research	23
2.2	Previous work related to this research	25
2.2.1	Ensuring the typability of an OQL query - Related Work	25
2.2.2	Intermediate representations of OQL - Related Work	27
2.2.3	Association Rules and Semantic Optimization - Related Work	29
2.3	Extending the previous work - Discussion	32
3	Type Inference of OQL	33
3.1	Introduction	33
3.2	Core OQL	35
3.3	Type and Schema Inference	38
3.3.1	Extended Type System	38
3.4	Type compatibility - Constraints	41
3.4.1	Collections - Membership Type Compatibility	42
3.5	Resolving constraints	46
3.6	Type Inference Rules	47
3.7	Inference Algorithm	49

3.8	Discussion	49
4	OQL calculus representation	51
4.1	Monoid Comprehension Calculus	51
4.1.1	Monoids	52
4.1.2	Monoid Homomorphisms	53
4.1.3	Monoid Comprehensions	54
4.2	Extended Calculus: syntax and monoids	55
4.2.1	Calculus Syntax	55
4.2.2	Monoids	59
4.2.3	Translation rules	61
4.3	Normalization of calculus expressions	65
4.3.1	Normalization Rules	65
4.3.2	Examples of transforming calculus expressions	69
4.3.3	Normalization Algorithm	71
4.4	Discussion	79
5	OQL algebraic representation	81
5.1	Algebraic Operators	81
5.2	Translating monoid calculus to algebraic expressions	89
5.2.1	Translation of Monoid Comprehensions	90
5.3	Benefits of the algebraic representation: Examples	97
5.4	Conclusion	101
6	Using Association Rules to Optimize Queries Semantically	103
6.1	Introduction	103
6.2	Converting query predicates to rule constraints	105
6.2.1	Heuristic H1	106
6.2.2	Heuristic H2	108
6.2.3	Discussion	109
6.3	Identifying associations between constraints	109
6.3.1	Algorithm	110
6.4	Identifying optimization solutions	116
6.4.1	Heuristic H1	116
6.4.2	Heuristic H2	117
6.5	Optimizing OQL queries	118
6.5.1	Heuristic H1	119

6.5.2	Heuristic H2	119
6.6	Discussion	119
6.7	Conclusion	120
7	Complexity and Optimization	123
7.1	Experimental Framework	123
7.2	Complexity of step 1	125
7.2.1	Initial version of step 1	125
7.2.2	Optimization of Step 1	126
7.3	Complexity of Step 2	128
7.3.1	Substep 2.1 - Absorbing simple constraints in path annotations	129
7.3.2	Substep 2.2 - Combining path annotations	131
7.4	Combining steps 1 and 2 optimizes the algorithm significantly	135
7.5	Overall performance of optimizing queries semantically	138
7.6	Discussion	142
8	Conclusion and Future Work	143
8.1	Summary	143
8.2	Future work	146
8.3	Conclusion	148
A	Inference Rules	149
B	Translating OQL to the monoid calculus	151
C	Edge traversals in the backtracking algorithm	155
C.1	Edge traversals given a graph with V vertices	155
C.2	Edge traversals given a graph with E edges	156
D	Implementation of OQL Processor	159

List of Figures

3.1	Syntax and Types for Core OQL	36
3.2	Subtyping for core OQL	37
3.3	Type Hierarchy	39
4.1	Syntax of the calculus OQL representation	56
4.2	Translation of <code>select</code> queries to their calculus representation	62
4.3	Normalization Algorithm	71
4.4	Normalization of the OQL construct <code>not (e)</code>	72
4.5	Normalization of monoid comprehensions	72
4.6	Normalizing comprehensions from qualifier q_i , where q_i is a generator ($u \leftarrow e$) . . .	73
4.7	Normalizing comprehensions from qualifier q_i , where q_i is a filter (e)	74
4.8	Normalizing comprehensions from qualifier q_i , where q_i is a binding ($u \equiv e$)	74
4.9	Normalizing the head of a comprehension	75
5.1	Syntax of the algebraic OQL representation	82
5.2	Translation of calculus to algebraic constructs in Group A (figure 5.1)	89
5.3	Implementing \mathcal{S} for various calculus constructs	93
6.1	Association rules relevant to the class <code>Employee</code>	106
6.2	Graph of Constraints	108
6.3	An abstract view of the algorithm	109
6.4	Paths from relaxed constraints to the target constraint $C_{4.R_{41}}$	111
6.5	Truth table which shows the correctness of formula 6.2	112
6.6	Correlations of query constraints	116
6.7	Correlations between query constraints and new indexed constraint	118
7.1	Graph of constraints forming a tree towards the target constraint	125
7.2	Cost of finding incomplete paths from sources to targets in step 1	127
7.3	Relating the cost of step 2.1 to the number of absorbed constraints	130

7.4	Cost of sorting initial and new path annotations	132
7.5	Cost of finding new combinations in step 2.2	134
7.6	Cost of forming combinations in step 2.2	134
7.7	Cost of optimized algorithm	137
7.8	Relating the query execution benefits to the optimization cost and the evaluation benefits (top view)	139
7.9	Relating the query execution benefits to the optimization cost and the evaluation benefits (1st side view)	140
7.10	Relating the query execution benefits to the optimization cost and the evaluation benefits (2nd side view)	140
C.1	The target constraint has n_0 incoming links	156

Chapter 1

Introduction

Query processing and its optimization have been two of the most popular areas of research in the database community. A lot of experience has accumulated, mostly on the execution of relational queries and their optimization using syntactic transformations. Recently, much interesting work has focused on semantic query optimization, but relatively few frameworks have been built to test the effectiveness of these techniques. To my knowledge, there is no research context that has developed and implemented all the phases of processing an object query language; further, relatively few attempts have been made towards using data mining techniques, e.g. association rules, in order to optimize queries semantically.

The comparative analysis of relational and object-oriented database systems is one of the most controversial database issues [CLC95, DD88]. The efficiency of data storage and query execution is the main advantage of the relational systems; however the advocates of the object-oriented world claim that their paradigm is more expressive. The vast experience of optimization techniques and index usage is a strong argument used by the supporters of RDBS; the advent of object-oriented programming languages and their natural connection with OODBS is an important counter-argument. Object-relational databases try to offer the best aspects of the two database paradigms [CMN95]. They allow the user to create persistent objects, but they use a relational physical organisation to store them.

The limited experience in developing purely object-oriented database systems means that many query processors have limited capabilities and several deficiencies. The work reported in this thesis was motivated by the lack of query engines that possess two valuable properties: i) the expressive power to encompass all the features of the object-oriented paradigm and ii) the flexibility to benefit from experience gained from relational systems, such as the use of semantic knowledge to speed up query execution.

The objective of this work is to explore all the phases of developing an object query processor

and to provide algorithms for applying two semantic optimization heuristics. The implementation of this framework enables us to measure the cost and the benefits of putting semantic knowledge in the service of query execution. It should be pointed out that the conclusions drawn from the experiments are equally applicable to relational query processors.

1.1 Relevant Research Areas

Developing a query processing framework and enriching it with semantic optimization techniques involves using techniques and resolving issues from disparate research areas.

Various aspects of type theory are studied in order to ensure the typability of a query. In particular, type checking is a means of exploring whether a certain query can be executed against a particular schema without causing any type-errors. Type inference is used to infer the type of a query without accessing any particular schema information. It can be used to identify the requirements that a schema should satisfy to be compatible with a certain query.

Part of this dissertation also involves work in translating a query into various intermediate representations before it actually gets executed. This entails not only translating a query from one form to another, but also optimizing it at each stage. A lot of work has been done in this area in the relational framework, but relatively few attempts have been made to do the same in an object-oriented context.

Semantic optimization is a third important aspect of our query processing framework. The underlying idea is to optimize a query based on semantic knowledge in the form of association rules. This first requires identifying the rules using standard data mining methods; these rules must then be combined into forms suitable for the optimization purposes.

Finally, assessing the performance benefits of semantic optimization requires estimating the cost and the scalability of the methods developed. In the last part of this thesis we look at the complexity of the algorithms involved.

1.2 Research Context

The query language used in our framework is OQL, the Object Query Language proposed as part of the Object Data Standard (version 3.0) published by the ODMG [CBB⁺00]. This standard provides specifications for storing and retrieving objects in databases. Its objective is to encompass all the important features of the OO paradigm and to ensure portability across database systems that comply with it. It consists of specifications of the following components:

- the Object Model defines the type system of the object paradigm,

- the Object Definition Language (ODL) is a specification language used to define the classes of a particular schema,
- the Object Interchange Format (OIF) is a specification language used to dump and load the state of the database to and from a file,
- the Object Query Language (OQL) is a functional language that allows the user to query objects,
- C++, Java, Smalltalk bindings enable developers to write code in the respective language to manipulate persistent objects.

A general criticism of a previous release (1.2) of the Object Standard can be found in [Ala97]. In this thesis, we focus on OQL, since it is the query language studied throughout this dissertation. It is a functional language, in which queries can be composed to an arbitrary depth, as long as their operands and results respect the type system defined in the Object Model. Many of the constructs of OQL are syntactically similar to those of SQL-92; however, OQL has a number of object-oriented features. First, it deals with complex objects, i.e. objects containing or being related to other objects, or containing collections (or structures) of objects. Secondly, it enables users to navigate over the objects using path expressions of arbitrary length. Thirdly, OQL allows the comparison of objects based on their object identity instead of their content values. Fourthly, it enables method invocations within a query. Fifthly, it manipulates a series of collections (sets, bags, lists, arrays, dictionaries) with the aid of additional language constructs. Sixthly, it handles polymorphic collections, i.e. collections with objects belonging to different classes. The late binding mechanism enables users to perform generic operations on the elements of these collections. The casting mechanism allows them to declare explicitly the class of an object, going down the class hierarchy.

The expressive power of OQL makes it an interesting query language to develop. In addition to its independent use as a query language, it can potentially be embedded in any object-oriented programming language. The object-oriented nature of OQL makes the gap between the query and the programming language much narrower than in the relational context. In certain cases, it allows the programmer to perform exactly the same operation by either using the commands available in the binding of the programming language, or by embedding an OQL query in it.

1.3 Existing Object-Oriented Database Systems

There currently exist several companies who claim to have designed OODB systems compliant with the ODMG standard. For instance, SunMicrosystems, POET, Object Design and Objectivity and

Versant have developed Java Blend, Poet Object Server, ObjectStore, Objectivity/DB and Versant Developer Suite respectively.

Java Blend enables users to write database applications for any database entirely in the Java programming language. Developers do not need to know how data are structured in a database, they should just be aware of the business logic underlying their application. The main idea behind building this product was to avoid using any query language resembling SQL.

ObjectStore enables the database to be accessed through two language bindings, for Java and for C++. It can either be used as a middle-tier Data Server that links to a conventional relational database, or as a real object-oriented database system. The rationale behind its development is similar to Java Blend. It does not provide an independent object query language, but enables programmers to query, update or delete objects using the language binding of their choice. Although this allows developers to access the database without switching to a query language, it deprives them of the expressiveness and the flexibility of OQL. Query processing in ObjectStore is studied in [OHMS92].

Objectivity/DB also includes C++ and Java bindings that conform to the second release of the ODMG standard (ODMG 2.0). In addition, it provides a separate query facility called SQL++. This is an ANSI-standard SQL engine extended with object-oriented features of Objectivity/DB. The design of SQL++ was intended to preserve the investment in traditional SQL-based tools; the effect of this was a loss of expressive power compared to OQL, such as the failure to handle multiple collections, late binding and class casting.

Versant Developer Suite (VDS) 6.0 is an object database management system that provides a tight binding between the application programming language (C++ or Java) and the underlying object database model. Versant has also its own proprietary Versant Query Language (VQL). VQL does not only have select query capabilities, but also supports insert, update and delete operations. Queries have an SQL-like syntax extended to include a number of object-oriented concepts. However, VQL does not have the expressive power of OQL.

A more detailed description is given for POET, since it has been used as the basis for developing the query framework of this research. POET has actually contributed in defining the ODMG standard. It provides two language bindings (for Java and C++), as well as an OQL query processor. However, only a subset of OQL constructs is implemented; neither method invocation nor the group-by operation are supported in queries. The type system of the ODMG object model is not entirely respected. For instance, the result of a query in POET is either of type `Object` or of type `Set(Object)`. According to the standard, it could be any simple type, class or collection like bag, list or array. The weaknesses of the POET OQL implementation motivated in part the development of a new OQL engine. Since in most other respects – object model, language bindings – POET conforms to the ODMG standard, its storage system and Java binding were selected as tools for building and testing the new query processor.

1.4 Type Checking vs Type Inference

The ODMG 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 [Ala99, BT00]. Their common objective is to ensure that the query execution will not fail due to type errors. This is also reflected by the following statement in the ODMG standard:

”OQL is a functional language where operators can freely be composed, as long as the operands respect the type system” ([CBB⁺00], page 89).

The core of the type checking model [BT00] is a set of typing rules which derive the type of OQL constructs based on the types of their constituent parts (operands or subqueries). If a subquery is just an identifier its type is looked up either in the query environment, if it is a local range variable, or in the database schema, if it is a named object, a query definition or an extent. Hence, knowledge of the schema type information of a database is essential for type checking.

On the other hand, type inference derives the most general type of a query without accessing any schema information. Whilst doing that, it establishes restrictions on the properties of a schema, in order that a certain query may be well-typed with respect to it.

This technique could be exploited in distributed database applications where we often need to check the typability of a query against multiple schemata. It would also be useful in order to transform queries into equivalent forms that fit different schema specifications. Finally, a type inference model could be exploited in contexts where we have a limited knowledge of schema information (semi-structured or unstructured databases).

OQL type inference is particularly interesting, because of the complexity of the type model of the language. Besides class types and simple types, the type model includes structures, as well as multiple collection types (sets, bags, lists, arrays and dictionaries). The wide variety of OQL constructs leads to multiple combinations of these types, and thus to a rather complex type system. Type inference is studied in detail as a part of this thesis.

1.5 Calculus and Algebraic Representations

Most query languages, including OQL, are declarative, i.e. they declare what has to be provided as the result, without describing how this is going to be achieved. Intermediate representations aid in bridging the gap between the user-level language and the query engine programs ultimately used to produce the query result. There are calculus-style and algebra-style intermediate representations that serve different purposes, but both need to reflect the special characteristics of the query language. Calculus-style formalisms support the concept of tuple (for relational models) or object (for

object models) variables, while algebra-style formalisms are set- (resp. collection-) oriented. The former allow for nesting of expressions in a natural way; the latter use operators that reflect mutually independent (context free), but well defined execution steps [SdBB96].

In this thesis, the monoid comprehension calculus is adopted as the initial internal representation of OQL. This calculus provides a uniform way of expressing a great variety of queries; for instance, operations over multiple collection types, group-by operations, existential or universal quantification have very similar calculus representations. The calculus framework of this dissertation has extended previous research [Afs98, FM95, FM00], in order to deal with several weak points of OQL and to resolve a number of ambiguities. One of these extensions concerns the introduction of new normalization rules wrt [FM00]. Normalization is the transformation (unnesting) of monoid expressions to syntactically equivalent ones which are either more efficient, or provide opportunities for further calculus or algebraic optimization.

A set of algebraic operators (e.g. Reduce, Select, Nest) are proposed as the second internal representation of OQL. In order to translate a query from the calculus to its algebraic form, one needs to apply a series of translation rules. Having obtained the algebraic form of a query the next step is to provide transformation rules that convert an algebraic form to a more efficient form. Note that the term *translation* is used to express the conversion from one internal representation to another; the term *transformation* denotes the conversion of an expression to an optimized one within the same internal representation [SdBB96]. The algebraic operators presented in this thesis are similar to those proposed by Fegaras and Maier [FM00]; however, a great deal of the semantics and the applied transformations are different in the two contexts. In addition, the respective algorithms for the translation of calculus to algebraic forms differ to a great extent.

1.6 Semantic optimization using association rules

In general, optimization methods transform a query from one representation to a more efficient one, or to a form that provides opportunities for further optimization. Calculus normalization and algebraic optimization, which have traditionally dominated the query evaluation process, involve transformations that are valid for any database state. In contrast, semantic transformations are valid for any state satisfying a set of constraints (semantic rules) [SSS92]. These rules differ from database integrity constraints, in that they are not strictly imposed, that is they do not reject any update that invalidates them. Such an update would rather cause the relaxation or the deletion of the rules it invalidates.

Association rules reflect correlations between attribute values in a database. They have mainly been studied in the context of sales transactions databases. A lot of work has focused on the market-basket problem, i.e. on finding association rules of the form *if a set of items appears in a transaction, then another different set of items is likely to appear as well*. These *boolean* association rules have

been extended to deal with quantitative and categorical values.

Two measures – support and confidence – express the validity of association rules. Support is the fraction of data instances satisfying a certain rule, while confidence is the ratio of instances satisfying the rule to the instances satisfying its antecedent. Semantic rules can be viewed as association rules with 100% confidence. Semantic optimization could be achieved using association rules, as long as their exceptions are taken into account in the query transformations. We assume that these exceptions are maintained in the presence of updates [AFP00, CLK97]. Considering the exceptions of association rules results in valid but not necessarily advantageous semantic transformations.

Siegel has proposed a series of heuristics for the semantic optimization of queries [SSS92]. The following two are studied in detail in this thesis:

- if a selection on attribute A is implied by another selection condition on attribute B and A is not an index attribute, then the selection on A can be removed from the query;
- if a relation R in the query has a restricted attribute A and an unrestricted cluster index attribute B, then look for a rule where the restriction on A implies a restriction on B.

A number of algorithms proposed in this dissertation enable the application of both heuristics for the optimization of OQL queries. In addition, an experimental framework has been set up, in order to measure the complexity of the optimization algorithms and the performance benefits of the query transformations.

1.7 Dissertation outline

The rest of this thesis is structured as follows: Chapter 2 presents first a series of applications for which this research is particularly useful. These applications are complex enough to benefit from an object-oriented paradigm and are rich in meaningful correlations in order to provide the opportunity for semantic optimization. Further, this chapter reviews previous work related to the various phases of query processing and optimization. In doing so, it points out the main differences between previous work and the research presented in this thesis.

Chapter 3 focuses on the first step of query processing, that is the type inference of OQL queries. Note that there is no discussion about checking the type of an OQL query with respect to a certain schema. The type-checking algorithm is quite similar to the ones proposed for relational query languages and is formally presented in [BT00]. Instead this chapter focuses on the following problem: given an OQL query find its most general type as well as the type requirements that a schema should satisfy in order to be compatible with the query.

Chapter 4 discusses the translation of an OQL query to its first internal representation, the monoid comprehension calculus. We argue against the claim of Fegaras and Maier that nearly all OQL

expressions have a direct translation into their monoid calculus, with the exception of indexed OQL collections ([FM00],section 3). We extend the calculus syntax and modify the translation rules, in order to provide a full mapping of OQL queries to their calculus forms. Further, we present additional rules for the normalization of monoid calculus expressions taking into account the new calculus expressions.

Chapter 5 is the equivalent of chapter 4 on the algebraic level. That is, it first gives an overview of the algebraic operators used at this level and explains their functionality; it then presents the rules for the translation of the calculus representation to the algebraic one. The reader will notice the differences between the algebraic representation of this thesis and the one proposed in [FM00].

Chapter 6 introduces the interesting topic of semantic optimization. It gives a series of algorithms for applying two optimization heuristics proposed by Siegel et al. [SSS92]. It then discusses the contexts in which each of these techniques is expected to speed up query execution.

Chapter 7 evaluates the complexity of the semantic optimization techniques and assesses their benefits in query execution. These issues are first addressed theoretically, by reasoning about the cost of each step of the optimization algorithms. Complexity is then evaluated practically, by setting up a series of experiments that monitor the behaviour of the algorithms under different input conditions.

Chapter 8 discusses the insights and the experience gained from this research. It presents some weaknesses and suggests directions for future work. Finally, it concludes this dissertation by pointing out its principal contributions.

Chapter 2

Applications and Related Work

This chapter discusses two important issues relevant to this work. First, it presents application frameworks that can benefit from the algorithms developed in this research. These applications are expected to be rich in semantic associations and complex enough to benefit from their deployment in an object-oriented framework. Second, this chapter reviews previous research done on the various stages of query execution and optimization. It focuses mostly on the work that has directly influenced this research, and highlights the novel points of this dissertation.

2.1 Applications that motivate this research

The applications that are particularly interesting in our context are often labelled as ‘data-rich’. This term has multiple interpretations: first, it requires that an application has large amounts of data that need to be managed in an intelligent way. Secondly, it implies that this data is inherently complex, and hence it would benefit from being structured in an object-oriented database. Thirdly, this term suggests that the data is rich in semantic correlations; the latter could be used either for query optimization or for decision-support purposes.

These applications can be found in both commercial and scientific contexts. Starting with the commercial world, most organisations are collecting increasing volumes of data about business processes and resources. Traditionally, they organise this data in relational databases, representing interconnections between tables by means of foreign keys. They query the information in order to produce long reports referring to aggregates of costs, earnings and various other statistical information. More complicated queries, aiming at risk analysis or prediction of future trends, are usually performed in a semi-manual way. Due to the volume of data, queries that range over entire tables (or extents) take a long time to evaluate; they are typically included in batch jobs and performed overnight.

For instance, consider a banking application that keeps data about all customers that have had a

loan contract with a bank. It would be very convenient if the bank information was organized in an object-oriented database, since one would be able to navigate easily from customers to their loans and deposits, their repayments and other obligations, as well as other customers they are related to, e.g. guarantors or partners in the same company, or suppliers. Apart from inverse relationships and path navigation, the banking application offers a lot of opportunity for inheritance, since products usually form hierarchies and so do customers or contracts. Further, this application provides many opportunities for identifying useful semantic knowledge. For example, in order to estimate whether a current customer is likely to repay his/her loan, it is sensible to look at the behaviour of previous customers with similar characteristics. This entails queries that range over large amounts of data relating the propensity of arrears of customers based on their loan application data. Identifying common patterns of customer behaviour, in the form of association rules, would be of great help to the execution of similar queries.

Moving to the scientific world, medical applications are representative examples that could benefit from the proposed methods. Given a database of patient histories (symptoms, diagnoses and treatments), interesting rules can be deduced from the association of exhibited symptoms and prescribed treatments with the evolution of a patient's condition. In order to reach slightly refined or more general conclusions, a doctor might want to address queries that are likely to be optimized by the application of the derived rules. It is worth representing the medical data in an object-oriented way for two reasons: first, symptoms, diseases and treatments have many attributes and two-way relationships with each other; doctors are particularly interested in these relationships and often write queries that navigate across them. In relational databases, this would require many expensive join operations, whilst in their object counterparts, it would involve pointer traversals. Second, the classification of symptoms or diseases under more general categories (using the inheritance mechanism of the object model) might enable high level correlations among medical data to be identified. This would enable the optimization of a wider variety of queries.

In general, an optimization is considered beneficial when the time saved by executing the optimized query is greater than the cost of applying the optimization. However, in applications where queries are executed frequently, it is possible to apply the optimization once at compile time and execute the optimized query several times without any further overhead. The semantic optimization methods proposed in this dissertation obey this general rule, thus favouring applications that apply a particular query repeatedly. An important consideration is the dynamics of the database; we should be able to monitor the validity of prespecified semantic transformations in the presence of database updates.

2.2 Previous work related to this research

As shown in the previous section, there are large families of commercial (e.g. banking) and scientific (e.g. medical) applications suitable for applying semantic optimization techniques. This fact has motivated a lot of research in the area of query processing and optimization that has been used as the basis of the current work. The remainder of this chapter outlines previous research efforts from the areas involved in this work: type theory, intermediate query representation, data mining and semantic optimization.

2.2.1 Ensuring the typability of an OQL query - Related Work

The first step towards processing an OQL query is checking its typability, i.e. ensuring that it can be executed without causing any compile-time type errors. This can be done either with respect to a specific schema (type checking) or just based on the structure of the query independent of any schema information (type inference). Both problems have been studied broadly for various languages during the last 30 years [BO96, Dam85, Gun92, Mil78, Sul00, Tof87]. Typability issues have been considered for query languages in both relational and object oriented frameworks.

A discussion about type-checking an OQL query can be found in the ODMG standard [CBB⁺00]. Unfortunately, this presentation of the type system is neither formal nor complete. There is obviously a need for a more formal, mathematical presentation; the informal English presentation may be eloquent but it leaves a lot of space for ambiguity and misinterpretation.

Riedel and Scholl set out to describe the type system underlying the ODMG object model formally [RS97]. However, their approach, reflected in the CROQUE-version of ODL, differs considerably from the type model as defined in the ODMG standard [CBB⁺00]. They propose a type-based object model, where object types need not be named, since they are distinguished uniquely by their set of functions (attributes, relationships and methods). Classes are not considered as the types of object instances, but rather as extents. Hence, they are treated in queries differently from the way that they are treated in programming languages. Although Riedel and Scholl define types and subtyping formally, they provide no formal means of checking the type of an OQL query, e.g. a complete set of type rules. Finally, they study a much earlier version of the Standard (version 1.2) so some of their work is no longer relevant.

Following this work, Alagić [Ala99] and Bierman and Trigoni [BT00] made two independent efforts towards checking the type of an OQL query. Alagić considers a formalisation of the OQL type system, and claims that 'several negative results are proved about the ability to type-check queries'. This claim is later questioned by the work of Bierman and Trigoni who give a set of typing rules adequate for checking the type of all OQL constructs. Alagić also demonstrates a much more serious problem with the Java language binding. He gives a type rule [[Ala99], Rule 26], which appears to

be a rule for typing Java directly, i.e. he considers OQL as a language extension of Java. Clearly there is a problem here as OQL has parameterised types (e.g. `bag(Employee)`), but Java does not (all it has is covariant arrays).

Using the host language type system to type queries has been the dream of OODBMS designers [ZM90]. However, nowhere in the Standard does the ODMG propose that the Java type system be used to type OQL. The type system of the Object Model should be the basis for OQL type checking. The Standard gives a number of language bindings that enable the user to invoke OQL queries within a programming language environment. It is the role of these bindings to provide the mapping between the programming language type system and that of the Object Model.

Bierman and Trigoni address the typing problem from a lower level [BT00]. They give a detailed account of the typing rules needed to check the type of all OQL constructs, and not only of the select queries (as is the case in Alagić's work). Hence they provide rules that deal with grouping and various collection operations (e.g. `union`, `intersect`, `except`).

Regarding the order operation, the Standard does not specify what happens when the sort criterion in the order clause does not distinguish completely between results. Presumably the order of elements in the resulting list is system-dependent. To alleviate this problem, Alagić enforces the limitation that the ordered elements should implement an interface with an ordering method. Bierman and Trigoni suggest that the resultant type of a `select group by` clause should be a list of bags of elements rather than a list of elements. Each bag should include all the elements with the same values for the sort criteria. If the sort criteria define an absolute order for the elements, then all the bags are singletons. An additional `order_flatten` operator could be potentially added as part of the OQL language in order to flatten the list of bags of elements into a list of elements in a system-specific way.

As far as type inference is concerned, the canonical reference for work on type systems for database programming languages is the work of Buneman and Ohori [BO96]. The goals of their type inference algorithm are identical to those when inferring the type of OQL queries, as presented in the next chapter; 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 to a considerable extent. Buneman and Ohori introduce kinded types to infer the type of a record based on selections of fields on this record. Instead the type inference work in this dissertation uses the notion of *general* types; in this way, it is possible to express general type information not only for records, but also for parametric collection types, dictionary types, structures and classes. Moreover, a wider variety of constraints are introduced and hence a different algorithm is proposed to resolve them.

The work most related to the object nature of OQL arises from studying type systems for object

oriented programming languages, e.g. [PS91, PC94, AH95]. Palsberg and Schwartzbach [PS91] present an algorithm that infers the type of an untyped object-oriented language resembling Smalltalk. They define *type* as a finite set of classes, and *induced type* as the set of classes of all possible non-nil values to which a concrete program may evaluate in any of its executions. There are finitely many classes in a program and a type of an expression refers to the set of possible classes the expression is an instance of. However, inferring the type of an OQL query without accessing a specific schema raises the need to deal with *class templates* rather than with specific classes. Moreover, the OQL type system includes several atomic, collection and structure types. Thus, a general type cannot be defined as a set of classes.

Gunter [Gun92] studies polymorphism and type inference for a calculus language called ML₀. He presents an algorithm *W* that infers the most general type of an ML₀ expression. The main point where his approach differs from that of the current work is the unification mechanism. He uses substitution as the mechanism to unify two types that should be equal. However, in the OQL setting, which supports subtyping, the need to unify two types is not simply equivalent to the need for their equality. In certain cases, it concerns comparing two types using subtype relations or their ability to be element types in the same collection. Thus, the notion of substitution is not sufficient in resolving constraints amongst OQL types.

The language Cardelli [Car87] considers for type inference is also a very simple lambda calculus with constants. He presents inference rules that are applied from left to right and bottom-up. Starting with initially broad assumptions, he gradually refines them using inference rules. In case an inference rule requires two types to be equal, he unifies them and in this way he specializes some of the types of the identifiers. The approach used in chapter 3 is very similar, except that the inference rules generate two categories of constraint. The first consists of constraints that relate the types of two different expressions that coexist in the same *context* (query construct). The second concerns constraints that relate different information about the type of the same expression. This information might have been inferred by analyzing different *contexts* (parts of the query).

2.2.2 Intermediate representations of OQL - Related Work

The next step after ensuring the typability of an OQL query is to transform it into its internal representations, i.e. the calculus and the algebraic form.

Fegaras and Maier present a standard of effectiveness for a database calculus, before proposing monoid calculus as an intermediate representation for OQL queries [Feg94, Feg97, Feg98, FM95, FM00]. They prove that monoid comprehensions conform to this standard by being expressive enough to represent OQL's constructs, manipulable (w.r.t. type checking, rewriting et cetera), evaluable (giving opportunities for alternative execution plans) and uniform in handling collection and primitive types (as monoids). Provided that certain types are viewed as monoids, a query can also

be considered as a map from a given monoid to another monoid. These maps are known as monoid homomorphisms. Thus, a framework is proposed in which OQL queries are translated to monoid comprehensions, which in turn are defined in terms of monoid homomorphisms. The calculus representation is eventually translated into a series of algebraic operators, like `join`, `select`, `unnest` and `reduce`.

Grust and Scholl also support the concept of monoid comprehension as a calculus representation [GS98a]. Like Fegaras and Maier, they also underline that monoid calculus expressions are amenable to normalization. The existence of normal forms constitutes a major advantage for two reasons. First it enables *pattern-matching*; secondly, it reduces the number of transformation rules needed to deal with the repertoire of OQL operators. In spite of that, normalisation should be treated carefully, because it can sometimes hide optimization opportunities [GS96]. Grust and Scholl map monoid comprehensions to a combinator algebra (including `map`, `filter`, `cross`, `join`, `semijoin`, `antijoin` and aggregation operators). The hybrid representation of monoid comprehensions and algebraic combinators is treated as pure syntactic sugar for a single basic recursion combinator `foldr`. The latter is used in combination with `build` in order to enable the derivation of streaming programs (cheap deforestation) [GS98b].

A hybrid representation of monoid comprehensions and a monoid-aware nested algebra (e.g. `semijoins`, `nestjoins`) is discussed in [GKG⁺97]. The core of the adopted algebra is the operator `fold[\oplus ; f](E)`, which applies `f` to the elements of `E`. It then uses the merge operation \oplus of some monoid `M` to collect the result. The algebraic operator `fold` along with certain variations like `crossfold` (for the evaluation of joins) are also supported in [Van92].

Buneman et al. propose a comprehension syntax as the starting point for the development of database languages [BLS⁺94]. They claim that comprehensions provide a natural way of expressing structural recursion and a powerful paradigm for manipulating collection types.

Straube et al. also present a framework for the translation of an object query language to calculus and algebraic representations [SO95a]. A query in their object calculus is of the form $\{o \mid \psi(o)\}$ where `o` is an object variable and $\psi(o)$ is a formula built from range or predicate atoms. They discuss safety criteria for the computability of queries and they give an algorithm for the translation of safe calculus expressions to algebra operator trees. Finally, they present equivalence-preserving rewrite rules in order to optimize algebraic expressions. The expressiveness and the type system of the query language considered in their context are less complex than those of OQL.

A sensible question that might be raised at this point is the following. Why are the algebraic representations of extended relational and object-oriented query languages (`fold`, algebraic combinators and `foldr`, monoid homomorphisms) not as uniform as those concerning the respective calculus representation (monoid comprehensions)? It is tempting to highlight that the main difference of the algebraic approaches lies in the concept of structural recursion. On the one hand, `foldr` directly

corresponds to structural recursion on the insert representation (sri). It is, therefore, naturally associated with sequential processing of collections. On the other hand, `fold` and monoid homomorphisms express structural recursion on the union representation (sru) and are consequently associated with the parallel processing of collections.

Opting for translating comprehensions to monoid homomorphisms has the advantage that monoid homomorphisms have the potential for parallel processing via divide-and-conquer methods [Afs98, HIT97]. This option has already been proved successful for relational query languages in [Afs98]. The query evaluation framework - developed as a part of the architecture - makes use of monoid comprehensions which are first translated to monoid homomorphisms and then (after several rewriting steps) to algorithmic skeletons, e.g. `hom`, `split`, `pipe` and `farm`. The latter are constructs with well-defined parallel behaviour and, in particular, the `hom` skeleton models the behaviour of the monoid homomorphism and exhibits its potential for parallelism.

However, the use of monoid homomorphisms would prevent us from efficiently evaluating execution plans using valuable experience from the relational technology. Reordering join operations, selecting different join algorithms, applying various (index-based) access methods and achieving semantic optimization are some of the optimizations that are hard to apply using the homomorphic representation.

Hence, queries in this dissertation are first translated to the monoid calculus and then to a set of algebraic operators that also provide the potential for parallel processing. Although this approach is similar to the one taken by Fegaras and Maier, it departs from theirs in several respects.

First, the normalization rules within the calculus representation have been extended. Secondly, the semantics of the algebraic operators have been altered to better suit the needs of query representation and optimization. Thirdly, the rules for the translation of the calculus to the algebraic formalism differ to a considerable extent. The calculus and the algebraic representations as well as the transformations performed within each level are described in chapters 4 and 5.

2.2.3 Association Rules and Semantic Optimization - Related Work

Within the intermediate query representations, there are a lot of opportunities for syntactic calculus or algebraic transformations. A lot of experience has been gained applying these optimizations in the relational or object-oriented context [CS96, Cha98, GLR97, Hel98, Ioa96, KPH98, PS96, PS97, RS93, SdBB96, SO95a, SO95b, WM99]. This thesis focuses on semantic query optimization, i.e. query transformations based on semantic knowledge rather than syntactic equivalence. In particular, semantic knowledge is represented by association rules of the form:

$$p_1 \wedge p_2 \wedge \dots \wedge p_n \Rightarrow p_0 \text{ (Exc)}$$

where p_i , $i = 0, \dots, n$ have the form `class_name.attr_name op value` and $op \in \{=, <, \leq, >, \geq\}$. `Exc` is the set of exceptions associated with the rule; it includes all the objects of class `class_name` that satisfy the antecedent $p_1 \wedge \dots \wedge p_n$, without satisfying the consequent p_0 . Such an association rule is often characterized by two measures of interest, confidence and support; we redefine these measures to fit the object-oriented paradigm. Confidence is the fraction of objects of class `class_name` satisfying the antecedent, that also satisfy the consequent. Support is the percentage of objects in the extent of `class_name` satisfying both the antecedent and the consequent.

A lot of research has been devoted to mining association rules that relate values of attributes in a relational table. The same techniques can equally be used in an object-oriented framework. Agrawal et al [AIS93] present an algorithm for identifying boolean association rules, i.e. correlations between boolean attributes. Further, Srikant and Agrawal presented techniques for discovering rules over quantitative and categorical attributes [SA96]. Later on, Miller and Yang studied the problem of mining association rules over interval data [MY97] using the Birch clustering algorithm [ZRL96]. A lot of attempts have been made to improve or to alter slightly the basic algorithm underlying the approaches discussed above [ASY98, CA97, FMMT96, KMR⁺94, NLH98, Par95, TUA⁺98, WYY00]. Chapter 7 presents an experimental framework intended to demonstrate the benefits of semantic optimization using association rules. The programs developed to identify these rules are based on standard algorithms for identifying quantitative rules [SA96] and rules over interval data [MY97].

Semantic knowledge that has been discovered using data mining techniques can be used for query optimization [HHCF96, Hsu96, HK96, Hsu97, HK98, HK00, Kin81b, LS95, MC96, SHKC93, YSP97, YS89]. A number of heuristics have been proposed that take advantage of association rules in order to optimize queries semantically [Bel96, CGM90, GGMR97, GGZ01, Kin81a, SSS92]. We refer briefly to the heuristics presented by Siegel et al [SSS92].

- *If a relation R in the query has a restricted attribute A and an unrestricted cluster index attribute B , then look for a rule where the restriction on A implies a restriction on B . This heuristic implies the introduction of the second constraint, which is indexed.*
- *If the only restriction on relation R is a join condition and R does not have an index on the join attributes, then look for a selection restriction on this relation. The new selection condition should be deduced from the selection conditions on the joining relation.*
- *If a selection condition on attribute A is implied by another selection condition on attribute B and A is not an index attribute, then the selection on A can be removed from the query.*
- *If a query contains a restriction on a range attribute, then try to locate a constraint on that attribute that is a contradiction to the existing restriction.*

- *If a restriction on a range attribute implies the answer to a query then the latter can be answered without accessing the database*

All of these heuristics make use of association rules to add or remove constraints, or to find contradictions or tautologies. In this dissertation we focus on the application of the first and third conditions, although the algorithms proposed can be adapted to serve all of the heuristics.

Chakravarthy et al. [CGM90, CGM88] also study the use of semantic knowledge for the transformation of queries into more efficient forms. They use a set of heuristics similar to those described above and propose techniques that take advantage of integrity constraints in the form of rules. Bell [Bel96] has studied the use of constraints in order to identify whether the result of a query contains distinct values. This knowledge is used in order to optimize `select · distinct . . .` queries or `group by` operations.

Beneventano et al. [BBSV97] have developed a semantic query optimizer as part of a broader set of object-oriented tools (ODB-tools). Their optimizer takes advantage of the constraint introduction and constraint elimination heuristics, as well as of additional ones, e.g. reducing the scope of a query to a smaller subclass. These heuristics are presented in detail by Grant et al. [GGMR97].

Cheng et al. describe the implementation of two semantic optimization techniques, namely predicate introduction and join elimination, in DB2 Universal Database [CGK⁺99, GLQ99]. After an input query is converted to an intermediate form called *query graph model*, it is transformed by a query rewrite engine [PHH92, PLH97] into a more efficient form using a set of integrity constraints.

Most of the approaches above ignore weak association rules, i.e. rules that are not a hundred per cent valid. In most cases rules represent integrity constraints, allowing no updates that invalidate them. Godfrey et al [GGZ01] consider the use of *absolute soft constraints* the exceptions of which are stored in materialized views. However, they do not consider combining soft constraints (and their exceptions) to infer new constraints. In general, previous approaches take advantage only of direct correlations as expressed by single rules. They do not combine rules to identify indirect associations that might be of use to the semantic optimization process.

In this dissertation, a series of algorithms is given that deals with these two issues. First, these algorithms take into account the exceptions associated with each rule, so that the result of the optimized query is accurate. Secondly, they combine existing association rules to identify useful indirect correlations along with their exceptions. Both direct and indirect correlations are combined to a series of optimization solutions; the criteria according to which the best solution is selected depend on the optimization heuristic.

2.3 Extending the previous work - Discussion

This chapter provided an overview of the applications and the previous work relevant to this research. The variety of applications is definitely a strong motivation for extending the existing work to address the problem of query processing and semantic optimization in an object-oriented framework. Although a lot of research has already been done in these areas, little work has focused on integrating the various stages of query processing, namely type inference, translation to calculus/algebraic representation and semantic optimization using association rules. Very few attempts have actually tried to measure experimentally the benefits of semantic optimization with respect to the optimization cost. To our knowledge there has been no implementation of the semantic optimization techniques that take account of indirect association rules with exceptions.

Chapter 3

Type Inference of OQL

This chapter presents 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. It is the result of joint work with Dr G.M. Bierman, presented in the proceedings of BNCOD 2001 [TB01]. The proposed 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. The 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. The resulting type inference system for OQL should be useful in any open distributed, or even semi-structured, database environment.

3.1 Introduction

The ODMG Standard [CBB⁺00] 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 [Ala99, BT00]. The work presented in this chapter builds on earlier work [BT00] 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)
{ attribute string      name;
  attribute string      position;
  attribute int          year_of_birth;
  attribute float        salary;
  attribute Department  department;}

class Department (extent Departments)
{ attribute string id;}

```

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 can we derive about the type of the result of the query, supposing that we have no specific schema in mind? Before addressing these questions in detail, we consider the setting where type inference could be important.

As one example, this information could be exploited in distributed database applications. Suppose time-critical queries are addressed against multiple schemata. If there are frequent updates on parts of these schemata, then many of the queries will inevitably fail to be executed. In order to avoid this situation, a user can register interest in specific updates of each schema—at least in those that would affect the critical queries—and resolve the type incompatibility at the time the queries get executed.

Type inference 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 [BHP94, Hul97]. 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 an 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.

Finally, the work discussed in this chapter enables the use of *general* query definitions. Version 2.0 of the ODMG standard allowed the user to write query definitions without explicitly defining the types of their arguments. This feature was changed in version 3.0, presumably either to simplify the task of type-checking or to complete the latter as much as possible at compile time. However, the ability to define general queries, which possibly run against several schemata without any change, is lost. Our inference algorithm enables the use of query definitions with unspecified argument types. It can infer the most general types of these arguments in exactly the same way as it derives the schema requirements (class/extent information, named objects or query definitions).

3.2 Core OQL

As mentioned previously, the work in type inference is based on earlier work in checking the type of an OQL query wrt specific schema information [BT00]. This section gives an overview of the grammar, the types and the subtyping relation of a core OQL, i.e. a fragment of the language defined in the Standard, but which has the same expressive power [BT00].

An OQL **program** consists of a number (maybe zero) of named definitions followed by a query. The syntax for queries and definitions¹, as well as the OQL types are presented in figure 3.1. In what follows, $\text{Col}(\sigma)$ denotes 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.

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

Queries $q ::= b \mid f \mid i \mid c \mid s$
 $\mid x$
 $\mid \text{bag}(q, \dots, q) \mid \text{set}(q, \dots, q) \mid \text{list}(q, \dots, q) \mid \text{array}(q, \dots, q)$
 $\mid \text{struct}(\ell: q, \dots, \ell: q)$
 $\mid C(\ell: q, \dots, \ell: q) \mid q.\ell \mid (C)q$
 $\mid q[q] \mid q \text{ in } q \mid q() \mid q(q, \dots, q)$
 $\mid \text{forall } x \text{ in } q: q \mid \text{exists } x \text{ in } q: q$
 $\mid q \text{ binop } q \mid \text{unop}(q)$
 $\mid \text{select } [\text{distinct}] q$
 $\text{from } (q \text{ as } x, \dots, q \text{ as } x)$
 $\text{where } q$
 $[\text{group by } (\ell: q, \dots, \ell: q)]$
 $[\text{having } q]$
 $[\text{order by } (q \text{ asc} \mid \text{desc}, \dots, q \text{ asc} \mid \text{desc})]$

Definitions $d ::= \text{define } x \text{ as } q$
 $\mid \text{define } x(x, \dots, x) \text{ as } q$

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 $\sigma ::= \text{int} \mid \text{float} \mid \text{bool} \mid \text{char} \mid \text{string} \mid \text{void}$
 $\mid \sigma \times \dots \times \sigma \rightarrow \sigma$
 $\mid \text{bag}(\sigma) \mid \text{set}(\sigma) \mid \text{list}(\sigma) \mid \text{array}(\sigma)$
 $\mid \text{struct}(\ell: \sigma, \dots, \ell: \sigma)$
 $\mid C$

We assume a distinguished class name `Object`.

Figure 3.1: Syntax and Types for Core OQL

Sub-typing

$$\begin{array}{c}
\frac{}{\mathbf{C} \leq \mathbf{Object}} \text{Top} \quad \frac{\mathbf{C} \sqsubseteq \mathbf{C}'}{\mathbf{C} \leq \mathbf{C}'} \text{Sub-Class} \\
\\
\frac{\sigma'_1 \leq \sigma_1 \cdots \sigma'_k \leq \sigma_k \quad \tau \leq \tau'}{\sigma_1 \times \cdots \times \sigma_k \rightarrow \tau \leq \sigma'_1 \times \cdots \times \sigma'_k \rightarrow \tau'} \text{Sub-Fun} \quad \frac{\sigma \leq \tau}{\text{Col}(\sigma) \leq \text{Col}(\tau)} \text{Sub-Coll} \\
\\
\frac{\sigma_1 \leq \tau_1 \quad \cdots \quad \sigma_k \leq \tau_k}{\mathbf{struct}(\ell_1: \sigma_1, \dots, \ell_k: \sigma_k, \dots, \ell_{k+n}: \sigma_{k+n}) \leq \mathbf{struct}(\ell_1: \tau_1, \dots, \ell_k: \tau_k)} \text{Sub-Struct} \\
\\
\frac{}{\sigma \leq \sigma} \text{Sub-Refl} \quad \frac{\sigma \leq \sigma' \quad \sigma' \leq \sigma''}{\sigma \leq \sigma''} \text{Sub-Trans}
\end{array}$$

Figure 3.2: Subtyping for core OQL

This we shall write $\sigma \leq \tau$ and define as the least relation closed under the rules given in Figure 3.2.

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

An interesting feature of our subtype relation is the treatment of structures. A structure type $\sigma = \mathbf{struct}(\ell_1: \sigma_1, \dots, \ell_m: \sigma_m)$ is considered to be a subtype of $\tau = \mathbf{struct}(\ell_1: \tau_1, \dots, \ell_n: \tau_n)$ 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 [BT00]. 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 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 context of inferring the type of an OQL query, we have no information about the classes, the query definitions, and the named objects, as we have no schema information. The problem addressed

in this chapter 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?

3.3 Type and Schema Inference

This section presents 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.

3.3.1 Extended Type System

It turns out to be convenient to extend the notion of type given in Figure 3.1; 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 3.1) with the following so-called *general* types.

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

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

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

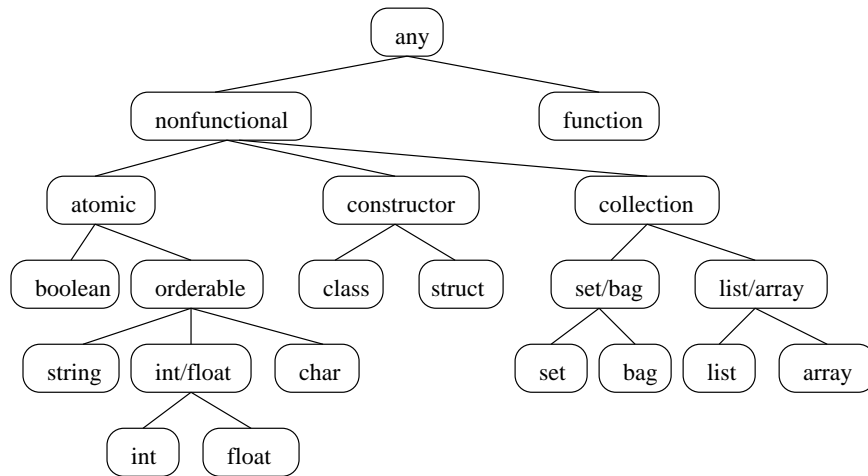


Figure 3.3: Type Hierarchy

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

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

This extended type system is coupled with the type hierarchy illustrated in Figure 3.3. 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.

A type is said to be *specific*, if it can be derived by the type system given in Figure 3.1. Otherwise, it is said to be *general*. All the leaves of the hierarchy tree (in Figure 3.3) 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 .

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

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

coll_i , $i = 1, 2$ are nodes in the subtree (figure 3.3) 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.

$$\frac{\text{constr}_1 \text{ is child of constr}_2 \quad \sigma_1 \subseteq \tau_1 \cdots \sigma_k \subseteq \tau_k}{\text{constr}_1(\ell_1: \sigma_1, \dots, \ell_k: \sigma_k, \dots, \ell_{k+n}: \sigma_{k+n}) \subseteq \text{constr}_2(\ell_1: \tau_1, \dots, \ell_k: \tau_k)} \text{GSR - Constr}$$

constr_i , $i = 1, 2$ are nodes in the subtree (figure 3.3) 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.

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

atom_i , $i = 1, 2$ are nodes in the subtree (figure 3.3) 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 constr_1 , constr_2 are nodes in the subtree (figure 3.3) with root constructor and that constr_1 is a child of constr_2 or $\text{constr}_1 = \text{constr}_2$. coll_1 and coll_2 are nodes in the subtree (figure 3.3) with root collection and coll_1 is a child of coll_2 or $\text{coll}_1 = \text{coll}_2$. Although LUB is always defined, this is not the case for GLB. There is not always a type that is more specific than each possible pair of OQL types. For instance, it is not possible to define the GLB of a class and an atomic type.

$$\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 \\
& \quad (\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(\ell_1: \sigma_1, \dots, \ell_k: \sigma_k), \text{constr}_2(\ell_1: \sigma'_1, \dots, \ell_k: \sigma'_k, \dots, \ell_{k+n}: \sigma'_{k+n})) \\
& \quad \stackrel{\text{def}}{=} \text{constr}_1(\ell_1: \text{GLB}(\sigma_1, \sigma'_1), \dots, \ell_k: \text{GLB}(\sigma_k, \sigma'_k), \dots, \ell_{k+n}: \sigma'_{k+n}) \\
& \text{LUB}(\text{constr}_1(\ell_1: \sigma_1, \dots, \ell_k: \sigma_k), \text{constr}_2(\ell_1: \sigma'_1, \dots, \ell_k: \sigma'_k, \dots, \ell_{k+n}: \sigma'_{k+n})) \\
& \quad \stackrel{\text{def}}{=} \text{constr}_2(\ell_1: \text{LUB}(\sigma_1, \sigma'_1), \dots, \ell_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}}{=} \\
& \quad \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}}{=} \\
& \quad \text{GLB}(\tau_1, \tau'_1) \times \dots \times \text{GLB}(\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 \\
& \quad \tau_1, \tau_2 \in \{\text{atomic, constructor}(), \text{collection}(\text{any})\}
\end{aligned}$$

3.4 Type compatibility - Constraints

The inference algorithm we present in section 3.7 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 ($\tau = \text{Merge_Result}(\tau_1, \tau_2)$). In section 3.5, 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 should 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 can be compared for inequality (`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 3.6, 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</code> (τ_1, \dots, τ_n)	6. $\tau_0 = \text{Merge_Result}(\tau_1, \tau_2)$
2. <code>Greater_Less_Than_Compatible</code> (τ_1, \dots, τ_n)	7. $\tau_0 = \text{Distinct_Result}(\tau_1)$
3. $\tau = \text{Member_Type}(\sigma)$	8. $\tau_1 = \tau_2$
4. $\tau = \text{Constructor_Member_Type}(\sigma, \ell)$	9. $\tau_1 \subseteq \tau_2$
5. $\tau_0 = \text{Arith_Result}(\tau_1, \tau_2)$	

The constraint `Equality_Compatible`(τ_1, \dots, τ_n) is analysed in section 3.4.1. The second constraint `Greater_Less_Than_Compatible`(τ_1, \dots, τ_n) is useful for ensuring typability for queries like `q1 < q2`. $\tau = \text{Member_Type}(\sigma)$ implies that σ is a collection type and that τ is the type of its members. $\tau = \text{Constructor_Member_Type}(\sigma, \ell)$ is used to denote that σ is a class or struct type with at least one member ℓ of type τ . If σ is a class type then ℓ can be any of its properties, relationships or methods. $\tau_0 = \text{Arith_Result}(\tau_1, \tau_2)$ is needed for the type inference of queries of the form `q1 op q2` where `op` $\in \{+, -, /, *\}$. Likewise, $\tau_0 = \text{Merge_Result}(\tau_1, \tau_2)$ arises as a constraint from inferring the type of `union`, `intersect` or `except` expressions. $\tau_0 = \text{Distinct_Result}(\tau_1)$ 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 . Finally, type equality and GSR (Generalisation-Specialisation Relationship) are handled by constraints 8 and 9.

3.4.1 Collections - Membership Type Compatibility

The first two constraints refer to type compatibility wrt 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(q1, ..., qn)`, which involves 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 Standard 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 would be `struct(x : int, y : int)` and `constructor(x : int)` respectively. Their least upper bound (LUB) 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 $\text{coll}_1, \text{coll}_2 \in \{\text{collection}, \text{set}/\text{bag}, \text{list}/\text{array}, \text{set}, \text{bag}, \text{list}, \text{array}\}$ and either the collection constructors are the same or one is child of the other in the hierarchy tree and σ is compatible with τ then $\text{coll}_1(\sigma)$ is compatible with $\text{coll}_2(\tau)$.
- Any two class types `class_id1` and `class_id2` 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 the constructor

type $\text{constr}_1(\ell_1: \sigma_1, \dots, \ell_n: \sigma_n, \ell_{11}: \sigma_{11}, \dots, \ell_{1k}: \sigma_{1k})$ is considered to be compatible with $\text{constr}_2(\ell_1: \tau_1, \dots, \ell_n: \tau_n, \ell_{21}: \tau_{21}, \dots, \ell_{2m}: \tau_{2m})$.

- If σ_i is compatible with τ_i , $i = 1, \dots, n$ and $\text{Constructor_Member_Type}(\text{class_id}, \ell_i) = \tau_i$, $i = 1, \dots, n$, then class_id and $\text{constructor}(\ell_1: \sigma_1, \dots, \ell_n: \sigma_n, \ell_{11}: \sigma_{11}, \dots, \ell_{1k}: \sigma_{1k})$ are compatible types, 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 a constructor type $\sigma = \text{constr}_1(\ell_1: \sigma_1, \dots, \ell_n: \sigma_n, \ell_{11}: \sigma_{11}, \dots, \ell_{1k}: \sigma_{1k})$ is compatible with $\tau = \text{constr}_2(\ell_1: \tau_1, \dots, \ell_n: \tau_n, \ell_{21}: \tau_{21}, \dots, \ell_{2m}: \tau_{2m})$, where $\text{constr}_1, \text{constr}_2 \in \{\text{constructor}, \text{struct}\}$ and constr_1 is a child of (or the same as) constr_2 then $\text{CUB}(\sigma, \tau) = \text{constr}_1(\ell_1: \text{CUB}(\sigma_1, \tau_1), \dots, \ell_n: \text{CUB}(\sigma_n, \tau_n))$.
- If $\sigma = \text{class_id}_1$ and $\text{Constructor_Member_Type}(\text{class_id}_1, \ell_i) = \sigma_i$, $i = 1, \dots, n$ and $\tau = \text{constructor}(\ell_1: \tau_1, \dots, \ell_n: \tau_n, \ell_{21}: \tau_{21}, \dots, \ell_{2m}: \tau_{2m})$ then $\text{CUB}(\sigma, \tau)$ is the least superclass of class_id_1 , say class_id_2 , satisfying the following condition: For all ℓ'_j , ℓ'_j is a property or a relationship of class_id_2 , if $\text{Constructor_Member_Type}(\text{class_id}_2, \ell'_j) = \phi_j$ then there exists k , $1 \leq k \leq n$, s.t. $\ell'_j = \ell_k \wedge \text{CUP}(\tau_k, \sigma_k) \subseteq \phi_j$, $\forall j = 1, \dots, m$, $m \leq n$.

- If σ is compatible with τ and $\sigma, \tau \subseteq \text{atomic}$ 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 $\text{set}(q_1, \dots, q_n)$. However, the OQL construct q_1 in q_2 raises another issue of a slightly different nature. The Standard [§4.10.8.3] states that if the type of q_2 is $\text{coll}(\tau)$ then the type of q_1 should be τ . We explain why this is not the case in our context. Suppose that the type of q_2 is inferred to be $\text{bag}(\text{struct}(x: \text{int}, y: \text{string}))$; then according to the Standard q_1 should have the type $\text{struct}(x: \text{int}, y: \text{string})$. Since a value of type $\text{struct}(x: \text{int})$ could potentially be added in the collection q_2 (that is, since $\text{struct}(x: \text{int})$ and $\text{struct}(x: \text{int}, y: \text{float})$ are compatible types), there is no reason why q_1 could not be of type $\text{struct}(x: \text{int})$ or even $\text{struct}()$.

The same situation occurs when dealing with a collection of objects of different classes. Suppose $\text{lub_class}(\ell_1: \sigma_1, \dots, \ell_n: \sigma_n)$ 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 $\text{object_class}(\ell'_1: \sigma'_1, \dots, \ell'_m: \sigma'_m)$ is a member of the collection, even if its class is not a subclass of $\text{lub_class}(\ell_1: \sigma_1, \dots, \ell_n: \sigma_n)$. This allows more queries to (safely) type-check, for example:

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

This query does not type-check according to the Standard. In order to be type-correct, there should be an explicit type cast $(\text{School_Teacher})x.father$. 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 the type of $x.father$ should be more specific than the member type or the collection `School_Teachers`. Instead we ensure that $x.father$ could potentially be a member of `School_Teachers`, by adding the constraint $\text{Equality_Compatible}(\sigma, \tau)$, where σ is the type of $x.father$ and τ is the member type of `School_Teachers`.

3.5 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 [Rob65]. 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)) ⊆ 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 $\tau_1 \subseteq \text{set/bag}(\tau)$ and $\tau_2 \subseteq \text{set/bag}(\tau)$ are existing 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, `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.

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 (*resolution*) 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 $\text{type_var} = \tau$ 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.

3.6 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 this section, we present the inference rules for a substantial part of OQL; the rules for the remaining OQL constructs are given in Appendix A.

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

$$\begin{array}{ccc} \frac{}{\mathcal{H}; \mathcal{C} \vdash b : \text{bool} \Rightarrow \mathcal{C}} & \frac{}{\mathcal{H}; \mathcal{C} \vdash i : \text{int} \Rightarrow \mathcal{C}} & \frac{}{\mathcal{H}; \mathcal{C} \vdash f : \text{float} \Rightarrow \mathcal{C}} \\ \frac{}{\mathcal{H}; \mathcal{C} \vdash c : \text{char} \Rightarrow \mathcal{C}} & \frac{}{\mathcal{H}; \mathcal{C} \vdash s : \text{string} \Rightarrow \mathcal{C}} & \frac{}{\mathcal{H} \cup \{x : \sigma\}; \mathcal{C} \vdash x : \sigma \Rightarrow \mathcal{C}} \end{array}$$

There are several rules to deal with various collections (sets, bags, lists, arrays). We just give one representative rule, that concerns the query construct $\text{set}(q_1, \dots, q_n)$. 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 .

$$\begin{array}{c} \frac{\mathcal{H}; \mathcal{C} \vdash q_1 : \sigma_1 \Rightarrow \mathcal{C}_1 \dots \mathcal{H}; \mathcal{C}_{n-1} \vdash q_n : \sigma_n \Rightarrow \mathcal{C}_n}{\mathcal{H}; \mathcal{C} \vdash \text{set}(q_1, \dots, q_n) : \text{set}(\text{CUB}(\sigma_1, \dots, \sigma_n)) \Rightarrow \mathcal{C}_n \wedge \{\text{EqualityCompatible}(\sigma_1, \dots, \sigma_n)\}} \\ \frac{\mathcal{H}; \mathcal{C} \vdash q_1 : \sigma \Rightarrow \mathcal{C}_1}{\mathcal{H}; \mathcal{C} \vdash \text{first}(q_1) : \phi \Rightarrow \mathcal{C}_1 \wedge \{\text{Member.Type}(\sigma) = \phi\} \wedge \{\sigma \subseteq \text{list/array}(\phi)\}} \\ \frac{\mathcal{H}; \mathcal{C} \vdash q_1 : \sigma \Rightarrow \mathcal{C}_1 \quad \mathcal{H}; \mathcal{C}_1 \vdash q_2 : \tau \Rightarrow \mathcal{C}_2}{\mathcal{H}; \mathcal{C} \vdash q_1[q_2] : \phi \Rightarrow \mathcal{C}_2 \wedge \{\tau = \text{int}\} \wedge \{\text{Member.Type}(\sigma) = \phi\} \wedge \{\sigma \subseteq \text{list/array}(\phi)\}} \end{array}$$

$$\frac{\mathcal{H}; \mathcal{C} \vdash q_1: \sigma \Rightarrow \mathcal{C}_1 \quad \mathcal{H}; \mathcal{C}_1 \vdash q_2: \tau \Rightarrow \mathcal{C}_2}{\mathcal{H}; \mathcal{C} \vdash q_1 \text{ in } q_2: \text{bool} \Rightarrow \mathcal{C}_2 \wedge \{\text{Equality_Compatible}(\sigma, \text{Member_Type}(\tau))\}}$$

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. As is explained earlier, 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 .

$$\frac{\mathcal{H}; \mathcal{C} \vdash q_1: \sigma_1 \Rightarrow \mathcal{C}_1 \quad \dots \quad \mathcal{H}; \mathcal{C}_{n-1} \vdash q_n: \sigma_n \Rightarrow \mathcal{C}_n}{\mathcal{H}; \mathcal{C} \vdash \text{class_name}(\ell_1: q_1, \dots, \ell_n: q_n): \text{class_name} \Rightarrow \mathcal{C}_n \wedge \{\sigma_1 \subseteq \text{Constructor_Member_Type}(\text{class_name}, \ell_1)\} \wedge \dots \wedge \{\sigma_n \subseteq \text{Constructor_Member_Type}(\text{class_name}, \ell_n)\}}$$

$$\frac{\mathcal{H}; \mathcal{C} \vdash q_1: \sigma_1 \Rightarrow \mathcal{C}_1 \quad \dots \quad \mathcal{H}; \mathcal{C}_{n-1} \vdash q_n: \sigma_n \Rightarrow \mathcal{C}_n}{\mathcal{H}; \mathcal{C} \vdash \text{struct}(\ell_1: q_1, \dots, \ell_n: q_n): \text{struct}(\ell_1: \sigma_1, \dots, \ell_n: \sigma_n) \Rightarrow \mathcal{C}_n}$$

$$\frac{\mathcal{H}; \mathcal{C} \vdash q_1: \tau \Rightarrow \mathcal{C}_1}{\mathcal{H}; \mathcal{C} \vdash q_1. \ell: \sigma \Rightarrow \mathcal{C}_1 \wedge \{\text{Constructor_Member_Type}(\tau, \ell): \sigma\}}$$

The inference rules for the existential and the universal quantification follow. It is worth noting that the variable x is bound in query q_2 to the member type of the collection q_1 .

$$\frac{\mathcal{H}; \mathcal{C} \vdash q_1: \sigma_1 \Rightarrow \mathcal{C}_1 \quad \mathcal{H} \cup \{x: \phi\}; \mathcal{C}_1 \vdash q_2: \sigma_2 \Rightarrow \mathcal{C}_2}{\mathcal{H}; \mathcal{C} \vdash \text{exists } x \text{ in } q_1: q_2: \text{bool} \Rightarrow \mathcal{C}_2 \wedge \{\sigma_2 = \text{bool}\} \wedge \{\text{Member_Type}(\sigma_1) = \phi\}}$$

$$\frac{\mathcal{H}; \mathcal{C} \vdash q_1: \sigma_1 \Rightarrow \mathcal{C}_1 \quad \mathcal{H} \cup \{x: \phi\}; \mathcal{C}_1 \vdash q_2: \sigma_2 \Rightarrow \mathcal{C}_2}{\mathcal{H}; \mathcal{C} \vdash \text{forall } x \text{ in } q_1: q_2: \text{bool} \Rightarrow \mathcal{C}_2 \wedge \{\sigma_2 = \text{bool}\} \wedge \{\text{Member_Type}(\sigma_1) = \phi\}}$$

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.

$$\frac{\mathcal{H}; \mathcal{C} \vdash q_1: \sigma \Rightarrow \mathcal{C}_1}{\mathcal{H}; \mathcal{C} \vdash q_1(): \phi \Rightarrow \mathcal{C}_1 \wedge \{\sigma = \text{unit} \rightarrow \phi\}}$$

$$\frac{\mathcal{H}; \mathcal{C} \vdash q_0: \sigma_0 \Rightarrow \mathcal{C}_0 \quad \dots \quad \mathcal{H}; \mathcal{C}_{n-1} \vdash q_n: \sigma_n \Rightarrow \mathcal{C}_n}{\mathcal{H}; \mathcal{C} \vdash q_0(q_1, \dots, q_n): \phi \Rightarrow \mathcal{C}_n \wedge \{\sigma_0 = \tau_1 \times \dots \times \tau_n \rightarrow \phi\} \wedge \{\sigma_1 \subseteq \tau_1\} \wedge \dots \wedge \{\sigma_n \subseteq \tau_n\}}$$

The rules concerning the query constructs $q \text{ binop } q$ or $q \text{ unop } q$ are omitted for space reasons.

We finish by giving the rule for a simple `select` query; the judgements dealing with the remaining OQL expressions are presented in appendix A.

$$\begin{array}{c}
\mathcal{H}; \mathcal{C} \vdash q_1 : \sigma_1 \Rightarrow \mathcal{C}_1 \\
\mathcal{H} \cup \{x_1 : \tau_1\}; \mathcal{C}_1 \vdash q_2 : \sigma_2 \Rightarrow \mathcal{C}_2 \\
\quad \dots \\
\mathcal{H} \cup \{x_1 : \tau_1, \dots, x_{n-1} : \tau_{n-1}\}; \mathcal{C}_{n-1} \vdash q_n : \sigma_n \Rightarrow \mathcal{C}_n \\
\mathcal{H} \cup \{x_1 : \tau_1, \dots, x_n : \tau_n\}; \mathcal{C}_n \vdash q_{01} : \sigma_{01} \Rightarrow \mathcal{C}_{01} \\
\mathcal{H} \cup \{x_1 : \tau_1, \dots, x_n : \tau_n\}; \mathcal{C}_{01} \vdash q_{00} : \sigma_{00} \Rightarrow \mathcal{C}_{00} \\
\hline
\mathcal{H}; \mathcal{C} \vdash \text{select } q_{00} \text{ from } q_1 \text{ as } x_1, \dots, q_n \text{ as } x_n \text{ where } q_{01} : \text{bag}(\sigma_{00}) \Rightarrow \mathcal{C}_{00} \wedge \\
\{\text{Member.Type}(\sigma_1) = \tau_1\} \wedge \dots \wedge \{\text{Member.Type}(\sigma_n) = \tau_n\} \wedge \{\sigma_{01} = \text{bool}\}
\end{array}$$

3.7 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 resolution) or simplify them all in the end after having applied all the inference rules. If the resolution 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} .

3.8 Discussion

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 [BDHS96]. The problem is addressed 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.

Chapter 4

OQL calculus representation

This chapter presents the translation of OQL into its first intermediate representation the monoid comprehension calculus. This calculus formalism has a number of advantages. Firstly, it has the expressive power to cover all OQL constructs. Secondly, being a calculus-style formalism, it supports the concept of object variables and it allows for nesting of expressions in a similar way to OQL. It would be hard to translate OQL directly into a collection-oriented algebraic representation. Thirdly, using forms known as monoid comprehensions, the monoid calculus captures OQL queries in a uniform way. This feature enables the optimization of calculus expressions using only a few rules.

The calculus representation adopted in this dissertation is based on the one proposed by Fegaras and Maier [FM00]. It has also been influenced by the work of Afshar [Afs98] in establishing a theoretical basis for his monoid query language (MQL). However, the calculus formalism adopted in this thesis addresses certain important issues that have not been considered previously. It highlights certain weak points of OQL and it adds new calculus expressions to deal with them; it proposes different ways of mapping certain OQL queries into their calculus representation (wrt [FM00]); finally, it introduces a series of new rules to transform calculus expressions into more efficient forms and gives a detailed algorithm that determines the order in which these rules must be applied.

4.1 Monoid Comprehension Calculus

This section gives an overview of the monoid calculus representation, as presented in previous work [Afs98, FM00]. It defines the notions of monoid, free monoid, monoid comprehension and monoid homomorphism and presents examples of how OQL is mapped to its calculus representation. The role of this section is to set up the theoretical basis for a formal description of our approach regarding the calculus intermediate representation of OQL.

4.1.1 Monoids

One of the benefits of using the monoid comprehension calculus is the uniform representation of OQL queries. For instance, this feature is obvious in the manipulation of multiple collection types. OQL is an expressive language that allows for the use of not only set collections, but also of bags, lists and arrays. The monoid calculus abstracts the similarities of these collections and allows for their uniform representation as collection monoids. Likewise, OQL includes a variety of operations on simple types, like integers and floats. In the context of our calculus, these can also be considered as primitive monoid types and hence they can be treated accordingly.

Definition 1 A monoid of type T is a pair $(\oplus, \mathcal{Z}_\oplus)$, where \oplus is an associative function of type $T \times T \rightarrow T$ and \mathcal{Z}_\oplus is the left and right identify of \oplus [FM00].

An alternative notation used for monoids in [Afs98] is (T, \oplus, \mathcal{Z}) . Function \oplus is known as the *merge* operation that takes as parameters two elements of type T and returns an element of type T . The value \mathcal{Z}_\oplus is the zero element of the monoid $(\oplus, \mathcal{Z}_\oplus)$, i.e. the value such that $\mathcal{Z}_\oplus \oplus x = x \oplus \mathcal{Z}_\oplus = x$ for every x . For instance, integers can be represented by either of the monoids $(+, 0)$ or $(\times, 1)$. The first monoid merges integers based on the addition operation; the zero element is the integer 0. The second monoid merges integers by multiplying them and the zero (neutral) element is therefore the integer 1. Likewise booleans can be represented by two different monoids, merging boolean values using the \wedge or the \vee operations. These monoids on simple types are known as primitive monoids.

In order to represent collections, we need to define the notion of collection monoids [FM00]. These are monoids for which we need to define an additional *unit* function \mathcal{U}_\oplus . This function takes as parameter an element x and returns a singleton collection containing x . Hence, a collection monoid is denoted by a triple $(\oplus, \mathcal{Z}_\oplus, \mathcal{U}_\oplus)$. It may represent any collection parametric type $T_\oplus(\tau)$, such that τ supports comparison for equality. For instance, the monoid $(\cup, \{\}, \lambda x. \{x\})$ corresponds to the collection type $\text{set}(a)$, where a is the type of the elements x .

Monoids are said to be commutative and/or idempotent, if their merge operations are commutative and/or idempotent respectively. For instance, the set monoid $(\cup, \{\}, \lambda x. \{x\})$ is both commutative and idempotent, the integer monoid $+, 0$ is just commutative, whilst the list monoid $(++, [], \lambda x. [x])$ has neither of the two properties.

Based on these properties, Afshar distinguishes three monoid categories:

1. the (unrestricted) category of monoids, M ;
2. the category of commutative monoids, CM ;
3. the category of commutative and idempotent monoids, CIM .

For the purposes of translating OQL queries to calculus formalism, we need to define an additional category, that of idempotent (but not necessarily commutative) monoids (IM).

Fegaras and Maier define a partial order \preceq between monoids based on their properties. They define a mapping ψ from monoids to the set $\{\mathbf{C}, \mathbf{I}\}$ as: $\mathbf{C} \in \psi(\oplus)$ if and only if \oplus is commutative and $\mathbf{I} \in \psi(\oplus)$ if and only if \oplus is idempotent. The partial order \preceq is defined as follows:

$$\otimes \preceq \oplus \equiv \psi(\otimes) \subseteq \psi(\oplus)$$

As it will be explained later on, both the monoid categories by Afshar, and the partial order by Fegaras and Maier are useful in order to determine when a mapping from one monoid to another, i.e. a monoid homomorphism, is well-defined.

4.1.2 Monoid Homomorphisms

Homomorphisms are defined in [Afs98] and [FM95] in slightly different ways. The definition by Afshar [Afs98] is given below.

Definition 2 *A monoid homomorphism between a monoid $(\mathcal{T}, \times, \mathbf{e})$ and a monoid $(\mathcal{M}, *, \mathbf{I})$ is a function $\mathbf{f} : \mathcal{T} \rightarrow \mathcal{M}$ which satisfies the following [Afs98]:*

$$\begin{aligned} \mathbf{f} \mathbf{e} &= \mathbf{I} \\ \mathbf{f}(\mathbf{a} \times \mathbf{b}) &= \mathbf{f}(\mathbf{a}) * \mathbf{f}(\mathbf{b}) \end{aligned}$$

Fegaras and Maier [FM95] define homomorphisms only on collection monoids. They regard homomorphisms as means of expressing OQL queries, none of which corresponds to a homomorphism from a primitive monoid. Hence although the previous definition by Afshar is more general [Afs98], the following one by Fegaras and Maier is sufficient in the context of translating OQL [FM95].

Definition 3 *A homomorphism $\text{hom}[\oplus, \otimes](\mathbf{f})$ from the collection monoid $(\oplus, \mathcal{Z}_{\oplus}, \mathcal{U}_{\oplus})$ to any monoid $(\otimes, \mathcal{Z}_{\otimes})$, where $\oplus \preceq \otimes$ is defined by the following inductive equations [FM00]:*

$$\begin{aligned} \text{hom}[\oplus, \otimes](\mathbf{f})(\mathcal{Z}_{\oplus}) &= \mathcal{Z}_{\otimes} \\ \text{hom}[\oplus, \otimes](\mathbf{f})(\mathcal{U}_{\oplus}(\mathbf{a})) &= \mathbf{f}(\mathbf{a}) \\ \text{hom}[\oplus, \otimes](\mathbf{f})(\mathbf{x} \oplus \mathbf{y}) &= (\text{hom}[\oplus, \otimes](\mathbf{f})(\mathbf{x})) \otimes (\text{hom}[\oplus, \otimes](\mathbf{f})(\mathbf{y})) \end{aligned}$$

In the latter definition, the function $\text{hom}[\oplus, \otimes](\mathbf{f})$ maps \mathcal{Z}_{\oplus} to \mathcal{Z}_{\otimes} , \mathcal{U}_{\oplus} to \mathbf{f} and the merge operation \oplus to \otimes . The third inductive equation reflects the "divide-and-conquer" mechanism implicit in homomorphisms. That is, one may apply the function \mathbf{f} on a collection $\mathbf{x} \oplus \mathbf{y}$, by applying it separately on its parts \mathbf{x} and \mathbf{y} and merging the results using the operation \otimes .

In order for a monoid homomorphism to be well-defined, Fegaras and Maier require the constraint $\oplus \preceq \otimes$ be hold. For instance, a homomorphism may map a list to a bag or to a set, but it

could not possibly map a bag to a list or a set to a bag. Fegaras and Maier set additional restrictions to ensure that homomorphisms are well-defined. If a monoid $(\oplus, \mathcal{Z}_\oplus)$ is not commutative it should be anti-commutative, i.e. $\forall x \neq y \neq \mathcal{Z}_\oplus : x \oplus y \neq y \oplus x$. Likewise, if a monoid is not idempotent it should be anti-idempotent, i.e. $\forall x \neq \mathcal{Z}_\oplus : x \oplus x \neq x$. However, these restrictions are not always satisfied in their framework. For instance, the primitive monoid $(\times, 1)$ is not idempotent. However there exists an element $x = 0 \neq 1$ such that $0 \times 0 = 0$.

In this thesis, we use the notions of monoid category and free monoid (formally presented in [Afs98]) in order to express the necessary (but not sufficient) conditions for the well-definedness of monoid homomorphisms. A slightly rephrased definition of free monoids (found in [[Afs98]],ch.6) is given below:

Definition 4 A free Monoid $(\mathcal{F}, *, e)$, where \mathcal{F} is generated by the injection map $i : \mathcal{S} \rightarrow \mathcal{F}$, is such that [Afs98]:

- $\forall f : \mathcal{S} \rightarrow \mathcal{M}$ – where f is a set function and \mathcal{M} is a monoid
- $\exists h : \mathcal{F} \rightarrow \mathcal{M}$ – where h is a monoid homomorphism
- s.t. $h \circ i = f$

In order for a homomorphism $f : \mathcal{M} \rightarrow \mathcal{N}$ to be well-formed, \mathcal{N} should lie in the same category as \mathcal{M} . In addition to that, if \mathcal{N} is a free monoid then \mathcal{M} should also be free.

4.1.3 Monoid Comprehensions

Fegaras and Maier use homomorphisms to define formal semantics for monoid comprehensions. A monoid comprehension has the form $\oplus\{e \parallel q_1, \dots, q_n\}$, where \oplus is called the accumulator of the comprehension, e is the head and q_i are the qualifiers of the comprehension. A qualifier takes either of the following forms:

- $u \leftarrow e'$: a *generator* binds elements of the collection e' to the variable u
- $u' \equiv e''$: a *binding* binds the variable u' to the value of the expression e''
- e''' : a *filter* may include variables bound in previous generators or bindings.

Definition 5 A monoid comprehension over a primitive or collection monoid \oplus is defined by the following inductive equations: (\otimes is the collection monoid associated with the expression u and is possibly different from \oplus) [FM95]

$$\begin{aligned} \oplus\{e \parallel \} &= \begin{cases} \mathcal{U}_\oplus(e) & \text{for a collection monoid } \oplus \\ e & \text{for a primitive monoid } \oplus \end{cases} \\ \oplus\{e \parallel x \leftarrow u, \bar{q}\} &= \text{hom}[\otimes, \oplus](\lambda x. \oplus\{e \parallel \bar{q}\})u \\ \oplus\{e \parallel \text{pred}, \bar{q}\} &= \text{if pred then } \oplus\{e \parallel \bar{q}\} \text{ else } \mathcal{Z}_\oplus \end{aligned}$$

where $\bar{q} = q_1, \dots, q_n$, $n \geq 0$.

An example is given that transforms a select query into its calculus counterpart.

```
select x
from Countries as x
where x.population > 50,000,000
```

The OQL query above is translated to the following comprehension:

$$\cup \{x \mid x \leftarrow \text{Countries}, x.\text{population} > 50,000,000\}$$

Fegaras and Maier present the syntax of their monoid comprehension calculus, as well as the translation rules for most OQL constructs. The interested reader may refer to [FM95, FM00] for the detailed rules that map queries into their calculus representation.

The calculus adopted in this thesis departs from their proposal in certain respects. Hence some of the translation rules from OQL to calculus form and certain optimization rules are different. The remaining sections highlight the points at which the current approach differs from their work and explains our reasons for proposing a number of modifications and extensions.

4.2 Extended Calculus: syntax and monoids

4.2.1 Calculus Syntax

The grammar of the monoid comprehension calculus adopted in this thesis is given in figure 4.1. The calculus expressions that are new in this framework (wrt [FM00]) are:

```
< class_name > ( $\ell_1 = e_1, \dots, \ell_k = e_k$ ),
e1[e2], not (e),
bagof(e), listof(e), arrayof(e)
orderof(e), elementof(e)
```

The construct $\langle \text{class_name} \rangle (\ell_1 = e_1, \dots, \ell_k = e_k)$ has been added in order to deal with the construction of objects. The second one ($e_1[e_2]$) accesses the e_2 -th element of the collection e_1 , whilst $\text{not}(e)$ obviously handles boolean negation. The last four expressions $\text{bagof}(e)$, $\text{listof}(e)$, $\text{arrayof}(e)$ and $\text{orderof}(e)$ have been added in order to avoid the problem of ill-formed comprehensions. The partial function $\text{elementof}(e)$ is the inverse of unit , that is, it takes a singleton collection and returns its element. If the collection has no elements or more than one element, an error is generated. If e has type $\text{coll}(\tau)$, then the types of $\text{bagof}(e)$, $\text{listof}(e)$, $\text{arrayof}(e)$, $\text{orderof}(e)$ and $\text{elementof}(e)$ are $\text{bag}(\tau)$, $\text{list}(\tau)$, $\text{array}(\tau)$, $\text{list}(\tau)$ and τ respectively.

```

e, e1, ..., en ::= literal (integer, float, boolean, char, string)
| < id > (db object/class/extent name)
| zero[< monoid >]
| unit[< monoid >](e)
| merge[< monoid >](e1, ..., en)
| < monoid > {e || qualifiers}
| struct(l1 = e1, ..., lk = ek)
| < class_name > (l1 = e1, ..., lk = ek)
| e1 < comparison_operator > e2
| e1 < arithmetic_operator > e2
| e. < attr_rltp_method >
| e(e1, ..., ek)
| e1[e2]
| not(e)
| bagof(e)
| listof(e)
| arrayof(e)
| orderof(e)
| elementof(e)

```

where

- < comparison_operator > ∈ {<, >, ≤, ≥, =, ≠},
- < arithmetic_operator > ∈ {+, −, ×, /, mod},
- < attr_rltp_method > is the name of an attribute, relationship or method of a class,
- < monoid > is either a primitive or a collection monoid,
- zero[< monoid >], unit[< monoid >](e), and merge[< monoid >](e₁, ..., e_n) are alternative representations of the zero element \mathcal{Z}_{\oplus} , the unit function \mathcal{U}_{\oplus} and the merge operation \oplus of some < monoid > ($\mathcal{Z}_{\oplus}, \oplus$).

Figure 4.1: Syntax of the calculus OQL representation

The ODMG standard defines that every collection in the `from` clause of a `select` query is converted into a bag ([CBB⁺00],4.10.9). For instance, consider the query:

```
select x
from Students as x
where x.grade = 10
```

Although the extent `Students` has type `set(Student)`, the result of the `from` clause is a bag of students. Hence it would be wrong to translate the query above to the calculus form:

$$\text{bag}\{x \mid x \leftarrow \text{Students}, x.\text{grade} = 10\}$$

This is an ill-formed comprehension, since there is no homomorphic mapping from a set to a bag. To avoid this problem, we first apply the non-homomorphic function `bagof(-)` to the set `Students`. The resulting calculus expression is:

$$\text{bag}\{x \mid x \leftarrow \text{bagof}(\text{Students}), x.\text{grade} = 10\}$$

The non-homomorphic function `listof(-)` is useful in a number of contexts. First it deals with the non-determinism introduced when ordering the elements of a collection based on insufficient sort criteria. The Standard [CBB⁺00] states the following:

”The order by a set of functions (f_1, \dots, f_n) is performed as follows: First sort according to function f_1 , then for all elements having the same f_1 value sort them according to f_2 , and so on ([CBB⁺00],ch 4.10.9.4)”.

The question that arises is “what if there are several elements with the same values of f_1, \dots, f_n ?”. For instance, the result of the following query is not determined in the case that at least two students have the same grade.

```
select x.name
from Students as x
order by x.grade
```

To this question, the ODMG standard does not give an adequate solution. It is implied that whether the ordering is *complete* is decided at run-time. If it is not complete, then either the query produces a non-deterministic result, or alternatively, it raises an error and it is rejected. In the former case, it is possible that the user gets different results every time s/he submits the same query; in the latter case, the user needs to wait until execution time to be notified about a possible error.

We suggest a different way of handling non-determinism. First we define the collection monoid `oset`, as a means of handling an ordered collection of elements with no duplicates. `oset` is defined

over a parametric type τ , such that the values of this type can be compared for equality. Its zero element and the unit and merge operations are defined as follows:

- $\text{zero}[\text{oset}] = \|\ \|\$
- $\text{unit}[\text{oset}](e) = \|\ e \|\$
- $\text{merge}[\text{oset}](\|\ e_1, \dots, e_n \|\, \|\ e \|\)$ is evaluated as follows:
 - If there exists i , s.t. $e_i = e$, then $\text{merge}[\text{oset}](\|\ e_1, \dots, e_n \|\, \|\ e \|\) = $\|\ e_1, \dots, e_n \|\$$
 - Otherwise, $\text{merge}[\text{oset}](\|\ e_1, \dots, e_n \|\, \|\ e \|\) = $\|\ e_1, \dots, e_n, e \|\$$
- $\text{merge}[\text{oset}](\|\ e_1, \dots, e_n \|\, \|\ e'_1, \dots, e'_m \|\) = $\text{merge}[\text{oset}](\text{merge}[\text{oset}](\|\ e_1, \dots, e_n \|\, \|\ e'_1 \|\), \|\ e'_2, \dots, e'_m \|\)$$

Inserting an existing element into an ordered set (*oset*) has no effect. Inserting a new element in an ordered set is identical to inserting it into a list. Thus, *oset* is idempotent, but not commutative. *oset* is free in the idempotent category of monoids (IM).

The following proposal for handling non-determinism departs from the Standard significantly:

- A `select · distinct · . . . · order by e_1, \dots, e_n` clause yields an *oset* of bags of elements. Each bag contains elements having the same values for all sort criteria e_1, \dots, e_n . If no two elements have the same values for all the criteria, then the result is an *oset* of singleton bags.
- A `select · . . . · order by e_1, \dots, e_n` clause (without `distinct`) yields a list of bags of elements. Each bag contains elements having the same values for all sort criteria e_1, \dots, e_n . If no two elements have the same values for all the criteria, then the result is a list of singleton bags.

Intuitively, the result of a `select · [distinct] · . . . · order by . . .` query is expected to be a list or an ordered set of elements. However, the solution proposed above returns a list or an ordered set of bags. Thus, we need a new OQL construct `order_flatten` that flattens this result into the expected collections. Non-determinism is inevitably introduced by flattening a list of non-singleton bags into a list of elements. The ordering of elements in the same bag is determined by system-dependent algorithms (`listof(set/bag(elements))`), and is therefore considered arbitrary in an open distributed world. Although this approach does not avoid non-determinism, it has the advantage that it avoids applying the user's sort criteria and the arbitrary sorting algorithm in one step. The user can see first the result of his sorting expression. On his conscious request (by the use of `order_flatten`), he can then get a flattened list or ordered set, especially for presentation or printing purposes.

The following table illustrates the result of applying the `order_flatten` on different nested collections. It is defined in a way corresponding to the discussion about `flatten` in the ODMG standard [CBB⁺00]. In what follows the notation $\text{coll}_1/\text{coll}_2$ (e.g. `set/bag`) denotes alternatives, not general types.

coll_1	coll_2	$\text{order_flatten}(\text{coll}_1(\text{coll}_2(\text{elem})))$
<code>set/bag</code>	<code>set/bag</code>	<code>coll₂(elem)</code>
<code>list/oset/array</code>	<code>set/bag</code>	<code>coll₁(elem)</code>
<code>set/bag</code>	<code>list/oset/array</code>	<code>coll₁(elem)</code>
<code>list/oset/array</code>	<code>list/oset/array</code>	<code>coll₂(elem)</code>

This section presented the calculus expressions used in our framework, and discussed the role of the nonhomomorphic functions `bagof` and `listof`. The functionality of `arrayof` and `orderof` remains to be discussed in section 4.2.3. In order to translate OQL queries into this calculus representation, we first need to define the monoids involved in the translation. In the following two sections, we pay particular attention to the new monoids (wrt [FM00]) that have been introduced in this dissertation, e.g. the `group` and the `dictionary` monoids. We also present several translation rules from OQL to calculus that involve these monoids. The complete set of translation rules is given in Appendix B.

4.2.2 Monoids

The primitive and collection monoids used in our context are the following:

Primitive Monoids	Collection Monoids
<code>sum (commutative)</code>	<code>set (commutative, idempotent)</code>
<code>avg (commutative)</code>	<code>bag (commutative)</code>
<code>some (commutative, idempotent)</code>	<code>list</code>
<code>all (commutative, idempotent)</code>	<code>array</code>
<code>max (commutative, idempotent)</code>	<code>oset (idempotent)</code>
<code>min (commutative, idempotent)</code>	<code>dictionary</code>
	<code>group (commutative)</code>

Some of the monoids above are identical to those used in relevant calculus frameworks. The interested reader may refer to [FM00] for their descriptions. The remainder of this section defines only those monoids that have been introduced in the current framework.

The group monoid

The `group` monoid in our context is defined over an unordered collection of structures of type `struct(< ℓ_1 >: $\tau_1, \dots, < \ell_n >: \tau_n, \text{partition} : \text{bag}(\tau)$)`, such that types τ_i , are comparable types

and τ is any nonfunctional type. Informally, `group` is used to manipulate structures with at least one label called `partition` and at least one other label. All labels except `partition` are called *group-by attributes*. All structures in the collection have the same labels. In order to insert a new structure in the group monoid, we check whether there exists a structure such that its *group-by* attributes have the same values as the corresponding attributes of the new structure. If there exists such a structure, then its `partition` (which is a bag) is updated, by being merged with the `partition` of the new structure. Otherwise, the new structure is inserted, as it would be normally inserted into a bag.

More formally,

- $\text{zero}[\text{group}] = \langle\langle\rangle\rangle$
- $\text{unit}[\text{group}](\text{struct} \langle \ell_1 = e_1, \dots, \ell_j = e_j, \text{partition} = e \rangle) = \langle\langle \text{struct} \langle \ell_1 = e_1, \dots, \ell_j = e_j, \text{partition} = e \rangle \rangle\rangle$
- Consider $\text{struct}_i = \text{struct} \langle \ell_1 = v_{i1}, \dots, \ell_j = v_{ij}, \text{partition} = v_i \rangle, i = 1, \dots, n$ and $\text{struct}' = \langle \ell_1 = v_{k1}, \dots, \ell_j = v_{kj}, \text{partition} = v \rangle$.
 - If $1 \leq k \leq n$, $\text{merge}[\text{group}](\langle\langle \text{struct}_1, \dots, \text{struct}_n \rangle\rangle, \langle\langle \text{struct}' \rangle\rangle) = \langle\langle \text{struct}_1, \dots, \text{struct}_{k-1}, \text{struct} \langle \ell_1 = v_{k1}, \dots, \ell_j = v_{kj}, \text{partition} = \text{merge}[\text{bag}](v_k, v) \rangle, \text{struct}_{k+1}, \dots, \text{struct}_n \rangle\rangle$.
 - Otherwise, if none of the existing structures has the same values for the *group-by* attributes as struct' , $\text{merge}[\text{group}](\langle\langle \text{struct}_1, \dots, \text{struct}_n \rangle\rangle, \langle\langle \text{struct}' \rangle\rangle) = \langle\langle \text{struct}_1, \dots, \text{struct}_n, \text{struct}' \rangle\rangle$.
- $\text{merge}[\text{group}](\langle\langle \text{struct}_1, \dots, \text{struct}_n \rangle\rangle, \langle\langle \text{struct}'_1, \dots, \text{struct}'_m \rangle\rangle) = \text{merge}[\text{group}](\text{merge}[\text{group}](\langle\langle \text{struct}_1, \dots, \text{struct}_n \rangle\rangle, \langle\langle \text{struct}'_1 \rangle\rangle), \langle\langle \text{struct}'_2, \dots, \text{struct}'_m \rangle\rangle)$

`group` is commutative, since the order in which we merge structures does not make any difference. We just compare values of the *group-by* attributes, and if they match we merge their partitions. The latter are simple bags, so their merging is a commutative operation. `group` is not idempotent, since if we merge two singleton group monoids with identical structures, the `partition` of the result contains twice the elements of each initial partition.

`group` is more specific than the corresponding $\text{grouped}(V \times N)$ monoid defined in [Afs98]. It accumulates all elements with the same grouping values into a bag, whilst $\text{grouped}(V \times N)$ groups on V values into any monoid N . In the case of `group` the elements are structures, the grouping values are

paired with corresponding labels and all elements having the same grouping values are merged into a bag with label `partition`. Instead, the merge operation of `grouped(V × N)` accumulates tuples, that is neither the grouping values nor the grouped elements have corresponding labels. This is not possible in OQL, the type system of which does not support tuples.

The dictionary monoid

`dictionary` represents an unordered collection of structures of type `struct(key : σ, value : τ)`, such that values of type `σ` can be compared for equality and `τ` is any nonfunctional type. Before inserting a new structure into a dictionary, we check whether the latter already contains a structure with the same key. If this is the case, the existing structure `struct(key : k, value : old_value)` is replaced by the new one `struct(key : k, value : new_value)`. Otherwise, the latter is inserted as it would be inserted into a bag. The notation `(k, v)` is often used as an alternative representation of the expression `struct(key : k, value : v)`. More formally,

- $\text{zero}[\text{dict}] = \langle \rangle$
- $\text{unit}[\text{dict}]((k, v)) = \langle (k, v) \rangle$
- If there exists i , $1 \leq i \leq n$ such that $k_i = k$, then

$$\begin{aligned} & \text{merge}[\text{dict}](\langle (k_1, v_1), \dots, (k_i, v_i), \dots, (k_n, v_n) \rangle, \langle (k, v) \rangle) \\ &= \langle (k_1, v_1), \dots, (k_i, v), \dots, (k_n, v_n) \rangle \end{aligned}$$
- If there is no i , $1 \leq i \leq n$ such that $k_i = k$, then

$$\begin{aligned} & \text{merge}[\text{dict}](\langle (k_1, v_1), \dots, (k_n, v_n) \rangle, \langle (k, v) \rangle) \\ &= \langle (k_1, v_1), \dots, (k_n, v_n), (k, v) \rangle \end{aligned}$$

The merge operation of the monoid `dictionary` is idempotent. That is, for each pair `(k, v)` $\text{merge}[\text{dict}](\langle (k, v) \rangle, \langle (k, v) \rangle) = \langle (k, v) \rangle$. `dictionary` is not a commutative monoid, since $\text{merge}[\text{dict}](\langle (1, 2) \rangle, \langle (1, 3) \rangle) = \langle (1, 3) \rangle$, whilst $\text{merge}[\text{dict}](\langle (1, 3) \rangle, \langle (1, 2) \rangle) = \langle (1, 2) \rangle$.

Notice that the merge operation of `dictionary` resembles the corresponding operation of the monoid `array`, defined by Fegaras and Maier [FM00]. A significant difference is that the key in their context is always integer, since it reflects the position of the `value` in the array.

4.2.3 Translation rules

The previous subsections defined the monoids that are novel in this calculus framework. This section uses these monoids to present a set of translation rules from OQL constructs to their calculus representation. In fact, it only presents those translation rules that differ from the corresponding ones in

OQL Expressions ¹	Calculus Expressions
<pre>select [distinct] f from q₁ as x₁, ..., q_k as x_k where p</pre>	$\text{bag } [\text{set}] \{ f ^* \parallel x_1 \leftarrow \text{bagof}(q_1 ^*), \dots, \\ x_k \leftarrow \text{bagof}(q_k ^*), p ^* \}$
<pre>select [distinct] f from q₁ as x₁, ..., q_k as x_k where p order by o₁ [asc/desc], ..., o_j [asc/desc]</pre>	$\text{list } [\text{oset}] \{ v_1.\text{partition} \parallel v_1 \leftarrow \text{orderof}(\text{group}\{ \\ \text{struct}(\text{asc/desc}_1 = o_1 ^*, \dots, \text{asc/desc}_j = o_j ^*, \\ \text{partition} = \text{unit}[\text{bag}] (f ^*) \parallel x_1 \leftarrow \text{bagof}(q_1 ^*), \\ \dots, x_k \leftarrow \text{bagof}(q_k ^*), p ^*) \}) \}$
<pre>select [distinct] f from q₁ as x₁, ..., q_k as x_k where p group by l₁ : g₁, ..., l_j : g_j having h order by o₁ [asc/desc], ..., o_n [asc/desc]</pre>	$\text{list } [\text{oset}] \{ v_1.\text{partition} \parallel v_1 \leftarrow \text{orderof}(\text{group}\{ \\ \text{struct}(\text{asc/desc}_1 = o_1 ^*, \dots, \text{asc/desc}_n = o_n ^*, \\ \text{partition} = \text{unit}[\text{bag}] (f) \parallel v_2 \leftarrow \text{bagof}(\text{group}\{ \\ \text{struct}(l_1 = g_1 ^*, \dots, l_j = g_j ^*, \text{partition} = \\ \text{unit}[\text{bag}] (\text{struct}(x_1 = x_1, \dots, x_k = x_k)) \parallel x_1 \leftarrow \\ \text{bagof}(q_1 ^*), \dots, x_k \leftarrow \text{bagof}(q_k ^*), p ^*) \}, l_1 \equiv v_2.l_1, \\ \dots, l_j \equiv v_2.l_j, \text{partition} \equiv v_2.\text{partition}, h ^*) \}) \}$

Figure 4.2: Translation of select queries to their calculus representation

other calculus frameworks [FM00, Afs98]. The complete set of translation rules from OQL to the monoid calculus is given in Appendix B.

We use the notation $|q|^*$ to denote the calculus equivalent of query q ($| - |^* : \text{OQL} \rightarrow \text{Calculus}$).

First the calculus expression $e_1[e_2]$ enables accessing ordered collections. In order to translate the query $\text{last}(q)$, the primitive monoid sum is used. Since this monoid is not idempotent, the elements of the collection q are counted after transforming it into a bag.

$$\begin{aligned} |\text{first}(q)|^* &= |q|^*[0] \\ |\text{last}(q)|^* &= |q|^*[\text{sum}\{1 \parallel x \leftarrow \text{bagof}(|q|^*)\} - 1] \\ |q_1[q_2]|^* &= |q_1|^* [|q_2|^*] \end{aligned}$$

Second, structures and classes are translated to calculus as follows:

$$\begin{aligned} |\text{struct}(l_1 : q_1, \dots, l_n : q_n)|^* &= \text{struct}(l_1 = |q_1|^*, \dots, l_n = |q_n|^*) \\ |C(l_1 : q_1, \dots, l_n : q_n)|^* &= C(l_1 = |q_1|^*, \dots, l_n = |q_n|^*) \end{aligned}$$

The translation of a dictionary expression is straightforward.

¹ f, g_i, h, p, q_i and o_i represent OQL queries in figure 4.2.

$$|\text{dictionary}((q_1, q'_1), \dots, (q_n, q'_n))|^* = \text{merge}[\text{dict}](\text{unit}[\text{dict}]((|q_1|^*, |q'_1|^*), \dots, \text{unit}[\text{dict}]((|q_n|^*, |q'_n|^*))))$$

The translation of `select` queries is given in figure 4.2. All collections in the `from` clauses are first converted to bags. Moreover, the monoid group and a new function `orderof` are used to translate queries with a `group by` and an `order by` clause. The function `orderof` operates on a collection of structures of the form:

$$\text{struct}(\langle \text{sc}_1 \rangle = v_1, \dots, \langle \text{sc}_n \rangle = v_n, \text{partition} = v)$$

where $\langle \text{sc}_i \rangle ::= \text{asc}_i | \text{desc}_i$ and the types τ_i of v_i are sortable types, for all $i = 1, \dots, n$.

All structures in the collection have the same labels. The labels $\langle \text{sc}_i \rangle$ are called *sort criteria*. By convention, their names start with `asc` or `desc` and are followed by the serial number of the corresponding criterion. For instance, if there are two sort criteria in the structures, the first being ascending and the second descending, their names would be `asc1` and `desc2` respectively. The function `orderof` sorts these structures based on their sort criteria. It is assumed that there are no structures with the same values for the sort criteria. The reason why is that `orderof` always applies to a `group` comprehension, such that the `group-by` attributes coincide with the sort criteria. Being a nonhomomorphic function, it waits until all the elements of its operand are calculated, before it starts ordering them.

The non-homomorphic function `bagof` is applied to the result of `group`, since there is no homomorphic function from `group` to `bag` or between two `group` monoids with different grouping attributes (labels).

The result of a `select` query with an `order by` clause is not a `list` of elements, but a `list` or `oset` of bags of elements. As is explained in section 4.2.1 a new OQL construct `order_flatten` is introduced to flatten the result into a `list` or `oset` of elements. The calculus forms of the new construct are given below:

$$\begin{aligned} |\text{order_flatten}(q)|^* &= \text{list}\{x \parallel s \leftarrow |q|^*, x \leftarrow \text{listof}(s)\} & q : \text{list}(\text{bag}(-)) \\ |\text{order_flatten}(q)|^* &= \text{oset}\{x \parallel s \leftarrow |q|^*, x \leftarrow \text{listof}(s)\} & q : \text{oset}(\text{bag}(-)) \end{aligned}$$

The ODMG defines that the OQL construct `distinct` converts a bag into a set. However, the result of the `distinct` operation on a list remains a list. In this framework, `distinct` converts lists or arrays into ordered sets (`osets`).

$$\begin{aligned} |\text{distinct}(q)|^* &= \text{set}\{x \parallel x \leftarrow |q|^*\} & q : \text{set}/\text{bag} \\ |\text{distinct}(q)|^* &= \text{oset}\{x \parallel x \leftarrow |q|^*\} & q : \text{list}/\text{oset} \\ |\text{distinct}(q)|^* &= \text{oset}\{x \parallel x \leftarrow \text{listof}(|q|^*)\} & q : \text{array} \end{aligned}$$

When applied on an ordered collection, the non-homomorphic function `listof` converts its operand into a list respecting the prespecified order of its elements.

Another point at which this framework extends the Standard is handling collection operations, i.e. `union`, `except` and `intersect`. The Standard allows only sets and bags to be the operands of such constructs. It is interesting to look at a possible extension of the Standard, that allows the `union` operation between lists, arrays and ordered sets. However, we do not attempt to define the operations `except` and `intersect` between lists, ordered sets or arrays. The reason why is that when c_1, c_2 are ordered collections, it is difficult to satisfy the following property:

$$\text{union}(\text{except}(c_1, c_2), \text{intersect}(c_1, c_2)) = c_1 \quad (4.1)$$

The following table presents the result of a union operation over different collections.

```

union : set(a) * set(a) → set(a)
      : set(a)/bag(a) * bag(a) → bag(a)
      : oset(a) * oset(a) → oset(a)
      : array(a) * array(a) → array(a)
      : oset(a)/list(a)/array(a) * list(a) → list(a)
      : oset(a) * array(a) → list(a)

```

In addition, if sets or bags are merged with lists, osets or arrays, lists and arrays are converted into bags and osets into sets ($\text{set}(a) * \text{list}(a) \Leftrightarrow \text{set}(a) * \text{bag}(a)$). If arrays are merged with other collection types, they are converted into lists. The only issue that needs to be resolved at this point is how the union of two arrays is evaluated.

Manipulating arrays

The Standard does not define a union operation between arrays. However, it allows flattening an array of arrays by merging the elements of the nested arrays. It is not clear whether this implies a list concatenation, where `NULL` values are preserved, or a merge operation (\square) of vectors, as defined by Fegaras and Maier ([FM00], section 5.3). In the first case, two arrays `[1, NULL, 3]` and `[2, 5, NULL, 12, NULL]` would be merged into `[1, NULL, 3, 2, 5, NULL, 12, NULL]`. This would make the use of arrays trivial, since their manipulation in the context of OQL would be exactly the same as that of lists. In the second case [FM00], vectors (or arrays) are represented as sets of pairs (i, e) , where i is the index of element e . Merging the arrays of the previous example, represented as $\{(1, 1), (3, 3)\}$ and $\{(1, 2), (2, 5), (4, 12)\}$, would result in the array $\{(1, 2), (2, 5), (3, 3), (4, 12)\}$. That is, a new value in position i of the second array overwrites the old value in the same position of the first array. Although this approach reflects the way vectors are manipulated in programming languages (assigning a value to a vector position), it creates practical problems in the context of OQL.

Consider for instance the OQL query $\text{exists } x \text{ in } (q_1 \text{ union } q_2) : p(x)$, where q_1, q_2 are arrays. The calculus form of this query is expected to be:

$$\text{some}\{|p(x)|^* \parallel x \leftarrow \text{bagof}(M_{\square}(|q_1|^*, |q_2|^*))\}$$

The function `bagof` is used to ensure that the arrays are merged completely, before the elements of the resulting array are bound to x . Since `some` is idempotent, one could possibly use a similar function `setof`; in practice this would make no difference, since there are no duplicate pairs in the resulting array. x is bound to elements of the form (i, e) whereas we expect it intuitively to be bound to e . The practical problem that arises is that we would have to introduce different translation rules for the same OQL expression. For instance, $\text{exists } x \text{ in } q_1 : q_2$, would be translated to two different calculus forms, depending on whether q_1 is an array or not.

In order to avoid this problem, we define that when the non-homomorphic function `bagof(-)` (or `listof(-)`) applies to an array, it creates a bag (or list) with elements the values e corresponding to each pair (i, e) in the initial array. Hence, in this thesis we adopt the monoid introduced in [FM00] for the manipulation of vectors and we refer to it as `array`. This monoid is idempotent, but not commutative. By adjusting the definitions of `bagof` and `listof` to handle arrays, we manage to manipulate arrays and other collections in a uniform way. Finally, we introduce the function `arrayof : list/oset(a) → array`, which takes each element e in position i of its operand, creates a pair (i, e) and adds it to the resulting array.

Sections 4.2.1 and 4.2.2 presented the main novel points of the calculus representation adopted in this thesis. As has already been pointed out, calculus forms can easily be transformed into simpler and more efficient forms. The next section presents a set of rules for *normalizing* queries; these rules are used to *unnest* complex queries in a number of examples; finally a normalization algorithm determines the order in which the rules must be applied.

4.3 Normalization of calculus expressions

4.3.1 Normalization Rules

Fegaras and Maier have proposed ten rewrite rules that can be used to transform calculus expressions into canonical forms:

$$\oplus\{e \parallel u_1 \leftarrow \text{path}_1, \dots, u_n \leftarrow \text{path}_n, \text{pred}\}$$

where path_i is a bound variable, the name of a class extent, the name of an object, or an expression $\text{path}'.\text{name}$. Their set of rewrite rules is extended with new rules that handle the additional features of the current calculus representation.

This section presents the set of normalization rules for calculus expressions and pays particular attention to the rules first introduced in this thesis. In the following rules the letters M and N represent monoids.

1. $M\{e \parallel \bar{q}, u \leftarrow \text{zero}[N], \bar{s}\} \Rightarrow \text{zero}[M]$
2. $M\{e \parallel \bar{q}, u \leftarrow \text{unit}[N](e'), \bar{s}\} \Rightarrow M\{e \parallel \bar{q}, u \equiv e', \bar{s}\}$
3. $M\{e \parallel \bar{q}, u \leftarrow \text{merge}[N](e_1, e_2), \bar{s}\} \Rightarrow$
 $\text{merge}[M](M\{e \parallel \bar{q}, u \leftarrow e_1, \bar{s}\}, M\{e \parallel \bar{q}, u \leftarrow e_2, \bar{s}\})$

For this transformation to be valid, either \bar{q} should be empty or M should be commutative. Say for instance that both N and M are list monoids and \bar{q} evaluates to a list with two elements $[v_1, v_2]$. The left part of the rule first joins v_1 with the elements from $\text{merge}[N](e_1, e_2)$ and then joins v_2 with each one of them. The right part of the rule first joins v_1 with the elements from e_1 , then v_2 with the elements from e_1 , then v_1 with the elements from e_2 and finally v_2 with the elements from e_2 . The left part of the rule obviously evaluates to a different list than the right part. This problem does not occur if \bar{q} is empty. However, for reasons of confluence, i.e. to ensure that the order of applying normalization rules exhaustively does not affect the final result, we only apply this rule when M is commutative.

4. $(\text{struct}(A_1 = e_1, \dots, A_n = e_n)).A_i \Rightarrow e_i$
5. $(\text{class_name}(A_1 = e_1, \dots, A_n = e_n)).A_i \Rightarrow e_i$
6. $M\{e \parallel \bar{q}, u \leftarrow N\{e' \parallel \bar{r}\}, \bar{s}\} \Rightarrow M\{e \parallel \bar{q}, \bar{r}, u \equiv e', \bar{s}\}$
7. $M\{e \parallel \bar{q}, \text{some}\{\text{pred} \parallel \bar{r}\}, \bar{s}\} \Rightarrow M\{e \parallel \bar{q}, \bar{r}, \text{pred}, \bar{s}\}$

Notice that M should be idempotent. For example, if for some q_i derived from \bar{q} there are only two elements r_1, r_2 from \bar{r} that satisfy pred , then q_i will be considered once in the initial expression, but twice in the transformed one. If M is idempotent then this will make no difference in the final result.

8. $M\{e \parallel \bar{q}, \text{merge}[\text{some}](p_1, p_2), \bar{s}\} \Rightarrow$
 $\text{merge}[M](M\{e \parallel \bar{q}, p_1, \bar{s}\}, M\{e \parallel \bar{q}, p_2, \bar{s}\})$

For this transformation to be valid M should be idempotent (as in rule 7).

9. $M\{M\{e \parallel \bar{r}\} \parallel \bar{s}\} \Rightarrow M\{e \parallel \bar{s}, \bar{r}\}$

The monoid M should be primitive. If for example M was a set monoid, the result of the initial expression would be a set of sets of elements, whilst the result of the transformed one would just be a set of elements.

$$10. M\{\text{merge}[M](e_1, e_2) \parallel \bar{s}\} \Rightarrow \text{merge}[M](M\{e_1 \parallel \bar{s}\}, M\{e_2 \parallel \bar{s}\})$$

For this transformation to be valid M should be commutative and primitive (as in rules 3 and 9 respectively).

$$11. M\{e \parallel \bar{p}, u \equiv e', \bar{q}\} \Rightarrow M\{e[e'/u] \parallel \bar{p}, \bar{q}[e'/u]\}$$

$$12. M\{e \parallel \bar{q}, \text{merge}[\text{all}](p_1, \dots, p_n), \bar{r}\} \Rightarrow M\{e \parallel \bar{q}, p_1, \dots, p_n, \bar{r}\}$$

$$13. \text{zero}[\text{all}] \Rightarrow \text{true}$$

$$\text{zero}[\text{some}] \Rightarrow \text{false}$$

$$\text{zero}[\text{sum}] \Rightarrow 0$$

$$14. M\{e \parallel \bar{q}, \text{true}, \bar{r}\} \Rightarrow M\{e \parallel \bar{q}, \bar{r}\}$$

$$15. M\{e \parallel \bar{q}, \text{false}, \bar{r}\} \Rightarrow \text{zero}[M]$$

$$16. M\{e \parallel \cdot\} \Rightarrow \text{unit}[M](e), \text{ if } M \text{ is a collection monoid.}$$

$$M\{e \parallel \cdot\} \Rightarrow e, \text{ if } M \text{ is a primitive monoid.}$$

$$17. \text{merge}[M](\bar{m}, \text{zero}[M]) \Rightarrow \text{merge}[M](\bar{m})$$

$$\text{merge}[M](\bar{m}_i) \Rightarrow \bar{m}_i$$

$$\text{merge}[\text{some}](\bar{m}, \text{true}) \Rightarrow \text{true}$$

$$\text{merge}[\text{all}](\bar{m}, \text{false}) \Rightarrow \text{false}$$

$$\text{merge}[\text{some}](\bar{m}, \text{false}) \Rightarrow \text{merge}[M](\bar{m})$$

$$\text{merge}[\text{all}](\bar{m}, \text{true}) \Rightarrow \text{merge}[M](\bar{m})$$

$$18. \text{merge}[M](\bar{m}, \text{merge}[M](\bar{n})) \Rightarrow \text{merge}[M](\bar{m}, \bar{n})$$

$$19. \text{not}(\text{not } q) \Rightarrow q$$

$$20. \text{not true} \Rightarrow \text{false} \text{ and } \text{not false} \Rightarrow \text{true.}$$

$$21. \text{not } (q_1 \text{ comp_op } q_2) \Rightarrow q_1 \text{ comp_op}' q_2,$$

where $\text{comp_op}, \text{comp_op}' \in \{<, >, \leq, \geq, =, \neq\}$ and

$$\forall e_1, e_2. ((e_1 \text{ comp_op } e_2) \parallel (e_1 \text{ comp_op}' e_2) = \text{true}).$$

$$22. \text{not } (\text{merge}[\text{all}](p_1, \dots, p_n)) \Rightarrow \text{merge}[\text{some}](\text{not } p_1, \dots, \text{not } p_n)$$

$$\text{not } (\text{merge}[\text{some}](p_1, \dots, p_n)) \Rightarrow \text{merge}[\text{all}](\text{not } p_1, \dots, \text{not } p_n)$$

$$23. \text{not } (\text{all/some}\{\text{pred} \parallel \bar{r}\}) \Rightarrow \text{some/all}\{\text{not pred} \parallel \bar{r}\}$$

24. $\text{bagof}(\text{zero}[M]) \Rightarrow \text{zero}[\text{bag}]$
 $\text{bagof}(\text{unit}[M](e)) \Rightarrow \text{unit}[\text{bag}](e)$
 $\text{bagof}(\text{merge}[M](e_1, \dots, e_n)) \Rightarrow \text{merge}[\text{bag}](\text{bagof}(e_1), \dots, \text{bagof}(e_n)),$
 if $M \in \{\text{bag}, \text{list}\}$
 $\text{bagof}(M\{e \parallel \bar{r}\}) \Rightarrow \text{bag}\{e \parallel \bar{r}\},$ if $M \in \{\text{bag}, \text{list}\}$
 $\text{bagof}(e) \Rightarrow \text{bag}\{u \parallel u \leftarrow e\},$ if e has either of the forms id or $e_1.\text{id}$ or $e_1[e_2]$ and e is a list of elements
 $\text{bagof}(e) \Rightarrow e,$ if e has either of the forms id or $e_1.\text{id}$ or $e_1[e_2]$ and e is a bag of elements
 $\text{bagof}(\text{bagof}(e)) \Rightarrow \text{bagof}(e)$
 $\text{bagof}(\text{listof}(e)) \Rightarrow \text{bagof}(e)$
 $\text{bagof}(\text{arrayof}(e)) \Rightarrow \text{bagof}(e)$
 $\text{bagof}(\text{orderof}(e)) \Rightarrow \text{bagof}(e)$
25. $M\{e \parallel \bar{q}, u \leftarrow \text{bagof}(\text{merge}[N](e_1, e_2)), \bar{s}\} \Rightarrow M\{e \parallel \bar{q}, u \leftarrow \text{merge}[N](e_1, e_2), \bar{s}\},$ if M is idempotent and N is either set or oset .
 $M\{e \parallel \bar{q}, u \leftarrow \text{bagof}(N\{e' \parallel \bar{r}\}), \bar{s}\} \Rightarrow M\{e \parallel \bar{q}, u \leftarrow N\{e' \parallel \bar{r}\}, \bar{s}\},$ if M is idempotent and N is either set or oset .
 $M\{e \parallel \bar{q}, u \leftarrow \text{bagof}(e'), \bar{s}\} \Rightarrow M\{e \parallel \bar{q}, u \leftarrow e', \bar{s}\},$ if M is idempotent and e' has either of the forms id or $e_1.\text{id}$ or $e_1[e_2]$ and e' is a set or oset of elements .
26. $\text{listof}(\text{zero}[M]) \Rightarrow \text{zero}[\text{list}]$
 $\text{listof}(\text{unit}[M](e)) \Rightarrow \text{unit}[\text{list}](e)$
 $\text{listof}(\text{merge}[\text{list}](e_1, \dots, e_n)) \Rightarrow \text{merge}[\text{list}](e_1, \dots, e_n)$
 $\text{listof}(\text{list}\{e \parallel \bar{r}\}) \Rightarrow \text{list}\{e \parallel \bar{r}\}$
 $\text{listof}(e) \Rightarrow e,$ if e has either of the forms id or $e_1.\text{id}$ or $e_1[e_2]$ and e is a list of elements
 $\text{listof}(\text{listof}(e)) \Rightarrow \text{listof}(e)$
 $\text{listof}(\text{arrayof}(e)) \Rightarrow \text{listof}(e)$
 $\text{listof}(\text{orderof}(e)) \Rightarrow \text{orderof}(e)$
27. $M\{e \parallel \bar{s}, u \leftarrow \text{listof}(e'), \bar{r}\} \Rightarrow M\{e \parallel \bar{s}, u \leftarrow \text{bagof}(e'), \bar{r}\},$ if $M \in \{\text{set}, \text{bag}, \text{group}\}$ (monoid degeneration)
 $M\{e \parallel \bar{s}, u \leftarrow \text{orderof}(e'), \bar{r}\} \Rightarrow M\{e \parallel \bar{s}, u \leftarrow \text{bagof}(e'), \bar{r}\},$ if $M \in \{\text{set}, \text{bag}, \text{group}\}$ (monoid degeneration)
28. $M\{e \parallel \bar{p}, x \leftarrow \text{bagof}(\text{group}\{\text{struct}(\ell_1 = e_1, \dots, \ell_n = e_n, \text{partition} = e') \parallel \bar{q}\}), \bar{r}, y \leftarrow x.\text{partition}, \bar{s}\} \Rightarrow M\{e \parallel \bar{p}, \bar{q}, \bar{r}, y \leftarrow e', \bar{s}\},$ if x, ℓ_1, \dots, ℓ_n are not free variables in e, r, s and M is commutative, that is $M \in \{\text{set}, \text{bag}, \text{group}\}$.

29. $M\{e \parallel \bar{p}, x \leftarrow \text{bagof}(\text{merge}[\text{group}](\text{group}\{\text{struct}(\ell_1 = e_1, \dots, \ell_n = e_n, \text{partition} = e') \parallel \bar{q}_1, u \leftarrow q_{21}, \bar{q}_3\}, \text{group}\{\text{struct}(\ell_1 = e_1, \dots, \ell_n = e_n, \text{partition} = e') \parallel \bar{q}_1, u \leftarrow q_{22}, \bar{q}_3\}))\}, \bar{r}, y \leftarrow x.\text{partition}, \bar{s}\} \Rightarrow M\{e \parallel \bar{p}, \bar{q}_1, u \leftarrow \text{merge}[\text{bag}](q_{21}, q_{22}), \bar{q}_3, \bar{r}, y \leftarrow e', \bar{s}\},$
 if x, ℓ_1, \dots, ℓ_n are not free variables in e, r, s and M is commutative, i.e. $M \in \{\text{set}, \text{bag}, \text{group}\}$.
30. $\text{elementof}(\text{unit}[M](e)) = e$
 $\text{elementof}(\text{set}/\text{bag}/\text{list}/\text{array}(e)) = e$

The rules presented above contain eight of the ten transformations proposed by Fegaras and Maier (1,2,3,4,6,7,9,11). One of the omitted rules concerns the beta reduction: $(\lambda u.e_1)e \Rightarrow e_1[e/u]$. In this context, method application is denoted by using the name of a method followed by an optional set of arguments. The second transformation that has been excluded from the normalization rules is:

$$M\{e \parallel \bar{q}, u \leftarrow (\text{if } e_1 \text{ then } e_2 \text{ else } e_3), \bar{s}\} \Rightarrow \\ \text{merge}[M](M\{e \parallel \bar{q}, e_1, u \leftarrow e_2, \bar{s}\}, M\{e \parallel \bar{q}, \text{not}(e_1), u \leftarrow e_3, \bar{s}\})$$

This rule is valid if M is commutative or \bar{q} is empty. There is no equivalent of this rule, since ODMG OQL statements do not use or imply the *if* statement.

4.3.2 Examples of transforming calculus expressions

We give an example that applies the normalization rules for the optimization of a query with an order by clause. Assume that the initial OQL query (Q_1) is:

```
order_flatten
(select q
 from q1 as x1, ..., qk as xk
 where q'
 order by q''1, ..., q''n)
```

This would be translated to the following monoid calculus form:

$$\text{list}\{u \parallel v_1 \leftarrow \text{list}\{v_2.\text{partition} | v_2 \leftarrow \text{orderof}(\text{group}\{\text{struct}(\text{asc}_1 = |q_1|'', \dots, \text{asc}_n = |q_n|'', \text{partition} = \text{unit}[\text{bag}](|q|'')) \parallel x_1 \leftarrow \text{bagof}(|q_1|'', \dots, x_k \leftarrow \text{bagof}(|q_k|'', |q'|'')\}), u \leftarrow \text{listof}(v_1)\} \quad (4.2)$$

Based on transformations 6 and 11, expression 4.2 is converted to the following form:

$$\begin{aligned}
& \text{list}\{u \parallel v_2 \leftarrow \text{orderof}(\text{group}\{\text{struct}(\text{asc}_1 = |q_1''|^*, \dots, \\
& \quad \text{asc}_n = |q_n''|^*, \text{partition} = \text{unit}[\text{bag}](|q|^*) \parallel x_1 \leftarrow \text{bagof}(|q_1|^*), \\
& \quad \dots, x_k \leftarrow \text{bagof}(|q_k|^*), |q'|^*\}), u \leftarrow \text{listof}(v_2.\text{partition})\} \quad (4.3)
\end{aligned}$$

Consider the query Q_2 :

```

select x
from Q1 as x

```

Based on 4.3, Q_2 takes the following calculus form:

$$\begin{aligned}
& \text{bag}\{x \parallel x \leftarrow \text{bagof}(\text{list}\{u \parallel v_2 \leftarrow \text{orderof}(\text{group}\{\text{struct}(\text{asc}_1 = |q_1''|^*, \\
& \quad \dots, \text{asc}_n = |q_n''|^*, \text{partition} = \text{unit}[\text{bag}](|q|^*) \parallel x_1 \leftarrow \text{bagof}(|q_1|^*), \\
& \quad \dots, x_k \leftarrow \text{bagof}(|q_k|^*), |q'|^*\}), u \leftarrow \text{listof}(v_2.\text{partition})\})\} \quad (4.4)
\end{aligned}$$

After applying the transformations 6, 11 and 24, expression 4.4 becomes:

$$\begin{aligned}
& \text{bag}\{u \parallel v_2 \leftarrow \text{orderof}(\text{group}\{\text{struct}(\text{asc}_1 = |q_1''|^*, \dots, \\
& \quad \text{asc}_n = |q_n''|^*, \text{partition} = \text{unit}[\text{bag}](|q|^*) \parallel x_1 \leftarrow \text{bagof}(|q_1|^*), \\
& \quad \dots, x_k \leftarrow \text{bagof}(|q_k|^*), |q'|^*\}), u \leftarrow \text{listof}(v_2.\text{partition})\} \quad (4.5)
\end{aligned}$$

Expression 4.5 is transformed through 27 and 24 into 4.6:

$$\begin{aligned}
& \text{bag}\{u \parallel v_2 \leftarrow \text{bagof}(\text{group}\{\text{struct}(\text{asc}_1 = |q_1''|^*, \dots, \\
& \quad \text{asc}_n = |q_n''|^*, \text{partition} = \text{unit}[\text{bag}](|q|^*) \parallel x_1 \leftarrow \text{bagof}(|q_1|^*), \\
& \quad \dots, x_k \leftarrow \text{bagof}(|q_k|^*), |q'|^*\}), u \leftarrow v_2.\text{partition}\} \quad (4.6)
\end{aligned}$$

This is the point at which transformation 28 is useful. In a commutative monoid comprehension, if a group monoid forms partitions that are used only to be unnested in later generators, then `group` may be replaced by `bag`. Grouping is avoided in the first place. Thus, expression 4.6 is transformed to the simplified form:

$$\text{bag}\{u \parallel x_1 \leftarrow \text{bagof}(|q_1|^*), \dots, x_k \leftarrow \text{bagof}(|q_k|^*), |q'|^*, u \leftarrow \text{unit}[\text{bag}](|q|^*)\} \quad (4.7)$$

Rules 2 and 11 finally enable the transformation of 4.7 into the final form 4.8.

$$\text{bag}\{|q|^* \parallel x_1 \leftarrow \text{bagof}(|q_1|^*), \dots, x_k \leftarrow \text{bagof}(|q_k|^*), |q'|^*\} \quad (4.8)$$

The normalized calculus form (expression 4.8) of query Q_2 is far simpler than the initial expression 4.4.

```

case Query of
    ... : ...
    not e : normalize_not
    ... : ...
    M{head || q1, ..., qn} : normalize_compr
    ... : ...

```

Figure 4.3: Normalization Algorithm

4.3.3 Normalization Algorithm

In order to transform a calculus expression into a simpler form, it is not sufficient to know the rules that guide the transformations. It is important to determine the order in which the rules must be triggered. The rules that are eligible to be applied depend on the structure of the query. The normalization algorithm consists of a single *case* step that considers all possible constructs of OQL (see figure 4.3).

Figure 4.4 presents the normalization algorithm for the construct `not (e)`.

Most of the other OQL constructs are normalized in a similar way. Only the case of monoid comprehension expressions is slightly more complicated and therefore needs to be explained in detail (see figure 4.5).

`optimize_from_qualifier(i)` optimizes the comprehension starting from qualifier q_i to the last qualifier and then to the head of the comprehension. The transformations applicable to each qualifier depend on whether the latter is a filter, a generator, or a binding. We study these cases separately:

- Suppose qualifier q_i is a generator $u \leftarrow e$. The rules that may be applied are 1, 2, 3, 6, 25, 27, 28 and 29. The algorithm that normalizes a comprehension from its generator in position i is presented in 4.6.
- Suppose qualifier q_i is a filter e . The rules that may be applied are 7, 8, 12, 14 and 15. The algorithm that normalizes a comprehension from its filter in position i is presented in figure 4.7.
- Suppose qualifier q_i is a binding $u \equiv e$. The only applicable rule is 11. The algorithm that normalizes a comprehension from its binding in position i is given in figure 4.8.

When all qualifiers have been considered the algorithm invokes a process to optimize the head of the monoid comprehension. The applicable rules in this case are 9, 10 and 16. Details of the algorithm are given in figure 4.9.

```

normalize_not (Rules 19,20,21,22,23)

Input: not(e)
Algorithm

opt_e = optimize e
case opt_e has the construct
    not pred : apply rule 19
    true,false : apply rule 20
    e1 < comp_oper > e2 : apply rule 21
    merge[all/some](p1,...,pn) : apply rule 22
    all/some{pred || ...} : apply rule 23

if any rule was applied
    optimize the result of the transformation
    return the result of the optimization
else
    annotate not (opt_e) as normalized
    return not (opt_e)

```

Figure 4.4: Normalization of the OQL construct not (e)

```

normalize compr
(Rules 1,2,3,6,7,8,9,10,11,12,14,15,16,25,27,28,29)

Input: M{head || q1,...,qn}
Algorithm

if the comprehension has qualifiers (n > 0)
    optimize_from_qualifier(1)
else
    optimize_head
annotate result as normalized and return it

```

Figure 4.5: Normalization of monoid comprehensions


```

if e is annotated as normalized
  opt_e = e
else
  opt_e = optimize e
case opt_e has the construct
  zero[N] : apply rule 1;
  unit[N](e1) : apply rule 2;
               optimize_from_qualifier(i);
  merge[N](e1, e2) : try rule 3
                       if rule 3 was applicable
                       optimize its result;
  N{e' || r̄} : : apply rule 6;
               optimize_from_qualifier(i);
  bagof(merge[N](e1, ..., en)) : try rule 25;
                                   if rule 25 was applicable
                                   optimize_from_qualifier(i);
  bagof(set/oset{e' || r̄}) : try rule 25;
                             if rule 25 was applicable
                             optimize_from_qualifier(i);
  listof/orderof(e) : try rule 27;
                      if rule 27 was applicable
                      optimize_from_qualifier(i);
  bagof(group{...}) : try rule 28;
                     if rule 28 was applicable
                     optimize_from_qualifier(i);
  bagof(merge[group](...)) : try rule 29;
                              if rule 29 was applicable
                              optimize_from_qualifier(i);
if no transformation was applied in the case clause
  optimize_from_qualifier(i+1)
(or optimize head if i+1 > n)

```

Figure 4.6: Normalizing comprehensions from qualifier q_i , where q_i is a generator ($u \leftarrow e$)

```

if e is annotated as normalized
  opt_e = e
else
  opt_e = optimize e
case opt_e has the construct
  true : apply rule 14;
        optimize_from_qualifier(i);
  false : apply rule 15; optimize result;
  some{...|...} : try rule 7;
                 if rule 7 was applicable
                 optimize_from_qualifier(i);
  merge[some](...) : try rule 8;
                    if rule 8 was applicable
                    optimize its result;
  merge[all](...) : apply rule 12;
                   optimize_from_qualifier(i);
if no transformation was applied in the case clause
  optimize_from_qualifier(i+1)
(or optimize_head if i+1 > n)

```

Figure 4.7: Normalizing comprehensions from qualifier q_i , where q_i is a filter (e)

```

if e is annotated as normalized
  opt_e = e
else
  opt_e = optimize e
apply rule 11
optimize_from_qualifier(i)

```

Figure 4.8: Normalizing comprehensions from qualifier q_i , where q_i is a binding ($u \equiv e$)

```

if e is annotated as normalized,
  opt_head = head
else
  opt_head = optimize head
if opt_head has the form M{...} (where M is a primitive monoid)
  apply rule 9
else if opt_head has the form merge[M](...)
  (where M is a primitive and commutative monoid)
  apply rule 10
else if there are no qualifiers
  apply rule 16
if any transformation was applied
  optimize its result
  return optimized result
else
  return the comprehension as optimized so far

```

Figure 4.9: Normalizing the head of a comprehension

The normalization algorithm has the following properties: i) it always terminates; ii) it yields the same result as the exhaustive iterative application of rules and iii) it yields this result in the most efficient way.

We first provide an informal proof that the algorithm always terminates. Normalization rules perform either of the following tasks.

- They push up merge operations (Rules 3, 8, 10, 18 and 22). Neither of the remaining rules pushes down nor creates new merge operators. Hence, rules 3, 8, 10, 18 and 22 are applicable only a finite number of times.
- They unnest inner comprehensions pushing up their head and qualifiers (Rules 6, 7 and 9). The number of inner comprehensions can only be increased finitely by either of the following rules

$$1. M\{e \parallel \bar{q}, \text{merge}[\text{some}](p_1, p_2), \bar{s}\} \Rightarrow \\ \text{merge}[M](M\{e \parallel \bar{q}, p_1, \bar{s}\}, M\{e \parallel \bar{q}, p_2, \bar{s}\})$$

$$M\{e \parallel \bar{q}, u \leftarrow \text{merge}[N](e_1, e_2), \bar{s}\} \Rightarrow \\ \text{merge}[M](M\{e \parallel \bar{q}, u \leftarrow e_1, \bar{s}\}, M\{e \parallel \bar{q}, u \leftarrow e_2, \bar{s}\})$$

$$M\{\text{merge}[M](e_1, e_2) \parallel \bar{s}\} \Rightarrow \text{merge}[M](M\{e_1 \parallel \bar{s}\}, M\{e_2 \parallel \bar{s}\})$$

In this case, the number of additional comprehensions is bounded by the initial number of merge operators.

2. $\text{bagof}(e) \Rightarrow \text{bag}\{u \parallel u \leftarrow e\}$, if e has either of the forms id or $e_1.\text{id}$ or $e_1[e_2]$ and e is a list of elements. Based on this rule, the number of comprehensions can be increased at most by the initial number of nonhomomorphic functions that could be transformed to bagof .
 3. $M\{e \parallel \bar{p}, u \equiv e', \bar{q}\} \Rightarrow M\{e[e'/u] \parallel \bar{p}, \bar{q}[e'/u]\}$, where e is a comprehension. In this case, the number of comprehensions can be increased at most by the number of occurrences of range variables in the tree. This is finite since the number of range variables is increased only when pushing up merge operators.
- They eliminate, push down or simplify nonhomomorphic functions (Rules 24, 25, 26, 27 and 30). These rules are applicable a finite number of times, since the number of nonhomomorphic functions is finite and the depth of the tree increases finitely using rules 3, 8 and 11.
 - They avoid unnecessary grouping by pushing up the qualifiers of nested group comprehensions (Rules 28 and 29).
 - They push down not operators (Rules 19, 20, 21, 22 and 23).
 - They transform calculus expressions to much simpler forms (Rules 1, 2, 4, 5, 12, 13, 14, 15, 16, 17), e.g. by pushing down the zero element (Rule 1).

There is no rule that performing one of these actions causes the reverse effect of another. Hence, since each one can be applied a finite number of times, the algorithm will eventually not be able to apply any of the rules and will terminate. The same applies for the exhaustive iterative application of rules.

We show that our algorithm is also exhaustive in the sense that when it terminates no normalization rule is applicable in any of the subexpressions. The algorithm can be viewed as a Depth-First unnesting algorithm where the children of a node are first normalized, before the node itself is considered. In a monoid comprehension, each child (qualifier) is normalized as an independent expression first, and then with respect to the outer comprehension (parent node). If normalization of a qualifier with respect to the outer comprehension alters only the qualifier, the qualifier is renormalized. If no transformation takes place, the following qualifier is visited. Finally, if the normalization process results in a completely different outer expression, which is no longer a comprehension, the new expression is normalized from the beginning. Hence, the next qualifiers are considered after the current

one, or all qualifiers are visited as subexpressions of a different calculus expression. In this way, all subexpressions are eventually normalized.

The question that arises is whether the actual order of applying rules affects the result of the exhaustive iterative algorithm. If we proved that the exhaustive application of rules (in a random iterative way) is confluent, we could deduce that our algorithm also yields the same result as the random iterative process. The reason is that since our algorithm is also exhaustive, it could be viewed as one of the possible executions of the random iterative process.

All rules (except for Rule 11²) modify at most the first two levels of a subtree (subexpression). Hence, if two transformations do not occur close to each other, the order in which they are applied does not affect the final result. However, there are two cases in which applying one rule might affect the result of the other, or might even hide the opportunity to apply the second rule.

1. If the possible transformations occur in two consecutive edges in a path from the root to a leaf.
2. If the possible transformations occur in two sibling edges, i.e. in two edges starting from the same node.

To prove confluence, we need to consider applying all possible pairs of rules in each one of the cases above. For space reasons, we only study the application of one pair in each case.

1. Consider the calculus expression $M\{e \parallel \bar{p}, u \leftarrow \text{merge}[N](q, \text{merge}[N](r, s)), \bar{t}\}$. If M is commutative then both rules 3 and 18 are applicable in two consecutive edges of the calculus tree.

- If we applied rule 3 first then the expression would be transformed to:

$$\text{merge}[M](M\{e \parallel \bar{p}, u \leftarrow q, \bar{t}\}, M\{e \parallel \bar{p}, u \leftarrow \text{merge}[N](r, s), \bar{t}\}).$$

Another application of rule 3 would then yield the expression:

$$\text{merge}[M](M\{e \parallel \bar{p}, u \leftarrow q, \bar{t}\}, \text{merge}[M](M\{e \parallel \bar{p}, u \leftarrow r, \bar{t}\}, M\{e \parallel \bar{p}, u \leftarrow s, \bar{t}\})).$$

By applying rule 12, we would get the final result:

$$\text{merge}[M](M\{e \parallel \bar{p}, u \leftarrow q, \bar{t}\}, M\{e \parallel \bar{p}, u \leftarrow r, \bar{t}\}, M\{e \parallel \bar{p}, u \leftarrow s, \bar{t}\}).$$

- If we applied rule 18 first then the expression would be transformed to the following form:

$$M\{e \parallel \bar{p}, u \leftarrow \text{merge}[N](q, r, s), \bar{t}\}.$$

Using rule 3 we would then get the same final result:

$$\text{merge}[M](M\{e \parallel \bar{p}, u \leftarrow q, \bar{t}\}, M\{e \parallel \bar{p}, u \leftarrow r, \bar{t}\}, M\{e \parallel \bar{p}, u \leftarrow s, \bar{t}\}).$$

2. Consider the calculus expression $M\{e \parallel \bar{p}, u \leftarrow \text{merge}[N](q, r), v \leftarrow \text{zero}[M]\}$. In this case we can apply either of the rules 1 or 3 on two sibling edges.

²Rule 11 is just a replacement of a variable by a subexpression and the order in which it is applied should not matter as long as the order in which the remaining rules are applied does not alter the final result.

- If rule 1 was applied first then the result would be $\text{zero}[M]$.
- If rule 3 was applied first then the expression would take the form:
 $\text{merge}[M](M\{e \parallel \bar{p}, u \leftarrow q, v \leftarrow \text{zero}[M]\}, M\{e \parallel \bar{p}, u \leftarrow r, v \leftarrow \text{zero}[M]\})$.

If any of the possible rule sequences $[1, 1, 17, 17]$ or $[1, 17, 1, 17]$ were then applied the result would be $\text{zero}[M]$.

Since the result does not depend on the order in which we apply rules on consecutive, sibling or isolated edges, the random iterative application of rules is confluent and yields the same result as our algorithm. Confluence is proven with the aid of a tree representation of rule applications.

Assume that $S_i = [r_{i1}, \dots, r_{in_i}]$ represents a sequence of rule applications on a calculus expression e , such that, after it is used, no other rule is applicable. If $S = \{S_1, \dots, S_m\}$ is the set of all possible sequences applicable to e , S could be represented as a tree. The head of the tree is the expression e , the leaves represent the final forms of e and each sequence S_i is a complete path from the head (e) to the i^{th} final form of e . The internal nodes represent intermediate forms of e , whilst each edge can be viewed as a rule application.

We navigate the tree upwards starting from internal nodes e_{1i} , $i = 1, \dots, k_1$, whose direct or indirect children are either leaves or nodes with one child. All leaves under such an intermediate expression e_{1i} should have the same form; the reason is that no matter which rule was applied first to e_{1i} , we maintain the possibility of getting the same result (proof about consecutive, sibling or isolated transformations). We then consider internal nodes e_{2i} , $i = 1, \dots, k_2$, whose children are either leaves or internal nodes with one child or nodes in $\{e_{11}, \dots, e_{1k_1}\}$. In the same way as before, we prove that the leaves of each one of these nodes e_{2i} have the same form. The same procedure continues until we reach the root of the tree, which also has the same property. Hence, all possible exhaustive sequences of rule applications finally yield the same result. We conclude that the random iterative process is confluent and yields the same result as our algorithm.

When a calculus expressions is normalized (using our algorithm), it is annotated as such, and needs not be considered again. In addition to that, we consider normalization of monoid comprehensions, every time we visit a subexpression (qualifier) and thus we identify normalization opportunities as early as possible. Another approach would be to normalize all qualifiers first, before proceeding to consider the outer comprehension. However, this might burden the normalization process with redundant steps. By annotating subexpressions as normalized and unnesting expressions as early as possible, the algorithm normalizes a calculus expression in the least number of steps. Our algorithm triggers the execution of one of the optimal sequences of rules that could possibly be executed based on a random exhaustive process.

In this section, we provided a complete framework for normalizing OQL queries. This framework has been implemented and has been tested for a series of complex OQL queries.

4.4 Discussion

This chapter gave an overview of the first intermediate representation of OQL, the monoid calculus. In particular, it presented how OQL queries are translated into calculus expressions, and how the latter are transformed into more efficient forms. This framework is based on work by Afshar [Afs98] and Fegaras and Maier [FM00]. However, it has extended the latter with new monoids, several different translation rules, new normalization rules, and an algorithm that determines the order of applying these rules. It has also pointed out several weak points of the ODMG standard and has proposed ways of resolving them.

Chapter 5

OQL algebraic representation

The previous chapter described how OQL queries are translated into the monoid calculus representation; it also presented several transformations within this representation. This chapter presents the second intermediate representation of OQL queries, referred to as the object algebra. We do not use monoid homomorphisms as the basis of this algebra, since they are difficult to optimize and their natural implementation involves inefficient nested loops. Instead, we adopt a set of algebraic operators for two reasons: i) these operators are easy to rearrange and ii) they can be mapped into efficient physical operators used broadly in the relational systems.

The algebraic operators are similar to those proposed by Fegaras and Maier [FM00]. However, the mechanism for translating queries from calculus to algebraic form is significantly different. This chapter presents the translation algorithm in detail, highlights its advantages and gives several examples that demonstrate its functionality.

5.1 Algebraic Operators

We give an overview of the algebraic operators and explain their functionality. First, we present those operators that have also been defined in [FM00] and point out subtle differences in the two approaches. We then define a set of algebraic operators that are novel in this dissertation and explain their role in the translation process. The complete syntax of the algebraic representation of OQL is given in Figure 5.1. In particular, Group A (figure 5.1) contains the expressions that have a very similar calculus equivalent, whilst Group B consists of the algebraic operators described below.

A subtle point in which the current algebraic operators differ from the corresponding ones in [FM00] is the specification of the additional information mon (monoid). In general, operators iterate over the elements of a collection, process them in some way, and accumulate the results. The monoid field specifies the kind of merge operation used to accumulate the results. Fegaras and Maier specify

```

e ::= /* Group A */
    < literal > (integer, float, boolean, char, string)
    < id > (db object/extent name)
    zero[< monoid >]
    unit[< monoid >](e)
    merge[< prim_mon >](e1, ..., en)
    struct(l1 = e1, ..., lk = ek)
    class_name(l1 = e1, ..., lk = ek)
    e1 < comp_op > e2
    e1 < arithm_op > e2
    e1.attr_rltp_method
    e(e1, ..., en)
    e1[e2]
    not (e)
    -e

    /* Group B */
    mon projection
     $\pi$  operand

    operand1  $\overset{\text{mon}}{\bowtie}$  operand2
                 $\overline{v_1, v_2}$  predicate

    operand1  $\overset{\text{mon}_1 \text{ mon}_2}{\Rightarrow \bowtie}$  operand2
                 $\overline{v_1, v_2}$  predicate

     $\overset{\text{mon}}{\Gamma}$   $\overset{\text{gmon}}{\text{gname}}$   $\overset{\text{gproject}}{\text{g}\overline{v}}$ 
     $\overline{v}$  predicate operand

     $\overset{\text{mon}}{\sigma}$   $\overline{v}$  predicate operand |  $\overset{\text{mon}_1 \text{ mon}_2}{\overline{\sigma}}$   $\overline{v}$  predicate operand
     $\overset{\text{mon}}{\mu}$   $\overline{v}$  path pred operand |  $\overset{\text{mon}_1 \text{ mon}_2}{\overline{\mu}}$   $\overline{v}$  path pred operand

     $\overset{\text{mon}}{M}$   $\overset{\text{mon}'}{(e_1, \dots, e_n)}$  |  $\overset{\text{mon}}{\overline{M}}$   $\overset{\text{mon}'}{\text{parts}}$   $\overline{v}$  operand
     $\overset{\text{mon}}{B}$   $\overline{v}$  function operand |  $\overset{\text{mon}}{\overline{B}}$   $\overline{v}$  function argument operand

```

Figure 5.1: Syntax of the algebraic OQL representation

this information (accumulator) only for the Reduce and Nest operators. We claim that it would be equally useful for the other operators.

First, it would enable the optimization of the execution plan, as will be explained later on. Second, it may reduce the number of operators in the algebraic plan. For instance, to filter the elements of a bag into a set, based on a predicate, one would need a Select operator to apply the predicate, and a Reduce operator to transform the bag into a set. In our context, it suffices to use a Select operator with monoid (set).

In the remainder of this section we present the operators in detail.

1. $\overset{\text{mon}}{\underset{\bar{v}}{\pi}} \begin{array}{l} \text{projection} \\ \text{operand} \end{array}$

Project : This operator iterates over the elements of operand. It binds (the parts of) each element to the variables \bar{v} and calculates the expression *projection*. It corresponds to the following calculus expression:

$$\text{mon}\{\text{projection} \mid u \leftarrow \text{operand}, v_1 \equiv u.v_1, \dots, v_n \equiv u.v_n\}$$

where $\bar{v} = \{v_1, \dots, v_n\}$

The Project operator differs from the similar Reduce operator ([FM00]) in that it considers no predicate. That is the Project operator does not perform any filtering on the elements of operand.

Hereafter the notation $\bar{v} \leftarrow \text{operand}$ is used as a short form for the binding:

$$(u \leftarrow \text{operand}, v_1 \equiv u.v_1, \dots, v_n \equiv u.v_n)$$

2. $\overset{\text{mon}}{\underset{\bar{v}}{\sigma}} \begin{array}{l} \text{predicate} \\ \text{operand} \end{array}$

Select : This operator iterates over the operand and filters its elements using *predicate*. Those elements of operand that satisfy *predicate* are accumulated using the merge operation of monoid *mon*. The operator corresponds to the calculus expression:

$$\text{mon}\{\text{struct}(v_1 = v_1, \dots, v_n = v_n) \mid \bar{v} \leftarrow \text{operand}, \text{predicate}\}$$

where $\bar{v} = \{v_1, \dots, v_n\}$.

The question that arises at this point is how to identify the monoid *mon*. In the case of Project the monoid is defined from the calculus expression that we translate, in a straightforward way. The same holds for the Nest operator which is defined later on. The following example demonstrates how the monoid information is derived for all other operators, including Select.

Say that the query to be executed is

$$\text{exists } j \text{ in Priority_List : } j.\text{priority} = 6$$

where `Priority_List` is an ordered collection (list) containing objects of class `Job`. This query would be translated into the calculus form

$$\text{some}\{\text{true} | j \leftarrow \text{Priority_List}, j.\text{priority} = 6\}$$

One would expect the algebraic equivalent of this expression to include a `Select` operator that tests the predicate `j.priority = 6` iterating over the collection `Priority_List`. If this operator has no knowledge about the monoid information, it respects the ordered nature of the operand and therefore it filters its elements in the prespecified order. However, knowing that the result monoid of the comprehension (`some`) is a commutative monoid, it may consider `Priority_List` as a bag, and rearrange its elements or make use of indices in an optimal way for the filtering. In fact, since `some` is an idempotent monoid, the operator could even decide to remove duplicates from `Priority_List` before filtering its elements. Thus, the physical implementation of `Select` could benefit from the information `mon = set/bag`. Whilst optimizing the algebraic plan, it is possible to refine the monoid information of several operators, e.g. from `set/bag` to `set`, instead of leaving this decision to the physical implementation.

The example illustrated how the monoid information of `Select` is derived by looking at the result monoid of the monoid comprehension. The following table presents this relation in detail.

Comprehension Monoid	Operator Monoid
< IC monoid >	set/bag
< C monoid >	bag
< I monoid >	oset/list
< other >	list

where,

- IC is a notation for monoids both idempotent and commutative.
- C stands for idempotent but not commutative monoids.
- I stands for commutative but not idempotent monoids.
- other is the category for monoids that are neither commutative nor idempotent.

The table above is used to derive the monoid information for all algebraic operators, excluding `Project` and `Nest`.

$$3. \prod_{\bar{v}, v'}^{\text{mon}} \begin{array}{l} \text{path} \\ \text{pred} \\ \text{oper} \end{array}$$

Unnest : This operator iterates over the operand `oper` and for each element, it *unfolds* the path expression. Each element of the path expression is bound to the new variable v' . The name of this variable should be different from the names of \bar{v} . Hence the result of this operation is a collection of structures with labels named as the variables in \bar{v} or as v' . In fact, only the structures that satisfy `pred` are included in the result. The monoid information (`mon`) is inferred in the same way as in the operator `Select`.

Notice that it is not necessary to include the predicate in the `Unnest` operator. Instead one could use a `Select` operator having `Unnest` as its operand. However, it is a good decision to include a predicate, since `Unnest` often generates a much bigger collection than its operand. Hence, it is worthwhile to prune this collection before sending it to the output. The same applies for the `Join` operator.

The `Unnest` operator corresponds to the calculus expression:

$$\text{mon}\{\text{struct}(v_1 = v_1, \dots, v_n = v_n, v' = v') | \bar{v} \leftarrow \text{oper}, v' \leftarrow \text{path}, \text{pred}\}$$

where $\bar{v} = \{v_1, \dots, v_n\}$.

$$4. \text{operand}_1 \begin{array}{c} \text{mon} \\ \boxtimes \\ \bar{v}_1, \bar{v}_2 \text{ predicate} \end{array} \text{operand}_2$$

Join : This operator takes the cartesian product of the two operands and selects those elements that satisfy `predicate`. It accumulates the results using the merge operation of the monoid `mon`. The latter is derived as in the case of the `Select` operator. The physical implementation of `Join` may benefit from the `mon` information. Indeed, if `mon` = `set`/`bag`, then even if the operands are both lists, they can be considered as bags or sets. Should they be considered as bags, they can be rearranged (sorted) to facilitate their joining; if they have many duplicates, it is even more efficient to convert them to sets before joining them. In the latter case `mon` should be refined to `set`. If `mon` = `list`, the optimizer can use meta-information regarding the ordering of elements in the operands in order to apply optimization techniques studied in sequence query processing [SLR94].

The `Join` operator corresponds to the calculus expression:

$$\text{mon}\{\text{struct}(v_{11} = v_{11}, \dots, v_{1n} = v_{1n}, v_{21} = v_{21}, \dots, v_{2m} = v_{2m}) | \bar{v}_1 \leftarrow \text{operand}_1, \bar{v}_2 \leftarrow \text{operand}_2, \text{predicate}\}$$

where $\bar{v}_1 = \{v_{11}, \dots, v_{1n}\}$ and $\bar{v}_2 = \{v_{21}, \dots, v_{2m}\}$.

$$5. \quad \begin{array}{c} \text{mon}_1 \text{ mon}_2 \quad \text{path} \\ \text{pred} \\ \text{operand} \\ \hline \overline{v} \quad v' \end{array}$$

OuterUnnest : This operator is used to *unfold* path expressions in nested (inner) comprehensions. Its functionality is similar to **Unnest** with a single important difference. If for an element e in **operand**, the collection path is empty, or none of its elements satisfy **pred**, then e is still included in the result, paired with the value **NULL** [FM00]. There should be no confusion between this **NULL** value and the **UNDEFINED** value; the latter is defined in the ODMG Standard ([CBB⁺00],chapter 4.6) as the result of accessing a property of the `nil` object. **UNDEFINED** is a valid value for any literal or object type.

Notice that **OuterUnnest** is annotated with information about two monoids. The first monoid mon_1 denotes the merge operation used to accumulate the results. The second one mon_2 gives information about how **path** should be unnested. For instance, if $\text{mon}_1 = \text{bag}$ and $\text{mon}_2 = \text{list}$, then although the result is a bag, the elements of collection **path** should be annotated with a serial number, based on their order. Hence, if we need to group these elements later on into a different list, we would be able to know their order despite their being members of a bag. Since $\text{mon}_1 = \text{bag}$, the elements of **operand** should be annotated with information that distinguishes its duplicates (see discussion about dot equality in **Nest**).

OuterUnnest is used to translate path expressions that are generators in nested (inner) comprehensions. Hence, mon_1 is derived based on the monoid of the current (outer) comprehension, whilst mon_2 is inferred by the monoid of the nested (inner) comprehension.

$$6. \quad \begin{array}{c} \text{operand}_1 \quad \begin{array}{c} \text{mon}_1 \text{ mon}_2 \\ \text{predicate} \end{array} \quad \text{operand}_2 \\ \hline \overline{v}_1, \overline{v}_2 \end{array}$$

OuterJoin : This operator is useful in the case of nested comprehensions. It differs from **Join** in that if an element from **operand₁** cannot be matched with at least one element from **operand₂**, it is preserved in the result by being paired with the **NULL** value.

Notice that **OuterJoin** is annotated with information about two monoids. The first monoid mon_1 denotes the merge operation used to accumulate the results. The second one mon_2 gives information about the nature of **operand₂**. They are derived in a similar way as in the case of **OuterUnnest**.

$$7. \quad \begin{array}{c} \text{mon} \quad \begin{array}{c} \text{gmon} \\ \text{gname} \end{array} \quad \begin{array}{c} \text{gproject} \\ \text{gv} \end{array} \\ \text{predicate} \quad \text{operand} \\ \hline \overline{v} \end{array}$$

Nest : This operator groups elements based on certain grouping variables (\overline{gv}), or to be accurate, on dot equality between the values of these variables [FM00]. This operator handles the translation of inner monoid comprehensions. We give an example that shows the importance of dot equality in the context of the **Nest** operator.

Consider the translation of the following monoid calculus expression:

$$\text{bag}\{x|x \leftarrow X, \text{all}\{y.a|y \leftarrow Y, p(x, y)\}\}$$

The translation algorithm from calculus to algebra first creates an `OuterJoin` operator, to join X and Y based on $p(x, y)$. The reason why it would not use a `Join` operation is that Y is in an inner comprehension wrt X . For each x in X the algorithm then gathers (through the monoid `all`) all the values $y.a$. This is done using the `Nest` operation with:

- `mon` (monoid): `bag`
- `operand`: $X \xrightarrow[\text{x,y p(x,y)}]{\text{bag set/bag}} Y$
- `gmon` (group monoid): `all`
- `gname` (group name): `new_var`
- `gproject` (group project): `y.a`
- $\overline{g\bar{v}}$ (grouping variables): `x`
- \bar{v} (vars): `x, y`
- `predicate`: `true`

There is a subtle point here that needs special attention. If there are two pairs (x_1, y_1) , (x_2, y_2) in the `OuterJoin` result, such that $x_1 = x_2$, they may have been derived in two ways: i) by combining the same element x in X with two different elements y_1, y_2 in Y ; ii) by combining two different elements x_1 and x_2 in X (s.t. $x_1 = x_2$) with the elements y_1, y_2 respectively. The `Nest` operator groups together two (or more) different values $y_1.a, y_2.a$ if they correspond to x_1, x_2 , such that $x_1 = x_2$ and x_1, x_2 were the same element x in X . If both conditions are satisfied then dot equality holds between x_1 and x_2 ($x_1 \doteq x_2$) [FM00]. The values $y.a$ that correspond to the same x (dot equality wise) are merged using \wedge , i.e. the merge operation of the group monoid `all`. In fact, they are merged only if they satisfy the predicate of the `Nest` operation.

The result of `Nest` is a collection of structures with labels the grouping variables $\{x\}$ and the group name `new_var`. The `OuterJoin` operator between X and Y has the monoid information `mon = bag` and the same applies for the `Nest` operator.

The operators defined so far are similar to those proposed by Fegaras and Maier. The next three operators are new in our context.

$$8. \quad \overline{\overline{\sigma}}_{\bar{v}}^{\text{mon}_1 \text{ mon}_2 \quad \overline{g\bar{v}} \text{ predicate}} \text{operand}$$

OuterSelect : This operator is similar to **Select**. The only difference is that if predicate is not satisfied, then the values of the outer variables \overline{ov} are preserved; they are appended with NULL values for the local variables in $(\overline{v} - \overline{ov})$. mon_1 indicates the way elements of the outer comprehension are accumulated, whilst mon_2 corresponds to the local (inner) comprehension. The usefulness of this operator will become evident in the rules for the translation of calculus to algebraic expressions.

$$9. \overset{mon}{M}^{\quad mon'} (e_1, \dots, e_n)$$

Merge : This operator merges the operands e_1, \dots, e_n using the merge operation of mon' . The results are accumulated using the merge operation of mon . For instance, say we have the monoid comprehension $set\{e(x)|x \leftarrow merge[bag](X_1, X_2), \dots\}$, where both X_1 and X_2 are bags. The algebraic expression corresponding to the first generator is:

$$\overset{set/bag}{M}^{\quad bag} (X_1, X_2)$$

$$10. \overset{mon}{\overline{M}}^{\quad \begin{matrix} mon' \\ parts \end{matrix}}_{\overline{v} \quad v'} \text{operand}$$

NestedMerge : This operator iterates over the operand and for each element, it merges the parts using the merge operation of monoid mon' . The result of merging them is bound to the new variable v' . Each structure in operand is extended with the new label v' . The resulting structures are accumulated using the merge operation of mon .

$$11. \overset{mon}{B}^{\quad function} \text{operand}$$

Bulk : This operator applies function, e.g. `listof`, `bagof`, `arrayof` or `orderof`, to operand. It accumulates the elements of the resulting collection using the merge operation of mon . This operator applies non-homomorphic bulk functions. For instance, if operand is a set of elements then the function `bagof(operand)` waits until the input is complete (all the elements have been received) and generates a bag containing these elements, but without any duplicates. Bulk operators are expected to create bottlenecks in any execution plan based on a pipeline mechanism.

$$12. \overset{mon}{\overline{B}}^{\quad \begin{matrix} function \\ argument \end{matrix}}_{\overline{v} \quad v'} \text{operand}$$

NestedBulk : This operator is similar to **Bulk**. However, instead of applying function to operand, it applies it to the collection argument, which is evaluated for each element in

Monoid Calculus Form	Algebraic Form
<code>< literal ></code>	<code>< literal ></code>
<code>zero[< monoid >]</code>	<code>zero[< monoid >]</code>
<code>unit[< monoid >](e)</code>	<code>unit[< monoid >](e #)</code>
<code>merge[< prim_mon >](e₁, ..., e_n)</code>	<code>merge[< prim_mon >](e₁ #, ..., e_n #)</code>
<code>struct(l₁ = e₁, ..., l_k = e_k)</code>	<code>struct(l₁ = e₁ #, ..., l_k = e_k #)</code>
<code>C(l₁ = e₁, ..., l_k = e_k)</code>	<code>C(l₁ = e₁ #, ..., l_k = e_k #)</code>
<code>e₁ < comp_op > e₂</code>	<code> e₁ # < comp_op > e₂ #</code>
<code>e₁ < arithm_op > e₂</code>	<code> e₁ # < arithm_op > e₂ #</code>
<code>e.attr_rltp_method</code>	<code> e #.attr_rltp_method</code>
<code>e(e₁, ..., e_k)</code>	<code> e #(e₁ #, ..., e_k #)</code>
<code>e₁[e₂]</code>	<code> e₁ #[e₂ #]</code>
<code>not (e)</code>	<code>not (e #)</code>
<code>-e</code>	<code>- e #</code>

Figure 5.2: Translation of calculus to algebraic constructs in Group A (figure 5.1)

operand. Each element in operand is extended with the resulting collection value of the bulk function, which is bound to the variable v' . The resulting structures are accumulated using the merge operation of `mon`.

In this section, we presented the syntax of the algebraic representation of OQL. We reviewed the algebraic operators proposed by Fegaras and Maier, proposed slight changes and introduced five new operators. We are now in a position to discuss the translation process from the monoid calculus to the algebraic level.

5.2 Translating monoid calculus to algebraic expressions

The translation of calculus expressions that have a similar algebraic equivalent (in Group A of figure 5.1) is straightforward. All we do is replace the constituent parts of these expressions by their algebraic counterparts. This is shown in table 5.2. The function $| - |^\# : \text{calculus} \rightarrow \text{algebra}$ is the generic translation algorithm that takes a monoid calculus expression and converts it to an algebraic form. Its implementation depends on the calculus expression.

The remaining calculus forms to be considered are the nonhomomorphic functions `bagof`, `listof`, `arrayof` and `orderof`, as well as the merge operation and monoid comprehensions. If the argument of a nonhomomorphic function (e.g. `bagof`) has no free variables, its algebraic equivalent is derived

as follows:

$$|\text{function}(e)|^\# = \overset{\text{mon}}{B} \overset{\text{function}}{e}$$

where $\text{function} \in \{\text{bagof}, \text{listof}, \text{arrayof}, \text{orderof}\}$ and mon is the monoid corresponding to the comprehension which has $\text{function}(e)$ as a generator.

Likewise, if the arguments of a merge calculus expression have no free variables, the expression is translated to the algebraic form:

$$|\text{merge}[\langle \text{coll_mon} \rangle](e_1, \dots, e_n)|^\# = \overset{\text{mon}}{M} \overset{\text{coll_mon}}{|e_1|^\#, \dots, |e_n|^\#}$$

where mon is derived as explained in the case of $\text{function}(e)$, if merge is a generator. Otherwise $\text{mon} = \text{coll_mon}$.

If the calculus expressions translated so far are parts of a comprehension and not all of their variables are free, then their translation is not so straightforward. It is tackled below in the algorithm for the translation of monoid comprehensions.

5.2.1 Translation of Monoid Comprehensions

Intuitively, the algorithm builds an algebraic expression by traversing the qualifiers first and the head afterwards. Initially, the algebraic expression is null. As the qualifiers are processed, it gradually grows bigger, until it takes its final form at the point it *absorbs* the head. Before giving a formal description of the algorithm we demonstrate its functionality through a simple example:

$$\begin{aligned} & (\text{set}\{x.\text{address} | x \leftarrow \text{People}, x.\text{name} = \text{"Clare"}\}) () \Rightarrow \\ & (\text{set}\{x.\text{address} | \}_x) \left(\overset{\text{set/bag}}{\sigma}_x \overset{x.\text{name}=\text{"Clare"}}{(\text{People})} \right) \Rightarrow \\ & () \left(\overset{\text{set}}{\pi}_x \overset{x.\text{address}}{\left(\overset{\text{set/bag}}{\sigma}_x \overset{x.\text{name}=\text{"Clare"}}{(\text{People})} \right)} \right) \end{aligned}$$

As we process the qualifiers and the head of the comprehension, we distinguish between the range variables that have been encountered in an outer comprehension, and those bound in previous generators of the same comprehension. The former are called `outer_vars` whilst the latter `local_vars`. Likewise, we distinguish the algebraic query produced so far by the qualifiers of the outer comprehensions, from the one generated by the previous qualifiers of the same comprehension. The first is called

outer_query and the second local_query. This is a characteristic point at which our algorithm differs from that of Fegaras and Maier. Their algorithm gradually develops a single outer_query, while we insist on maintaining a local context of the nested comprehension too. Later on, we will explain the implications of the two different approaches. At certain critical points, the local_query is merged with the outer_query. The local_query becomes empty and the local_vars are included in the outer_vars. However, we keep storing the initial *outer* variables in a vector called ext_vars. It is worth noting that if a predicate refers to a combination of outer_vars and local_vars, we defer its evaluation until the local- and outer- queries are merged. Meanwhile, we maintain it in a vector named predicates.

The input of the translation algorithm consists of the outer_query already formed and the outer_vars encountered. Obviously, if the comprehension is an external one, both input values are null, empty respectively. The algorithm has a naturally functional implementation, going along with recursive descent of a calculus syntax tree. In the remainder of this section, we present a set of rules that formally define the translation algorithm for monoid comprehensions. These rules are preceded by the necessary notation.

Notation

Each side of the rules has the following constituent elements:

- the ext_vars (external variables), i.e. the variables from outer comprehensions,
- a monoid calculus expression (which has not been converted yet to its algebraic representation),
- the local_query already formed,
- the local_vars already encountered,
- the outer_query already formed,
- the outer_vars already encountered,
- the set of predicates that combine local and outer range variables, but have not been absorbed into either of the algebraic queries (local_query or outer_query).

Thus, each side of a rule has the form:

$$\langle \text{ext_vars} \rangle \langle \text{calculus} \rangle \begin{array}{l} \langle \text{local_query} \rangle \\ \langle \text{local_vars} \rangle \end{array} \begin{array}{l} \langle \text{outer_query} \rangle \\ \langle \text{outer_vars} \rangle \end{array} \text{predicates}$$

In the monoid calculus expressions, on the bottom of each qualifier, we denote its dependency wrt the outer and the local variables. If all variables in the qualifier are bound, then dep = null. If the

free variables are all included in the `local_vars`, then `dep = local`. However, if any of its free variables is included in the `outer_vars`, then `dep = outer`.

In the translation rules from calculus to algebra we broadly use the function \mathcal{S} . This function applies on a pair of queries, one algebraic and one calculus-like. It has a different implementation depending on the calculus (the second) query. Its goal is to simplify the calculus expression, i.e. to remove any nested queries from it, by incorporating their algebraic equivalent to the initial algebraic query.

We only use this function when some of the free variables of the calculus expression are in the range variables over the algebraic query. Hence, \mathcal{S} takes three input parameters: i) an algebraic expression (AE), ii) a calculus one (CE) and iii) the range variables over the algebraic expression (\bar{v}). It returns an extended algebraic expression (AE'), a simplified calculus one (CE') and an updated set of range variables (\bar{v}') over AE'.

$$\mathcal{S} \left\{ \begin{array}{c} \text{AE} \\ \text{CE} \end{array} \right\}_{\bar{v}} = \left\{ \begin{array}{c} \text{AE}' \\ \text{CE}' \end{array} \right\}_{\bar{v}'}$$

If expression CE is a comprehension, the extended algebraic query AE' is derived by calling the function $| - |^\#$ on the comprehension with `outer_vars = \bar{v}` and `outer_query = AE`. The resulting algebraic query will be a `Nest` operator and the resulting variables $\bar{v}' = \bar{v} \cup \{g\}$, where `g` is the name of the group (`gname`) in the `Nest` operator. Also, $\text{CE}' = g$. This will become evident to the reader after studying the translation rules in subsections 5.2.1, 5.2.1 and 5.2.1.

If the only expression that allowed nested queries was the monoid comprehension form, then we would not need to introduce the \mathcal{S} function. All we would need is the recursive application of the translation rules. However, in our context, we need to simplify all kinds of nested queries, beyond monoid comprehensions. Figure 5.3 shows how \mathcal{S} works for various calculus expressions.

We are now in a position to present the rules for the translation of monoid comprehensions to their algebraic forms.

Generators without any dependency from local or outer variables

$$\begin{array}{l} \langle \text{ext_vars}_1 \rangle \langle M\{\text{head}|u \leftarrow e, \bar{r}, p\}_{\text{dep=null}} \rangle \langle \text{loc_query}_1 \rangle \langle \text{out_query} \rangle \text{pred}_1 \\ \qquad \qquad \qquad \langle \text{loc_vars}_1 \rangle \langle \text{out_vars} \rangle \\ \Rightarrow \langle \text{ext_vars}_2 \rangle \langle M\{\text{head}|\bar{r}, p[\hat{u}]\} \rangle \langle \text{loc_query}_2 \rangle \langle \text{out_query} \rangle \text{pred}_2 \\ \qquad \qquad \qquad \langle \text{loc_vars}_2 \rangle \langle \text{out_vars} \rangle \end{array}$$

where

- `p` is the calculus expression in which all filters of the comprehension have been assembled, `p[\hat{u}]` is what remains from `p` after removing the parts of it that refer to `u` only (denoted as

$\text{IF } \mathcal{S}\left\{ \begin{array}{c} \text{AE} \\ e_1 \end{array} \right\}_{\bar{v}} = \left\{ \begin{array}{c} \text{AE}_1 \\ e'_1 \end{array} \right\}_{\bar{v}_1}, \dots, \mathcal{S}\left\{ \begin{array}{c} \text{AE}_{n-1} \\ e_n \end{array} \right\}_{\bar{v}_{n-1}} = \left\{ \begin{array}{c} \text{AE}_n \\ e'_n \end{array} \right\}_{\bar{v}_n}$	
$\mathcal{S}\left\{ \begin{array}{c} \text{AE} \\ \langle \text{id} \rangle \end{array} \right\}_{\bar{v}} = \left\{ \begin{array}{c} \text{AE} \\ \langle \text{id} \rangle \end{array} \right\}_{\bar{v}}$	$\mathcal{S}\left\{ \begin{array}{c} \text{AE} \\ \text{not}(e_1) \end{array} \right\}_{\bar{v}} = \left\{ \begin{array}{c} \text{AE}_1 \\ \text{not}(e'_1) \end{array} \right\}_{\bar{v}_1}$
$\mathcal{S}\left\{ \begin{array}{c} \text{AE} \\ \text{unit}[\langle \text{mon} \rangle](e_1) \end{array} \right\}_{\bar{v}} = \left\{ \begin{array}{c} \text{AE}_1 \\ \text{unit}[\langle \text{mon} \rangle](e'_1) \end{array} \right\}_{\bar{v}_1}$	
$\mathcal{S}\left\{ \begin{array}{c} \text{AE} \\ e_1.\text{attr.rltp.meth} \end{array} \right\}_{\bar{v}} = \left\{ \begin{array}{c} \text{AE}_1 \\ e'_1.\text{attr.rltp.meth} \end{array} \right\}_{\bar{v}_1}$	
$\mathcal{S}\left\{ \begin{array}{c} \text{AE} \\ \text{function}(e_1) \end{array} \right\}_{\bar{v}} = \left\{ \begin{array}{c} \text{mon} \quad \text{function} \\ \mathbf{B} \quad e'_1 \\ \bar{v}_1 \text{ new} \quad \text{AE}_1 \\ \text{new} \end{array} \right\}_{\bar{v}_1 \cup \{\text{new}\}}$	
$\mathcal{S}\left\{ \begin{array}{c} \text{AE} \\ e_1[e_2] \end{array} \right\}_{\bar{v}} = \left\{ \begin{array}{c} \text{AE}_2 \\ e'_1[e'_2] \end{array} \right\}_{\bar{v}_2}$	$\mathcal{S}\left\{ \begin{array}{c} \text{AE} \\ e_1 \text{ op } e_2 \end{array} \right\}_{\bar{v}} = \left\{ \begin{array}{c} \text{AE}_2 \\ e'_1 \text{ op } e'_2 \end{array} \right\}_{\bar{v}_2}$
$\mathcal{S}\left\{ \begin{array}{c} \text{AE} \\ \text{merge}[\langle \text{prim} \rangle](e_1, \dots, e_n) \end{array} \right\}_{\bar{v}} = \left\{ \begin{array}{c} \text{AE}_n \\ \text{merge}[\langle \text{prim} \rangle](e'_1, \dots, e'_n) \end{array} \right\}_{\bar{v}_n}$	
$\mathcal{S}\left\{ \begin{array}{c} \text{AE} \\ \text{merge}[\langle \text{coll} \rangle](e_1, \dots, e_n) \end{array} \right\}_{\bar{v}} = \left\{ \begin{array}{c} \text{mon} \quad \text{coll} \\ \mathbf{M} \quad e'_1, \dots, e'_n \\ \bar{v}_n \text{ new} \quad \text{AE}_n \\ \text{new} \end{array} \right\}_{\bar{v}_n \cup \{\text{new}\}}$	
$\mathcal{S}\left\{ \begin{array}{c} \text{AE} \\ \text{struct}(l_1 : e_1, \dots, l_n : e_n) \end{array} \right\}_{\bar{v}} = \left\{ \begin{array}{c} \text{AE}_n \\ \text{struct}(l_1 : e'_1, \dots, l_n : e'_n) \end{array} \right\}_{\bar{v}_n}$	

Figure 5.3: Implementing \mathcal{S} for various calculus constructs

$p[u]$), or to u and any of the loc_vars_1 (denoted as $p[u, \text{loc_vars}_1]$), or to u and any of the $(\text{loc_vars}_1 \cup \text{out_vars})$ (denoted as $p[u, \text{loc_vars}_1, \text{out_vars}]$).

- if there exists $p[u]$ such that $p[u] \neq \text{true}$ then

$$- \text{gener_query} = \sigma_{\bar{v}}^{\text{mon}} \text{simpl_p}[u] e'$$

$$\text{where } \mathcal{S} \left\{ \begin{array}{c} |e|\# \\ p[u] \end{array} \right\}_{\{u\}} = \left\{ \begin{array}{c} e' \\ \text{simpl_p}[u] \end{array} \right\}_{\bar{v}}$$

otherwise,

$$- \text{gener_query} = |e|\#$$

- $\text{pred}_2 = \text{pred}_1 \cup p[u, \text{loc_vars}_1, \text{outer_vars}]$
- $\text{loc_vars}_2 = \text{loc_vars}_1 \cup \{u\} = \bar{v}_1$
- if $\text{loc_query}_1 \neq \text{null}$ then

- if there exists $p[u, \text{loc_vars}_1] \neq \text{true}$ then

$$\mathcal{S} \left\{ \begin{array}{c} \text{loc_query}_1 \bowtie_{\bar{v}_1, \text{true}}^{\text{mon}} \text{gener_query} \\ p[u, \text{loc_vars}_1] \end{array} \right\}_{\bar{v}_1} = \left\{ \begin{array}{c} \text{alg_query}_1 \\ \text{simpl_pred}_1 \end{array} \right\}_{\bar{v}_2}$$

$$\text{loc_query}_2 = \sigma_{\bar{v}_2}^{\text{mon}} \text{simpl_pred}_1 \text{alg_query}_1$$

- Otherwise, if there is no $p[u, \text{loc_vars}_1]$,

$$\text{loc_query}_2 = \text{loc_query}_1 \bowtie_{\bar{v}_1, \text{true}}^{\text{mon}} \text{gener_query}$$

Otherwise, if $\text{loc_query}_1 = \text{null}$, then

- $\text{loc_query}_2 = \text{gener_query}$

- $\text{ext_vars}_2 = \begin{cases} \text{ext_vars}_1, & \text{if } \text{ext_vars}_1 \neq \text{empty vector} \\ \text{out_vars}, & \text{if } \text{ext_vars}_1 = \text{empty vector} \end{cases}$

Generators dependent on local (but not outer) variables

$$\begin{array}{l}
\langle \text{ext_vars}_1 \rangle \langle M\{\text{head}|u \leftarrow e, \bar{r}, p\}_{\text{dep=local}} \rangle \langle \text{loc_query}_1 \rangle \langle \text{out_query} \rangle \text{pred}_1 \\
\langle \text{loc_vars}_1 \rangle \langle \text{out_vars} \rangle \\
\Rightarrow \langle \text{ext_vars}_2 \rangle \langle M\{\text{head}|\bar{r}, p[\hat{u}]\} \rangle \langle \text{loc_query}_2 \rangle \langle \text{out_query} \rangle \text{pred}_2 \\
\langle \text{loc_vars}_2 \rangle \langle \text{out_vars} \rangle
\end{array}$$

where

$$\bullet \mathcal{S} \left\{ \begin{array}{c} \text{loc_query}_1 \\ e \end{array} \right\}_{\text{loc_vars}_1} = \left\{ \begin{array}{c} \text{alg_query}_1 \\ e' \end{array} \right\}_{\bar{v}}$$

$$\bullet \text{alg_query}_2 = \underset{\bar{v}}{\mu}_u^{\text{mon } e'} \text{alg_query}_1$$

$$\bullet \text{Say } p[\widehat{\text{out_vars}}] = \text{merge}[\text{all}](p[u], p[u, \text{loc_vars}_1]).$$

$$\text{If } \mathcal{S} \left\{ \begin{array}{c} \text{alg_query}_2 \\ p[\widehat{\text{out_vars}}] \end{array} \right\}_{\bar{v} \cup \{u\}} = \left\{ \begin{array}{c} \text{alg_query}_3 \\ \text{simpl_p}[\widehat{\text{out_vars}}] \end{array} \right\}_{\bar{v}_1}$$

$$\text{then } \text{loc_query}_2 = \underset{\bar{v}_1}{\sigma}^{\text{mon } \text{simpl_p}[\widehat{\text{out_vars}}]} \text{alg_query}_3$$

$$\bullet \text{loc_vars}_2 = \bar{v}_1$$

$$\bullet \text{pred}_2 = \text{pred}_1 \cup p[u, \text{out_vars}]$$

$$\bullet \text{ext_vars}_2 = \begin{cases} \text{ext_vars}_1, & \text{if } \text{ext_vars}_1 \neq \text{empty vector} \\ \text{outer_vars}, & \text{if } \text{ext_vars}_1 = \text{empty vector} \end{cases}$$

Generators dependent (among others) on outer variables

$$\begin{array}{l}
\langle \text{ext_vars}_1 \rangle \langle M\{\text{head}|u \leftarrow e, \bar{r}, p\}_{\text{dep=outer}} \rangle \langle \text{loc_query}_1 \rangle \langle \text{out_query}_1 \rangle \text{pred}_1 \\
\langle \text{loc_vars}_1 \rangle \langle \text{out_vars}_1 \rangle \\
\Rightarrow \langle \text{ext_vars}_2 \rangle \langle M\{\text{head}|\bar{r}, p[\hat{u}]\} \rangle \langle \text{loc_query}_2 \rangle \langle \text{out_query}_2 \rangle \text{pred}_2 \\
\langle \text{loc_vars}_2 \rangle \langle \text{out_vars}_2 \rangle
\end{array}$$

where

$$\bullet \text{If } \text{loc_query}_1 \neq \text{null then}$$

$$\begin{array}{l}
- \mathcal{S} \left\{ \begin{array}{c} \text{out_query}_1 \quad \underset{\bar{v}, \text{true}}{\text{mon}_1 \text{ mon}_2} \text{loc_query}_1 \\ \text{pred}_1 \end{array} \right\}_{\bar{v}} = \left\{ \begin{array}{c} \text{alg_query} \\ \text{simpl_pred}_1 \end{array} \right\}_{\bar{v}_1} \\
\text{where } \bar{v} = \text{out_vars}_1 \cup \text{loc_vars}_1
\end{array}$$

$$- \text{alg_query}_1 = \begin{array}{c} \text{mon}_1 \text{ mon}_2 \quad \text{ext_vars}_1 \\ \text{---} \mathcal{O} \quad \text{simpl_pred}_1 \\ \bar{v}_1 \end{array} \text{alg_query}$$

Otherwise,

$$- \text{alg_query}_1 = \text{out_query}_1 \text{ and } \bar{v}_1 = \text{out_vars}_1.$$

$$\bullet \text{ Say } \mathcal{S} \left\{ \begin{array}{c} \text{alg_query}_1 \\ e \end{array} \right\}_{\bar{v}_1} = \left\{ \begin{array}{c} \text{alg_query}_2 \\ \text{simpl_e} \end{array} \right\}_{\bar{v}_2} \text{ and}$$

$$\text{alg_query}_3 = \begin{array}{c} \text{mon}_1 \text{ mon}_2 \quad \text{simpl_e} \\ \text{---} \mu \quad \text{true} \\ \bar{v}_2 \text{ u} \end{array} \text{alg_query}_2$$

- Let $\text{preds_on_u} = \text{merge}[\text{all}](\text{p}[\text{u}], \text{p}[\text{u}, \text{loc_vars}_1], \text{p}[\text{u}, \text{loc_vars}_1, \text{out_vars}_1])$.

$$\mathcal{S} \left\{ \begin{array}{c} \text{alg_query}_3 \\ \text{preds_on_u} \end{array} \right\}_{\bar{v}_2 \cup \{\text{u}\}} = \left\{ \begin{array}{c} \text{alg_query}_4 \\ \text{simpl_p} \end{array} \right\}_{\bar{v}_3}$$

$$\bullet \text{ out_query}_2 = \begin{array}{c} \text{mon}_1 \text{ mon}_2 \quad \text{ext_vars}_2 \\ \text{---} \mathcal{O} \quad \text{simpl_p} \\ \bar{v}_3 \end{array} \text{alg_query}_4$$

$$\bullet \text{ out_vars}_2 = \bar{v}_3$$

$$\bullet \text{ loc_vars}_2 = \{\}$$

$$\bullet \text{ loc_query}_2 = \text{null}$$

$$\bullet \text{ pred}_2 = \{\}$$

$$\bullet \text{ ext_vars}_2 = \begin{cases} \text{ext_vars}_1, & \text{if } \text{ext_vars}_1 \neq \text{empty vector} \\ \text{out_vars}, & \text{if } \text{ext_vars}_1 = \text{empty vector} \end{cases}$$

Absorbing remaining predicates and head of the comprehension

$$\begin{array}{l} \langle \text{ext_vars}_1 \rangle \langle \text{M}\{\text{head}\} \rangle \langle \text{loc_query}_1 \rangle \langle \text{out_query}_1 \rangle \text{pred}_1 \\ \quad \quad \quad \langle \text{loc_vars}_1 \rangle \langle \text{out_vars}_1 \rangle \\ \Rightarrow \langle \text{ext_vars}_2 \rangle \langle \rangle \langle \text{loc_query}_2 \rangle \langle \text{out_query}_2 \rangle \text{pred}_2 \\ \quad \quad \quad \langle \text{loc_vars}_2 \rangle \langle \text{out_vars}_2 \rangle \end{array}$$

where

- If $\text{loc_query}_1 \neq \text{null}$ and $\text{out_query}_1 \neq \text{null}$ then

$$\mathcal{S} \left\{ \begin{array}{c} \text{out_query}_1 \\ \begin{array}{c} \text{mon}_1 \text{ mon}_2 \\ \text{---} \otimes \\ \bar{v}, \text{true} \\ \text{pred}_1 \end{array} \\ \text{loc_query}_1 \end{array} \right\}_{\bar{v}} = \left\{ \begin{array}{c} \text{alg_query}_1 \\ \text{simpl_pred}_1 \end{array} \right\}_{\bar{v}_1}$$

where $\bar{v} = \text{out_vars}_1 \cup \text{loc_vars}_1$.

Otherwise, if $\text{out_query}_1 \neq \text{null}$

$\text{alg_query}_1 = \text{out_query}_1$ and $\bar{v}_1 = \text{out_vars}_1$

Otherwise, if $\text{loc_vars}_1 \neq \text{null}$

$\text{alg_query}_1 = \text{loc_query}_1$ and $\bar{v}_1 = \text{loc_vars}_1$.

- If $\text{pred}_1 \neq \{\}$

$$- \text{alg_query}_2 = \frac{\text{mon}_1 \text{ mon}_2 \quad \text{ext_vars}_1 \quad \text{simpl_pred}_1}{\bar{v}_1} \text{alg_query}_1$$

Otherwise,

$$- \text{alg_query}_2 = \text{alg_query}_1.$$

- $\mathcal{S} \left\{ \begin{array}{c} \text{alg_query}_2 \\ \text{head} \end{array} \right\}_{\bar{v}_1} = \left\{ \begin{array}{c} \text{alg_query}_3 \\ \text{simpl_head} \end{array} \right\}_{\bar{v}_2}$

- If $\text{ext_vars}_1 = \{\}$ then the resulting algebraic expression is:

$$- \frac{\text{M} \quad \text{simpl_head}}{\bar{v}_2} \text{alg_query}_2$$

Otherwise (if $\text{ext_vars}_1 \neq \{\}$), the resulting algebraic expression takes the following form:

$$- \frac{\text{mon} \quad \text{M} \quad \text{simpl_head}}{\bar{v}_2} \begin{array}{c} \text{group_id} \quad \text{ext_vars}_1 \\ \text{null} \quad \text{alg_query}_2 \end{array}$$

5.3 Benefits of the algebraic representation: Examples

This section presents the advantages of this algebraic framework with respect to the one proposed by Fegaras and Maier [FM00]. A number of translation examples are given to show in practice where the novel points of the current approach are applicable.

1. This framework addresses the non-homomorphic transformations involved in OQL queries, e.g. `bagof` and `listof`. To do that, it introduces the operators `Bulk` and `NestedBulk`.

Consider, for instance the following OQL query:

```

select x.C + y.D
from X as x,
      (select distinct z.A
       from x.B as z) as y

```

(5.1)

If X is a bag or list and $x.B$ is a set, bag, list, or oset ($x \in X$), the query above corresponds to the following calculus expression:

$$\text{bag}\{x.C + y.D \mid x \leftarrow X, y \leftarrow \text{bagof}(\text{set}\{z.A \mid z \leftarrow x.B\})\}$$

The calculus expression above translates to the algebraic form:

$$\begin{array}{c} \text{bag} \\ \pi \\ x,y \end{array} \begin{array}{c} x.C+y.D \\ \mu \\ x,g' y \end{array} \begin{array}{c} g' \\ \text{true} \\ \text{bag} \\ \mathbf{B} \\ x,g g' \end{array} \begin{array}{c} \text{bagof} \\ g \\ \text{bag} \\ \Gamma \\ x,z \end{array} \begin{array}{c} \text{set} \\ g \\ \text{true} \end{array} \begin{array}{c} z.A \\ x \end{array} \left(\begin{array}{c} \text{bag set/bag} \\ \mu \\ x z \end{array} \begin{array}{c} x.B \\ \text{true} \\ (X) \end{array} \right)$$

- Nested comprehensions are processed locally as much as possible, before being merged with the algebraic expressions generated in outer comprehensions. This leaves space for further algebraic optimization.

Consider, for example, the query:

```

select x.B
from X as x
where
  min (select y1.A from Y1 as y1
       where exists y2 in Y2 : y1.C1 < y2.C2
       and x.C < y2.C2) > 4

```

(5.2)

If Y_1 and Y_2 are bags, the normalized calculus representation of the previous query is:

$$\text{bag}\{x.B \mid x \leftarrow X, \min\{y_1.A \mid y_1 \leftarrow Y_1, y_2 \leftarrow Y_2, y_1.C_1 < y_2.C_2, x.C < y_2.C_2\} > 4\}$$

This calculus expression takes the algebraic form:

$$\begin{array}{c} \text{bag} \\ \pi \\ x,g \end{array} \begin{array}{c} x.B \\ \sigma \\ g>4 \\ x,g \end{array} \begin{array}{c} \text{bag} \\ \Gamma \\ x,y_1,y_2 \end{array} \begin{array}{c} \min \\ g \\ \text{true} \end{array} \begin{array}{c} y_1.A \\ x \end{array} \left(X \begin{array}{c} \text{bag set/bag} \\ \bowtie \\ x,y_1,y_2 \quad x.C < y_2.C_2 \end{array} \left(Y_1 \begin{array}{c} \text{set/bag} \\ \bowtie \\ y_1,y_2 \quad y_1.C_1 < y_2.C_2 \end{array} Y_2 \right) \right)$$

Based on the translation rules by Fegaras and Maier [FM00], the sub-expression

$$E = X \underset{x, y_1, C_1, y_2, C_2}{\overset{\text{bag set/bag}}{=}} \bowtie \left(Y_1 \underset{y_1, y_2}{\overset{\text{set/bag}}{\bowtie}} Y_2 \right) \underset{x, y_1, y_2}{\overset{\text{set/bag}}{\bowtie}} Y_2$$

would be replaced by:

$$E' = \left(X \underset{x, y_1}{\overset{\text{set/bag}}{\bowtie}} Y_1 \right) \underset{x, y_1, y_2}{\overset{\text{set/bag}}{\bowtie}} Y_2$$

Notice that we have preserved the local context of the inner comprehension by using `Join` instead of `OuterJoin`. E is expected to be more efficient than E' , because the intermediate join between Y_1 and Y_2 has a filter that reduces the size of the result. In E' the first `OuterJoin` is likely to generate a big collection. Even if it was restricted by some predicate, being an `OuterJoin` it would preserve the values of the left operand independent of whether they matched any elements of the right operand. Moreover, E' cannot be transformed to E , since the two expressions are not syntactically equivalent.

3. This framework handles the merge operations of collection monoids, e.g. union.

Consider the query:

$$\begin{aligned} & \text{select } x.C + y.D \\ & \text{from } X \text{ as } x, (x.A \text{ union } x.B) \text{ as } y \end{aligned} \quad (5.3)$$

Assume that both attributes A and B of the extent X are set collections. The calculus representation of this query is:

$$\text{bag}\{x.C + y.D \mid x \leftarrow X, y \leftarrow \text{bagof}(\text{merge}[\text{set}](x.A, x.B))\}$$

This query cannot be translated in contexts where operators like `Merge` or `OuterMerge` are not supported (e.g. [FM00]). In this framework, it takes the algebraic form:

$$\underset{x, \text{new}, \text{new}', y}{\overset{\text{bag}}{\pi}} \underset{x, \text{new}, \text{new}', y}{\overset{\text{bag}}{\mu}} \underset{x, \text{new}, \text{new}'}{\overset{\text{new}'}{\text{true}}} \underset{x, \text{new}, \text{new}'}{\overset{\text{bag}}{B}} \underset{x, \text{new}}{\overset{\text{bagof}}{\text{new}}} \underset{x, \text{new}}{\overset{\text{bag}}{M}} \underset{x, \text{new}}{\overset{\text{set}}{x.A, x.B}} (X)$$

4. The `OuterSelect` operator allows for filtering the results of a nested comprehension before processing its head. This is not possible in other contexts, when the predicate is not simple, but needs to be unnested.

Consider the OQL query:

```

select x.D
from X as x
where min (select y1.A1 + y2.A2
           from Y1 as y1, Y2 as y2
           where (forall y3 in Y3 : y1.B1 < x.B + y3.B3) and
                 (exists z in x.E : z.C > y2.C2)) > 4

```

(5.4)

If X is a bag or list and Y_1, Y_2, Y_3 and $x.E$ are sets, bags or lists, this query corresponds to the following normalized calculus form:

$$\text{bag}\{x.D \mid x \leftarrow X, \min\{y_1.A_1 + y_2.A_2 \mid y_1 \leftarrow Y_1, y_2 \leftarrow Y_2, \\ \text{all}\{y_1.B_1 < x.B + y_3.B_3 \mid y_3 \leftarrow Y_3\}, z \leftarrow x.E, z.C > y_2.C_2\} > 4\}$$

This calculus form is translated to the algebraic expression:

$$\text{bag}_{x,g_1} \pi_{x,D} \sigma_{g_1 > 4} (E_1)$$

where

$$E_1 = \text{bag}_{x,y_1,y_2} \Gamma_{\min_{g_1} \frac{y_1.A_1 + y_2.A_2}{x}} \text{true} \text{ bag set/bag}_{x,y_1,y_2,z} \mu_{z.C > y_2.C_2} \text{ bag bag}_{x,y_1,y_2,g} \sigma_{x,g} (E_2)$$

where

$$E_2 = \text{bag}_{x,y_1,y_2,y_3} \Gamma_{\text{all}_g \frac{y_1.B_1 < x.B + y_3.B_3}{x,y_1,y_2}} \text{true} \left(\left(X \text{ bag set/bag}_{x,y_1,y_2} \bowtie \left(Y_1 \text{ set/bag}_{y_1,y_2} \bowtie Y_2 \right) \right) \text{ bag set/bag}_{x,y_1,y_2,y_3} \bowtie Y_3 \right)$$

5. Annotating operators with monoid information allows for the optimization and the efficient execution of the physical plan.

This has been discussed in section 5.1.

6. Annotating the Outer· operators with two monoids allows for unnesting non-commutative collections in a commutative context. This is an issue that has not been considered in previous research.

Say that the outer comprehension of a calculus expression has a commutative monoid, e.g. `bag`. Assume also that it has a predicate, part of which is a nested comprehension with a non-commutative monoid, say `list`. This inner comprehension includes variables bound in the generators of the outer comprehension. At some point, the generators of the inner comprehension are combined with those of the outer comprehension, through an `OuterJoin` or an `OuterUnnest` operator. In our context, these two operators are annotated with two monoids, one for the outer and one for the inner comprehension. This means that the elements of the inner comprehension are not simply unnested or joined. They are annotated with a serial number denoting their order in the original nested collections. This information is essential at the time a subsequent `Nest` operator groups these values again into a non-commutative monoid. If the `Outer` operators flattened everything, say into a bag, then the `Nest` operator would try to merge elements of a bag into a list in an arbitrary way.

5.4 Conclusion

This chapter presents a framework for the translation of calculus expressions to algebraic forms. This work is based on related work by Fegaras and Maier. However, the current approach has extended their work in several respects. First, new algebraic operators have been proposed to address the translation of calculus expressions that have not been considered in other contexts. The latter expressions include `bagof(e)`, `listof(e)`, `arrayof(e)`, `orderof(e)` and `merge[< coll_mon >]`. Second, operators that have been used in other contexts are annotated with monoid information. This is useful for three reasons: i) the optimization of the algebraic plan, ii) the efficient execution of physical operators and iii) the ability to outerjoin or outer-unnest elements of a nested comprehension even if its monoid has different properties from the monoid of the external comprehension. Third, the translation algorithm from calculus to algebra has been modified in two respects: i) generators of nested comprehensions are processed locally as long as possible and ii) predicates relating variables of inner and outer comprehensions are pushed down (using `OuterSelect`); this is not possible in [FM00] for predicates that contain nested queries.

Most of the issues that are addressed in this chapter are not found in relational contexts. Many of them are due to the manipulation of multiple collection types. This feature of OQL adds complexity to the algebraic framework and requires careful manipulation of the operators borrowed from the relational algebra.

Chapter 6

Using Association Rules to Optimize Queries Semantically

In chapters 3, 4 and 5, we studied various stages of OQL processing, namely type inference, translation to calculus form, normalization and translation to algebraic form. Unlike the calculus level, where the focus was on syntactic transformations (normalization rules), the objective at the algebraic level is to identify opportunities for optimizing queries semantically. This chapter presents joint work with Dr Ken Moody, which has also been reported in the SSDBM conference [TM01]. The goal of this work is to study the use of association rules to add or eliminate constraints in the `where` clause of a `select` query. In terms of the algebraic representation of OQL, the goal is to add or remove constraints in the predicate of a `Select` operator. We take advantage of the following heuristics presented by Siegel et al. [SSS92]: i) if a selection on attribute A is implied by another selection condition on attribute B and A is not an index attribute, then the selection on A can be removed from the query; ii) if a relation R in the query has a restricted attribute A and an unrestricted cluster index attribute B, then look for a rule where the restriction on A implies a restriction on B. The contribution of our work is twofold. First, we present detailed algorithms that apply these heuristics. Hence, our ideas are easy to implement. Second, we discuss conditions under which it is worth applying these optimization techniques, and we show the extent to which they speed up query execution.

6.1 Introduction

Rules that correlate values of data attributes in large databases have been investigated for two main purposes. The first is to aid management and decision-making in large companies with great amounts of data. The second relates to query optimization, by exploiting the semantic knowledge expressed in rules [CGM90, GGMR97, SSS92, Yoo94, YSP95]. Siegel et al. [SSS92] have proposed several

heuristics to guide rule generation and so help to transform queries to more efficient forms. In this chapter we investigate algorithms for applying two of these heuristics to optimize range queries. Although we are interested in optimizing OQL, our algorithms and optimization techniques are also applicable in any relational or object-relational database system. The form of queries under consideration is:

```
select e(x)
from Extent_X as x
where p1(x) and ... and pn(x)
```

The first heuristic (H1) aims to eliminate as many predicates (or constraints) $p_i(x)$ as possible, when they are implied by other existing predicates. The second heuristic (H2) finds implications from one or more existing predicates to a new predicate on a cluster index attribute. If the latter is added to the where clause of the query and is tested first, then the existing predicates will be tested on a smaller number of objects. In order to take advantage of either of these heuristics, we use association rules that are relevant to `Extent_X`.

In our context, association rules express semantic correlations between the values of the attributes of all objects belonging to a certain extent. They consist of the antecedent, which has one or more constraints, and the consequent, which has just one constraint: $\text{constr}_1 \wedge \dots \wedge \text{constr}_n \Rightarrow \text{constr}_0$. The constraints constr_i are of the form `class_name.attributei comp_opi valuei`, where for all $0 \leq i \leq n$, $\text{comp_op}_i \in \{<, >, <=, >=, =, \neq\}$. Association rules are typically ranked by two measures of interest: confidence and support. Confidence indicates how accurate or reliable a rule is, since it is defined as the fraction of all objects satisfying the antecedent that also satisfy the consequent. Support measures the popularity of a rule, since it reflects the percentage of objects of the extent satisfying both the antecedent and the consequent. A number of algorithms have been proposed for mining association rules over categorical or quantitative data [AIS93, SA96, MY97, Par95, TUA⁺98]. In the remainder of the chapter, we assume that the process of mining rules has been completed and that we already have a warehouse with a set of rules for each extent in the database. Further, in this warehouse we store the exceptions to each rule, i.e. those objects that satisfy its antecedent but not its consequent. We also assume that the rules and their exceptions are correctly maintained when the database is updated.

The automatic use of these rules by the OQL optimizer raises many issues. The first concerns mapping the predicates in the where clause of a `select` query to constraints that could potentially be found in the antecedent (or the consequent) of association rules. This is a step towards identifying the templates of association rules that could help in our optimizations; this task is slightly different for heuristics H1 and H2. The next step is to present a common graph algorithm that goes through all existing association rules and discovers the ones that fit the templates along with their exceptions. Another graph algorithm should combine these rules to discover the possible optimization solutions,

i.e. the sets of constraints that could be omitted from (H1) or added to (H2) the query. Using the optimal optimization solution, we then present the actual transformations of OQL queries for each of the two optimization techniques. Finally, we investigate the conditions under which it is worth applying the optimizations and show the extent to which they are expected to speed up query execution.

6.2 Converting query predicates to rule constraints

The first step towards applying either of the heuristics is to map the predicates found in the where clause of a select query to constraints found in the existing association rules. Consider the following OQL query:

```
select x from Employees as x
where x.salary > 35,000 and x.salary < 55,000 and
x.position = "associate" and
x.birthYear <= 1975 and x.birthYear >= 1955
and x.emplYear <= 1990 and
x.emplYear > 1970
```

(6.1)

We first need to gather all predicates that refer to the same attribute into one predicate, and form constraints of the form `attribute_name in range`. The resulting constraints in our example are given below:

```
salary in {(35,000,55,000)}  CA
position in {"associate"}    CB
birthYear in {[1955,1975]}  CC
emplYear in {(1970,1990)}   CD
```

Based on these query constraints, the next step is to search for *similar* constraints in the association rules. Suppose that the warehouse contains the association rules illustrated in figure 6.1, which are relevant to the class `Employee`. The constraints in the antecedent and the consequent of the rules are labeled with short names, e.g. $C_{k,R_{kj}}$ is a constraint limiting the values of attribute k to the range R_{kj} . Here a *range* is a specified set of values of some attribute. The set of ranges associated with a particular attribute is partially ordered under inclusion in the usual way.

To distinguish constraints that are derived from the where clause of our query from those in the antecedent or consequent of the association rules, we refer to the former as *query constraints*, and to the latter as *rule constraints*. For each heuristic we look for rule constraints corresponding to the query constraints.

$$\begin{array}{l}
\text{birthYear in } \{[1955, 1975]\}(C_{1.R_{11}}) \Rightarrow \text{salary in } \{[35, 000, 50, 000]\}(C_{2.R_{21}}) \\
\text{position in } \{\text{"associate"}\}(C_{4.R_{41}}) \Rightarrow \text{salary in } \{[30, 000, 60, 000]\}(C_{2.R_{22}}) \\
\text{position in } \{\text{"manager"}\}(C_{4.R_{42}}) \Rightarrow \text{birthYear in } \{[1955, 1980]\}(C_{1.R_{12}}) \\
\text{birthYear in } \{[1950, 1970]\}(C_{1.R_{13}}) \Rightarrow \text{position in } \{\text{"associate"}\}(C_{4.R_{41}}) \\
\text{emplYear in } \{[1975, 1990]\}(C_{5.R_{53}}) \Rightarrow \text{position in } \{\text{"associate"}\}(C_{4.R_{41}}) \\
\text{emplYear in } \{[1965, 1985]\}(C_{5.R_{52}}) \Rightarrow \text{efficiency in } \{[3, 6]\}(C_{3.R_{31}}) \\
\left. \begin{array}{l} \text{salary in } \{[30, 000, 65, 000]\}(C_{2.R_{23}}) \\ \text{position in } \{\text{"manager"}\}(C_{4.R_{42}}) \end{array} \right\} \Rightarrow \text{emplYear in } \{[1970, 1985]\}(C_{5.R_{51}}) \\
\left. \begin{array}{l} \text{salary in } \{[30, 000, 60, 000]\}(C_{2.R_{22}}) \\ \text{efficiency in } \{[3, 6]\}(C_{3.R_{31}}) \end{array} \right\} \Rightarrow \text{birthYear in } \{[1955, 1975]\}(C_{1.R_{11}}) \\
\left. \begin{array}{l} \text{salary in } \{[30, 000, 60, 000]\}(C_{2.R_{22}}) \\ \text{efficiency in } \{[3, 6]\}(C_{3.R_{31}}) \end{array} \right\} \Rightarrow \text{position in } \{\text{"associate"}\}(C_{4.R_{41}}) \\
\left. \begin{array}{l} \text{birthYear in } \{[1955, 1980]\}(C_{1.R_{12}}) \\ \text{salary in } \{[35, 000, 50, 000]\}(C_{2.R_{21}}) \end{array} \right\} \Rightarrow \text{efficiency in } \{[3, 6]\}(C_{3.R_{31}})
\end{array}$$

Figure 6.1: Association rules relevant to the class Employee

6.2.1 Heuristic H1

We first discuss rule constraints relevant to the constraint elimination heuristic. For each constraint C on `attribute_name` with range R , identify the following rule constraints:

Category 1 - relaxed constraints First, identify the constraints whose ranges are the next more general than (or the same as) the range of the original constraint C . If there exists such a constraint C_0 with a range R_0 , then for all other constraints C_i $i \neq 0$ on `attribute_name` with ranges R_i , if $R \subseteq R_i$ then $\text{not}(R_i \subseteq R_0)$. That is, R_0 is one of the least supersets of R among all ranges of constraints referring to the same attribute.

Category 2 - relaxed combinations of constraints Second, identify minimal combinations of constraints which have the property that the union of their ranges is a superset of the range of the original constraint C . The constraints in any combination must satisfy the following conditions: firstly, none of the constraints can be the same as (or with a wider range than the range of) any constraint identified previously in category 1; secondly, no combination must include all the constraints of a smaller combination; thirdly, if either of two constraints on the same attribute may take part in a combination, and the range of the one is a subset of the range of the other, then the one with the smaller range is selected.

Category 3 - tightened constraints Third, identify the constraints whose ranges are the next more specific than (or the same as) the range of the original constraint. If there exists such a constraint C_0 with a range R_0 , it means that for all other constraints C_i $i \neq 0$ on `attribute_name` with ranges R_i , if $R \supseteq R_i$ then $\text{not}(R_i \supseteq R_0)$. That is, R_0 is one of the greatest subsets of R among all ranges of constraints referring to the same attribute.

Category 4 - tightened combinations of constraints Fourth, identify maximal combinations of constraints which have the property that the intersection of their ranges is a subset of the range of the original constraint C . The constraints in any combination must satisfy the following conditions: firstly, none of the constraints can be the same as (or with a range contained in the range of) any constraint identified previously in category 3: secondly, no combination must include all the constraints of a smaller combination: thirdly, if either of two constraints on the same attribute may take part in a combination, and the range of the one is a superset of the range of the other, then the one with the wider range is selected.

In the example of query 6.1, the query constraints C_A , C_B , C_C and C_D correspond to the following rule constraints:

	<i>Relaxed Constraints</i>	<i>Tightened Constraints</i>
C_A	salary in {[30, 000, 60, 000]}($C_{2.R_{22}}$)	salary in {[35, 000, 50, 000]}($C_{2.R_{21}}$)
C_B	position in {"associate"}($C_{4.R_{41}}$)	position in {"associate"}($C_{4.R_{41}}$)
C_C	birthYear in {[1955, 1975]}($C_{1.R_{11}}$)	birthYear in {[1955, 1975]}($C_{1.R_{11}}$)
C_D		emplYear in {[1970, 1985]}($C_{5.R_{51}}$) emplYear in {[1975, 1990]}($C_{5.R_{53}}$)
	<i>Relaxed Combinations of Constraints</i>	<i>Tightened Combinations of Constraints</i>
C_C		{birthYear in {[1955, 1980]}($C_{1.R_{12}}$) birthYear in {[1950, 1970]}($C_{1.R_{13}}$)}
C_D	{emplYear in {[1970, 1985]}($C_{5.R_{51}}$) emplYear in {[1975, 1990]}($C_{5.R_{53}}$)}	

For each tightened combination of constraints we generate a new constraint whose range is the intersection of the ranges of the constraints in the combination. We add this new constraint to the category of tightened constraints and create a temporary rule with the initial constraints in the antecedent and the new tightened constraint in the consequent. For instance, based on the tightened combination {birthYear in {[1955, 1980]}($C_{1.R_{12}}$), birthYear in {[1950, 1970]}($C_{1.R_{13}}$)}, we create the constraint {birthYear in {[1955, 1970]}($C_{1.R_{1,23}}$)}. Repeating this procedure for all tightened combinations of constraints, we generate a new set of tightened constraints.

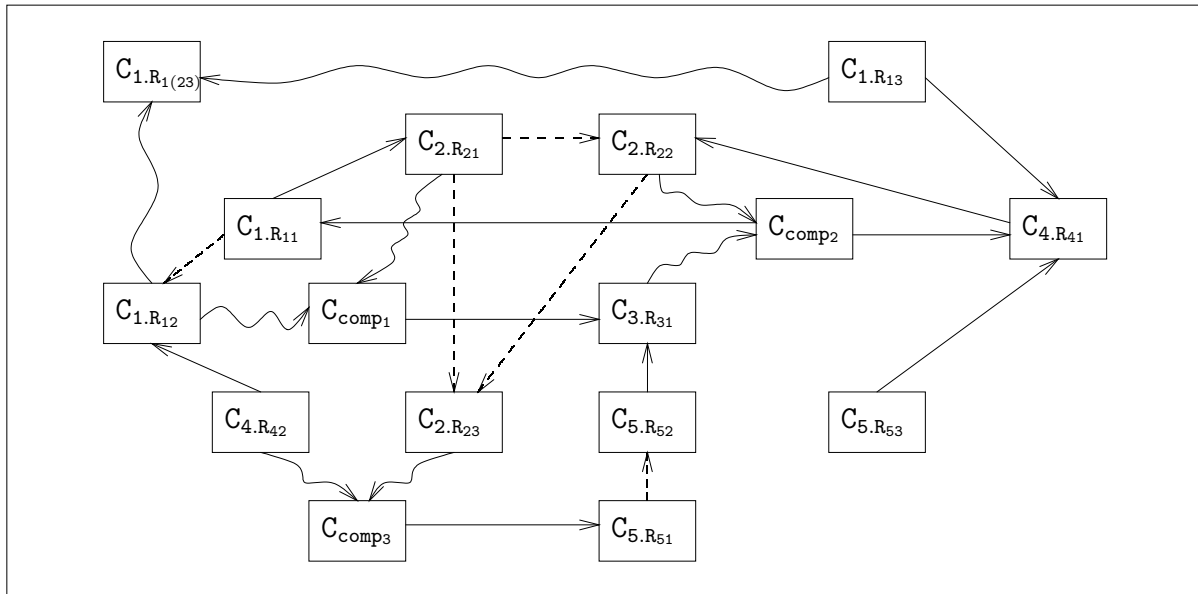


Figure 6.2: Graph of Constraints

The rationale behind finding the constraints in categories 1 to 3 is the following. By definition, the query constraints imply their corresponding relaxed constraints. Likewise, the tightened constraints imply the initial query constraints. Hence if we can find association rules that relate constraints in the first two categories to those in the third, it is equivalent to finding rules that correlate the initial query constraints with each other.

6.2.2 Heuristic H2

For the purposes of applying the constraint introduction heuristic, we identify the following rule constraints:

Category 1 - relaxed constraints These are identical to the relaxed constraints described in section 6.2.1.

Category 2 - relaxed combination of constraints This category is the same as its counterpart in section 6.2.1.

Category 3 - index constraints If the extent has a cluster index on attribute $attr_i$, then this category includes all the constraints on this attribute that occur as consequents in the existing rules. Assume that the cluster index of the extent *Employees* is the attribute *efficiency*. Based on the rules of figure 6.1, the only index constraints is the constraint $efficiency \text{ in } \{[3, 6]\}$.

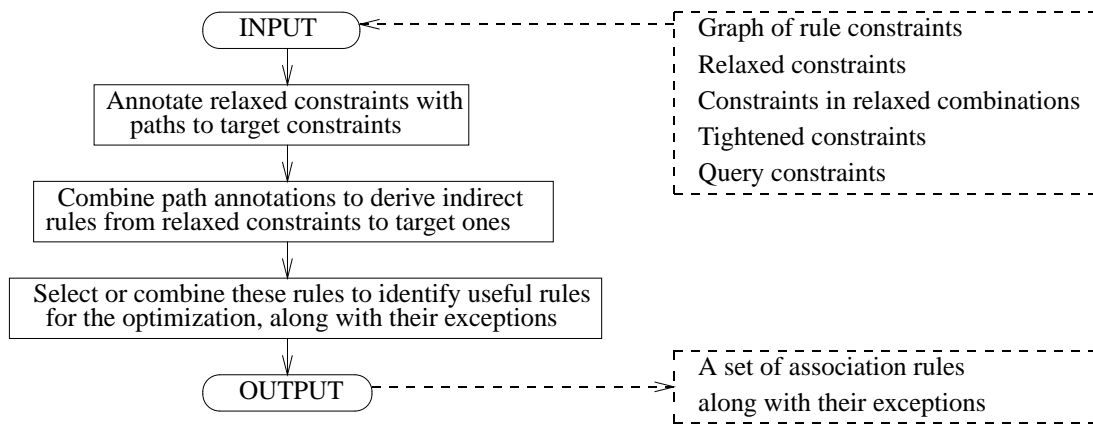


Figure 6.3: An abstract view of the algorithm

Rule constraints imply relaxed (combinations of) constraints, by definition. If we succeed in finding paths from relaxed (combinations of) constraints to an index constraint, then we can add the latter to the *where* clause of the query. This allows us to reduce the disc access overhead without altering the results of the query.

6.2.3 Discussion

Constraints in categories 1 and 2 play an identical role for heuristics H1 and H2. However, the other categories play a different role in the two cases. Tightened constraints are similar to existing query constraints that can potentially be eliminated. Index constraints are new constraints that can potentially be added to the *where* clause of the query. In the next section we present a graph algorithm that finds associations from the relaxed constraints to the tightened (H1) or the index (H2) constraints respectively. Since the algorithm is common to both heuristics, tightened and index constraints are hereafter referred to as *target constraints*.

6.3 Identifying associations between constraints

In this section, we present an algorithm that goes through the association rules in the warehouse and identifies direct or indirect correlations between the relaxed and the target constraints discussed in the previous section. Before presenting this algorithm, we discuss an alternative graphical representation of the association rules. Consider for example the association rules in figure 6.1. They can be represented using a directed graph of constraints, see figure 6.2.

A rule having a single constraint in its antecedent (e.g. $C_{1.R_{11}} \Rightarrow C_{2.R_{21}}$) is represented by two vertices, labeled with the constraints ($C_{1.R_{11}}$ and $C_{2.R_{21}}$) and linked together by an edge from the an-

tecedent constraint ($C_{1.R_{11}}$) to the consequent constraint ($C_{2.R_{21}}$). A rule with more than one constraint in its antecedent uses an additional vertex to denote the conjunction. The link \rightsquigarrow denotes the fact that a constraint in the antecedent of a rule implies the consequent only in combination with the other conjuncts forming the antecedent. The dashed edges link pairs of nodes that represent constraints on the same attribute, with the source constraint implying the destination constraint.

The notion of path in our context requires one extension to the standard notion for a directed graph. A path from a constraint C_{source} to a constraint C_{dest} is any sequence of vertices linked by directed edges, provided that no vertex represents a composite constraint. If a path contains a composite constraint, say $C = C_1 \wedge \dots \wedge C_n$, then there must be n paths from C_{source} , one to each of the conjuncts C_i . These subpaths merge at the node representing the composite constraint, and from there the path continues towards C_{dest} .

We consider also paths from a set of constraints $\{C_1, \dots, C_n\}$ to a single constraint. Such a path consists of subpaths starting from the source constraints $\{C_1, \dots, C_n\}$ which merge at various composite constraints before reaching the destination constraint.

We can now develop an algorithm that navigates over a graph of constraints, finds paths linking relaxed (combinations of) constraints to target constraints and combines them to derive association rules useful for the purposes of optimization. It consists of three basic steps as shown in figure 6.3. The algorithm is discussed in detail in section 6.3.1. At the end of each step we give an example to show how it can be applied for heuristic H1. It is equally applicable to heuristic H2.

Relaxed constraints and constraints in relaxed combinations often play the same role in the following algorithm, and we use the term *source constraint* in such contexts.

6.3.1 Algorithm

For each constraint $C_{k.R_{kj}}$ in the target constraints take the following steps:

Step 1

Navigate backwards from $C_{k.R_{kj}}$ in the graph of constraints, i.e. navigate against the direction of the edges and annotate any source constraints encountered with information about the reverse path traversed so far. On encountering a composite constraint navigate backwards from *all* of its conjuncts. If a particular constraint is encountered more than once, or if there is no incoming link, then backtracking from the current constraint terminates.

Step 1 may be implemented as a recursive method `backtrack` on a class `Constraint`. We demonstrate its functionality by an example. Consider the target constraint $C_{4.R_{41}}$ in the constraint graph (figure 6.2). We recursively apply `backtrack` to all the constraints leading to it, i.e. $C_{1.R_{13}}$, $C_{5.R_{53}}$ and C_{comp_2} . The argument for the method is the path traversed so far, i.e. $\{C_{4.R_{41}}\}$. $C_{1.R_{13}}$ is

<i>Relaxed Constraints</i>	<i>Path Annotations</i>
$C_{2.R_{22}}$ $C_{4.R_{41}}$ $C_{1.R_{11}}$	$\{C_{2.R_{22}} \rightsquigarrow C_{comp2} \rightarrow C_{4.R_{41}}\}$ $\{C_{2.R_{22}} \rightarrow C_{2.R_{23}} \rightsquigarrow C_{comp3} \rightarrow C_{5.R_{51}} \rightarrow C_{5.R_{52}} \rightarrow C_{3.R_{31}} \rightsquigarrow$ $\rightsquigarrow C_{comp2} \rightarrow C_{4.R_{41}}\}$ $\{C_{1.R_{11}} \rightarrow C_{2.R_{21}} \rightarrow C_{2.R_{22}} \rightsquigarrow C_{comp2} \rightarrow C_{4.R_{41}}\}$ $\{C_{1.R_{11}} \rightarrow C_{2.R_{21}} \rightarrow C_{2.R_{22}} \rightarrow C_{2.R_{23}} \rightsquigarrow C_{comp3} \rightarrow C_{5.R_{51}} \rightarrow$ $\rightarrow C_{5.R_{52}} \rightarrow C_{3.R_{31}} \rightsquigarrow C_{comp2} \rightarrow C_{4.R_{41}}\}$ $\{C_{1.R_{11}} \rightarrow C_{2.R_{21}} \rightsquigarrow C_{comp1} \rightarrow C_{3.R_{31}} \rightsquigarrow C_{comp2} \rightarrow C_{4.R_{41}}\}$ $\{C_{1.R_{11}} \rightarrow C_{1.R_{12}} \rightsquigarrow C_{comp1} \rightarrow C_{3.R_{31}} \rightsquigarrow C_{comp2} \rightarrow C_{4.R_{41}}\}$
<i>Constraints in Relaxed Combinations</i>	<i>Path Annotations</i>
$C_{5.R_{51}}$ $C_{5.R_{53}}$	$\{C_{5.R_{51}} \rightarrow C_{5.R_{52}} \rightarrow C_{3.R_{31}} \rightsquigarrow C_{comp2} \rightarrow C_{4.R_{41}}\}$ $\{C_{5.R_{53}} \rightarrow C_{4.R_{41}}\}$

Figure 6.4: Paths from relaxed constraints to the target constraint $C_{4.R_{41}}$

not a source constraint, so no path annotation is assigned to it. No constraint leads to it (there is no incoming link), so there is no further recursive call from it. However, $C_{5.R_{53}}$ is a constraint in a relaxed combination, so it is annotated with the path $\{C_{4.R_{41}}\}$. No recursive call occurs from it either. C_{comp2} is a composite constraint, so we recursively call `backtrack` on each of the conjuncts $C_{2.R_{22}}$ and $C_{3.R_{31}}$, passing as argument the updated path traversed $\{C_{comp2} \rightarrow C_{4.R_{41}}\}$. Recursive calls of `backtrack` continue until the constraint on which a call is made has no incoming link or the constraint appears in the argument path annotation. In figure 6.4, we present the path annotations assigned to source constraints by the time Step 1 is completed for the target constraint $C_{4.R_{41}}$.

Discussion of step 2

The objective of step 2 is to combine the annotations of source constraints to identify *complete* paths from sets of these constraints to a target constraint $C_{k.R_{kj}}$. The first step of the algorithm annotated the source constraints with complete or incomplete paths from them to a particular target constraint. For instance, $C_{1.R_{11}}$, which is a relaxed constraint, is annotated with four paths, as shown in figure 6.4. Paths that include \rightsquigarrow links are called incomplete; they are useful only in combination

A(e)	B(e)	C(e)	e in E(A \Rightarrow B)	e in E(B \Rightarrow C)	e in E(A \Rightarrow C)
T	T	T	F	F	F
T	T	F	F	T	T
T	F	T	T	F	F
T	F	F	T	F	T
F	T	T	F	F	F
T	T	F	F	T	F
F	F	T	F	F	F
F	F	F	F	F	F

Figure 6.5: Truth table which shows the correctness of formula 6.2

with other incomplete paths that include the related conjuncts. One of the objectives of step 2 is to combine incomplete paths in order to identify (complete) paths from a source constraint like $C_{1.R_{11}}$ to a target constraint like $C_{4.R_{41}}$. In fact, we can combine paths from more than one source constraint to a destination constraint.

Each edge \rightarrow stands for an association rule which possibly has a list of exceptions E. A series of consecutive edges \rightarrow corresponds to an indirect association from the source constraint to the destination one. The exceptions to this rule are evaluated from the exceptions of the association rules involved. In particular, if constraint A leads to constraint B based on a rule with exceptions $E_{A \Rightarrow B}$, and B leads to C based on a rule with exceptions $E_{B \Rightarrow C}$, then the exception set of the indirect path $A \rightarrow C$ (or the indirect rule $A \Rightarrow C$) is:

$$E = \{exc | exc \leftarrow E_{A \Rightarrow B}, \text{ not } C(exc)\} \cup \{exc | exc \leftarrow E_{B \Rightarrow C}, A(exc)\} \quad (6.2)$$

Note that $E_{A \Rightarrow B} \cap E_{B \Rightarrow C} = \{\}$, since all exceptions of $E_{A \Rightarrow B}$ satisfy A, but do not satisfy B. Therefore they are not exceptions of $E_{B \Rightarrow C}$. The correctness of formula 6.2 is proved by the truth table in figure 6.5.

Step 2

Path annotations derived from step 1 are combined to form complete paths by repeating steps 2.1 and 2.2:

Step 2.1 For all path annotations we omit all the initial constraints, until the first constraint which is followed by a \rightsquigarrow edge (this constraint is maintained). To omit these initial constraints, we must

first evaluate the exceptions involved in the rules until the \rightsquigarrow link (equation 6.2). The resulting paths are either empty or start with sub-constraints leading to composite constraints.

When this step is executed for the first time the resulting paths in our example are the following:

$$\begin{aligned}
C_{2.R_{22}} \quad P_1 &: \{C_{2.R_{22}} \rightsquigarrow C_{\text{comp}2} \rightarrow C_{4.R_{41}}\} \\
E_1 &= \{\} \\
C_{2.R_{22}} \quad P_2 &: \{C_{2.R_{23}} \rightsquigarrow C_{\text{comp}3} \rightarrow C_{5.R_{51}} \rightarrow C_{5.R_{52}} \rightarrow C_{3.R_{31}} \rightsquigarrow C_{\text{comp}2} \rightarrow C_{4.R_{41}}\} \\
E_2 &= E(C_{2.R_{22}} \Rightarrow C_{2.R_{23}}) = \{\} \\
C_{1.R_{11}} \quad P_3 &: \{C_{2.R_{22}} \rightsquigarrow C_{\text{comp}2} \rightarrow C_{4.R_{41}}\} \\
E_3 &= \{e \mid e \leftarrow E(C_{1.R_{11}} \Rightarrow C_{2.R_{21}}), \text{not } C_{2.R_{22}}(e)\} \\
C_{1.R_{11}} \quad P_4 &: \{C_{2.R_{23}} \rightsquigarrow C_{\text{comp}3} \rightarrow C_{5.R_{51}} \rightarrow C_{5.R_{52}} \rightarrow C_{3.R_{31}} \rightsquigarrow C_{\text{comp}2} \rightarrow C_{4.R_{41}}\} \\
E_4 &= \{e \mid e \leftarrow E(C_{1.R_{11}} \Rightarrow C_{2.R_{21}}), \text{not } C_{2.R_{22}}(e), \text{not } C_{2.R_{23}}(e)\} \\
&= \{e \mid e \leftarrow E(C_{1.R_{11}} \Rightarrow C_{2.R_{21}}), \text{not } C_{2.R_{23}}(e)\} \\
C_{1.R_{11}} \quad P_5 &: \{C_{2.R_{21}} \rightsquigarrow C_{\text{comp}1} \rightarrow C_{3.R_{31}} \rightsquigarrow C_{\text{comp}2} \rightarrow C_{4.R_{41}}\} \\
E_5 &= E(C_{1.R_{11}} \Rightarrow C_{2.R_{21}}) \\
C_{1.R_{11}} \quad P_6 &: \{C_{1.R_{12}} \rightsquigarrow C_{\text{comp}1} \rightarrow C_{3.R_{31}} \rightsquigarrow C_{\text{comp}2} \rightarrow C_{4.R_{41}}\} \\
E_6 &= \{\} \\
C_{5.R_{51}} \quad P_7 &: \{C_{3.R_{31}} \rightsquigarrow C_{\text{comp}2} \rightarrow C_{4.R_{41}}\} \\
E_7 &= \{\} \cup \{e \mid e \leftarrow E(C_{5.R_{52}} \Rightarrow C_{3.R_{31}}), C_{5.R_{51}}(e)\} \\
C_{5.R_{53}} \quad P_8 &: \{\} \\
E_8 &= E(C_{5.R_{53}} \Rightarrow C_{4.R_{41}})
\end{aligned}$$

If a path annotation is empty, it means that there is a complete path from the source constraint (e.g. $C_{5.R_{53}}$) to the destination one ($C_{4.R_{41}}$). A new rule is created from the former to the latter ($C_{5.R_{53}} \Rightarrow C_{4.R_{41}}$) and the exceptions evaluated so far for this path annotation (E_8) are assigned to the new association rule. The path annotation and its exceptions are deleted. If an empty annotation corresponds to more than one source constraint then the rule that is generated has a composite antecedent, i.e. it is of the form: $C_1 \wedge \dots \wedge C_n \Rightarrow C_0$.

When using equation 6.2 to determine exceptions during step 2.1 we must take note of the exception sets arising from earlier iterations of steps 2.1 and 2.2.

Step 2.2 For each composite constraint C_{comp} at the second position of some path annotation, we try to find a set of path annotations having at their first positions the component constraints of C_{comp} . These path annotations are combined into new path annotations as follows: i) the subconstraints at their first positions are omitted; ii) the resulting exceptions of the remaining paths are the union of the exceptions evaluated so far for each combined path that satisfy the

source constraints of the other combined paths. If $C = C_1 \wedge C_2$,

$$\begin{aligned} E(C_A \wedge C_B \Rightarrow C) = \\ \{e \mid e \leftarrow E(C_A \Rightarrow C_1), C_B(e)\} \cup \\ \{e \mid e \leftarrow E(C_B \Rightarrow C_2), C_A(e)\} \end{aligned} \quad (6.3)$$

When step 2.2 is executed for the first time in our example, the following path annotations are generated:

$$\begin{array}{ll} C_{2.R_{22}}, C_{5.R_{51}} & P_9 : \{C_{\text{comp}_2} \rightarrow C_{4.R_{41}}\} \\ & E_9 = \{e \mid e \leftarrow E_7, C_{2.R_{22}}(e)\} \\ C_{1.R_{11}}, C_{5.R_{51}} & P_{10} : \{C_{\text{comp}_2} \rightarrow C_{4.R_{41}}\} \\ & E_{10} = \{e \mid e \leftarrow E_3, C_{5.R_{51}}(e)\} \\ & \cup \{e \mid e \leftarrow E_7, C_{1.R_{11}}(e)\} \\ C_{1.R_{11}} & P_{11} : \{C_{\text{comp}_1} \rightarrow C_{3.R_{31}} \rightsquigarrow C_{\text{comp}_2} \rightarrow C_{4.R_{41}}\} \\ & E_{11} = \{e \mid e \leftarrow E_5, C_{1.R_{11}}(e)\} \\ & = E(C_{1.R_{11}} \Rightarrow C_{2.R_{21}}) \end{array}$$

If two path annotations which are combined do not have the same paths after the composite constraint (after the second position), then two combinations of paths should be generated, one for each annotation.

We retain path annotations P_1 to P_8 in order to combine them with new annotations at a later execution of step 2.2.

Substeps 2.1 and 2.2 are repeated until neither of them has any effect on the path annotations. This happens when all paths that have not been converted to rules in step 2.1 cannot be further combined in step 2.2.

Step 3

Not all the rules derived in step 2.1 are useful for optimization purposes. The aim of this step is to filter and combine the rules derived from the previous step, referred to as R_A , in order to generate rules relating query constraints with each other (H1), or query constraints to new index constraints (H2). The resulting set of rules is referred to as R_B . We first give a detailed account of step 3 in the context of heuristic H1. We then point out a detail that should be changed so that the step is also applicable for H2.

We look for (sets of) rules in R_A of the following kinds:

- A rule from a relaxed constraint to the target constraint C_i , i.e. of the form $C \Rightarrow C_0$. Such a rule is used to generate a new rule $C' \Rightarrow C'_0$, such that C' , C'_0 are the query constraints that refer to

the same attributes as C and C_0 respectively, C is one of the relaxed constraints of C' and C_0 is one of the target constraints of C'_0 . Indeed,

$$\left. \begin{array}{l} C' \Rightarrow C \\ C \Rightarrow C_0 \text{ (existing rule)} \\ C_0 \Rightarrow C'_0 \end{array} \right\} C' \Rightarrow C'_0$$

If $C \Rightarrow C_0$ has exceptions E , the exceptions of the new rule $C' \Rightarrow C'_0$ are $E' = \{e | e \leftarrow E, C'(e)\}$.

- A set of rules of the form $C_i \Rightarrow C_0$, $i = 1, \dots, n$, such that C_i are all constraints in the same relaxed combination of size n . This set of rules is used to generate the new rule $C' \Rightarrow C'_0$, in which C' , C'_0 are the query constraints that refer to the same attribute as the C_i and C_0 respectively. If the corresponding exception lists of the rules $C_i \Rightarrow C_0$, $i = 1, \dots, n$ are E_i , the exceptions of the new rule are $E = \{e | e \leftarrow E_1 \cup \dots \cup E_n, C'(e)\}$.
- A rule from a set of relaxed constraints to a target one, i.e. $C_1 \wedge \dots \wedge C_n \Rightarrow C_0$. This rule is converted to the corresponding rule $C'_1 \wedge \dots \wedge C'_n \Rightarrow C'_0$, such that C'_0, \dots, C'_n are the query constraints on the same attributes as C_0, \dots, C_n respectively. The exceptions of the new rule are those of the initial rule that satisfy $C'_1 \wedge \dots \wedge C'_n$, i.e. $E' = \{e | e \leftarrow E, C'_1(e), \dots, C'_n(e)\}$.
- A set of rules whose antecedents contain both relaxed constraints and constraints in relaxed combinations is useful if it consists of all rules of the form:

$$\underbrace{C_1 \dots C_k}_{\text{Relaxed constraints}} \quad \underbrace{C_{(k+1).j_1} \dots C_{(k+m).j_m}}_{\text{Constraints in relaxed combinations}} \Rightarrow C_0$$

where $1 \leq j_i \leq n_i$ for all $1 \leq i \leq m$.

Constraints $C_{(k+i).1}, \dots, C_{(k+i).n_i}$ form the relaxed combination i . The set contains $N = n_1 \times \dots \times n_m$ rules, such that each constraint in a relaxed combination occurs in all possible combinations with constraints from the other relaxed combinations. In some of the rules some (or all) of the relaxed constraints $C_1 \dots C_n$ may not be present.

If such a set of rules occurs in R_A we derive $C'_1 \wedge \dots \wedge C'_k \wedge C'_{k+1} \wedge \dots \wedge C'_{k+m} \Rightarrow C'_0$, such that $C'_1, \dots, C'_{k+m}, C'_0$ are the query constraints that refer to the same attributes as $C_1, \dots, C_{(k+m).j_m}, C_0$ respectively. If the original rules have exception lists E_1, \dots, E_N , then the resulting rule has exceptions $E = \{e | e \leftarrow E_1 \cup \dots \cup E_N, C'_1(e), \dots, C'_{k+m}(e)\}$.

Evaluating exceptions appears very expensive in this case. However, in the common situation that $m = 1$, the cost of finding exceptions to the new rule is similar to the cost in the second case of step 3.

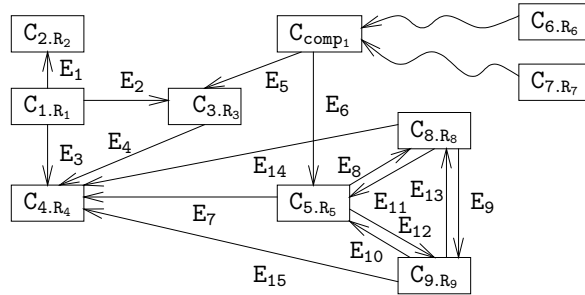


Figure 6.6: Correlations of query constraints

Note that in the context of heuristic H2, C_0 is a constraint on the cluster index attribute; it does not correspond to any of the existing query constraints. Hence, step 3 is applicable to H2, provided that in all four cases above, $C'_0 = C_0$.

6.4 Identifying optimization solutions

In the previous section, we gave an algorithm that navigates over a constraint graph and extracts a set of useful association rules along with their exceptions. This section combines these rules to identify all possible solutions to the constraint elimination or constraint introduction problem.

Consider an OQL `select` query with say 15 predicates. We combine all predicates referring to the same attribute into a single constraint $\text{attr}_i \text{ in range}_i (C_{i.R_i})$. Assume that we derive the following set of constraints: $C_{1.R_1}, \dots, C_{12.R_{12}}$, say. Since there is just one range for each attribute, ranges are presented as R_i instead of R_{ij} . We first identify the rule constraints of possible interest for H1 or H2, then apply the graph algorithm presented in section 6.3. The next step is discussed separately for each heuristic.

6.4.1 Heuristic H1

Say that the graph algorithm (section 6.3) results in the associations illustrated in figure 6.6. Not all query constraints need be related; in our example only the first nine are. Note that the association graph (figure 6.6) is transitive; if there is a direct link from constraint A to constraint B, and another one from B to C, then A and C are also directly connected. This is a property of the graph algorithm in section 6.3. We now define an algorithm that produces a set of constraint elimination solutions $\{(M_i, E_i)\}$, where M_i, E_i stand for Maintained Constraints and Exceptions respectively. Each pair (M_i, E_i) is a solution with the following explanation:

We may eliminate all but the constraints M_i from the `where` clause of the query, given that we take into account the exceptions E_i .

Algorithm

Step 1 Identify all constraints that do not have any incoming links. These constraints should be maintained (not eliminated), since no other constraint implies them. In figure 6.6, these constraints are $C_{1.R_1}$, $C_{6.R_6}$, $C_{7.R_7}$.

Step 2 Identify all cyclic paths that have no incoming link from any external constraint. If two cyclic paths have at least one common constraint, they are considered as a single cyclic path. This case does not occur in our example. Form combinations of constraints by choosing one constraint from every cyclic path identified. In our example there is just one such path containing the constraints $C_{5.R_5}$, $C_{8.R_8}$, $C_{9.R_9}$; therefore the combinations consist of just one element: $\{C_{5.R_5}\}$, $\{C_{8.R_8}\}$, $\{C_{9.R_9}\}$. Extend each of these combinations with the constraints without incoming links, identified in step 1. The resulting combinations are $\{C_{5.R_5}, C_{1.R_1}, C_{6.R_6}, C_{7.R_7}\}$, $\{C_{8.R_8}, C_{1.R_1}, C_{6.R_6}, C_{7.R_7}\}$ and $\{C_{9.R_9}, C_{1.R_1}, C_{6.R_6}, C_{7.R_7}\}$. Each combination is a minimal set of constraints that implies all the remaining constraints in the graph (figure 6.6).

Step 3 For each combination evaluate the exceptions that are involved in removing the implied constraints. If a constraint is implied by more than one constraint in a combination, then consider the exceptions of the strongest implication (the implication with the fewest exceptions). For example, the exceptions corresponding to the combination $\{C_{5.R_5}, C_{1.R_1}, C_{6.R_6}, C_{7.R_7}\}$ are evaluated by forming the union of the exception lists: E_1 for the elimination of C_2 ; $\text{min_set}\{E_2, E_5\}$ for the elimination of C_3 ; $\text{min_set}\{E_3, E_6, E_7\}$ for the elimination of C_4 ; E_8 for the elimination of C_8 ; E_{12} for the elimination of C_9 .

Step 4 Form the final solutions to the elimination problem by adding to each combination the remaining query constraints that are not related to each other ($C_{10.R_{10}}$, $C_{11.R_{11}}$, $C_{12.R_{12}}$). The exceptions corresponding to each solution are filtered so that they satisfy all the maintained constraints, i.e. all the constraints in the combination.

Note that all the solutions identified in the algorithm have the same number of constraints. We assume that the cardinality of the extent against which the query is run is far greater than any of the exception lists annotating the optimization solutions. Therefore, there is no advantage in eliminating only a subset of the constraints in a solution in order to decrease the number of exceptions involved.

6.4.2 Heuristic H2

Assume that in the context of H2 the graph algorithm (section 6.3) results in the associations illustrated in figure 6.7. The constraints of the form $C_{i.R_i}$ are existing query constraints, while the

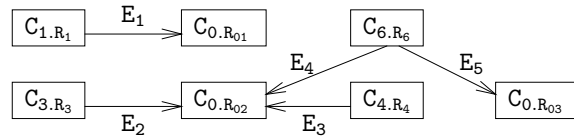


Figure 6.7: Correlations between query constraints and new indexed constraint

constraints $C_{0.R_{0j}}$, $j = 1, \dots, 3$ are new constraints on the cluster index attribute. The additional subscript j is needed because the constraints on the index attribute may have different ranges. We identify all possible constraint introduction solutions as follows:

Algorithm

For each constraint on the index attribute $C_{0.R_{0j}}$

Step 1 find the least expensive association from a query constraint to $C_{0.R_{0j}}$, i.e. find the incoming link $C_{i.R_i} \Rightarrow C_{0.R_{0j}}$ with the fewest exceptions. For example if the cardinalities of E_1, \dots, E_5 are n_1, \dots, n_5 and $n_4 \leq n_2 \leq n_3$, then the least expensive rule for $C_{0.R_{0j}}$ is $C_{6.R_6} \Rightarrow C_{0.R_{02}} (E = E_4)$.

Step 2 filter out all the exceptions that do not satisfy at least one of the constraints $C_{1.R_1}, \dots, C_{12.R_{12}}$. We do not need to test the exceptions for the antecedent of the corresponding rule, since they satisfy it by definition.

The constraint introduction solutions identified in our example are the following:

<i>Introduced Constraints</i>	<i>Exceptions</i>
$C_{0.R_{01}}$	$\{e e \leftarrow E_1, C_{2.R_2}(e), \dots, C_{12.R_{12}}(e)\}$
$C_{0.R_{02}}$	$\{e e \leftarrow E_4, C_{1.R_1}(e), \dots, C_{5.R_5}(e), C_{7.R_7}(e), \dots, C_{12.R_{12}}(e)\}$
$C_{0.R_{03}}$	$\{e e \leftarrow E_5, C_{1.R_1}(e), \dots, C_{5.R_5}(e), C_{7.R_7}(e), \dots, C_{12.R_{12}}(e)\}$

6.5 Optimizing OQL queries

In the previous sections, a series of algorithms were given to find a collection of constraint elimination or constraint introduction solutions. In this section, we show how the original query is transformed to its optimized form using the optimal solution. Consider the OQL query:

```

select x
from Extent_X as x
where C1.R1 and ... and Cn.Rn
  
```

6.5.1 Heuristic H1

Let $\langle \{C_{i1}, \dots, C_{im}\}, E_i \rangle$ be the maintained constraints and the exceptions of a solution i . Only the maintained constraints of the optimization solution should be tested on the objects of the whole extent; however, all the constraints should be tested on the exception cases and the objects that satisfy the maintained constraints but not the ones omitted should be removed from the result. Hence, the original query should be converted to the following one:

```
(select x from Extent_X as x
  where Ci1 and ... and Cim)
  except Ei
```

If we assume that tests on the query constraints take roughly the same time, the optimal solution is the one with fewest exceptions E_i .

6.5.2 Heuristic H2

Let $\langle C_{0.R_{0i}}, E_i \rangle$ be the index constraint and the corresponding exceptions of a constraint introduction solution. Instead of testing the query constraints $C_{1.R_1}, \dots, C_{n.R_n}$ on the entire extent $Extent_X$, we apply them only on the results of the subquery

```
select x
from Extent_X as x      (Q0i)
where C0.R0i
```

Since $C_{0.R_{0i}}$ is a constraint on a cluster index attribute, the select operation is expected to be quite fast. The original query is transformed to its more efficient form:

```
(select x
  from Q0i as x
  where C1.R1 and ... and Cn.Rn)
  union Ei
```

The exceptions E_i are merged to the result because they satisfy the query constraints, but not the new index constraint $C_{0.R_{0i}}$. The optimal solution is the one that introduces the index constraint with the highest selectivity and the least exceptions, i.e. the solution that minimizes $|Q_{0i}| + |E_i|$.

6.6 Discussion

We now look at two different scenarios, and estimate the extent to which heuristics H1 and H2 speed up query execution.

The first scenario concerns frequently executed queries. Assume that the association rules which are used by algorithm 6.3 are not modified. We may optimize a query once at compilation time, then execute its optimized form. It is worth optimizing provided that the execution time of the optimized query is less than the execution time of the original query. For heuristic H1 this happens only if the time saved by omitting some constraints is greater than the time needed to remove the exceptions from the result (except operation). The more the eliminated constraints and the fewer the exceptions, the better the optimization. It is worth pointing out that the time saved by the elimination of constraints is CPU-related. Since query execution is dominated by data access time, this optimization is not expected to alter performance significantly. It would help only in contexts rich in associations with few exceptions, in which users express many constraints in their queries.

Heuristic H2 is expected to bring more significant benefits. Firstly, this optimization involves a union operation, which is much cheaper than the `except` operation used in H1. Secondly, instead of retrieving all the objects of an extent from the database, we need only look at the subset retrieved through an indexed constraint. Hence, we save a considerable amount of data access time, spending a negligible amount of CPU time in evaluating the additional constraint.

The second scenario concerns queries which are executed only once. In this case, the time required for optimization is important. This time depends on the algorithm that finds associations between relaxed constraints and tight constraints (see section 6.3), since this is the most expensive step in the optimization process. This algorithm finds paths in a directed graph, and combines the exceptions associated with each edge of the path to derive the total exceptions for the path. In chapter 7 we propose a number of optimizations for this algorithm and study its complexity in detail.

6.7 Conclusion

The use of association rules for query optimization is relevant to both relational and object-oriented database systems. There has been a lot of research on generating association rules and maintaining them in the presence of updates. Research has also focused on finding heuristics that take advantage of rules in order to optimize a query. Most of this work ([CGM90, GGMR97]) has considered integrity rules, rather than association rules with exceptions. Semantic optimization heuristics were also applied without considering indirect associations. In this chapter, we propose algorithms that apply two optimization heuristics presented by Siegel et al., taking account of both exceptions and indirect associations. We show how to use these heuristics to optimize an OQL query.

The complexity of the optimization process is closely related to the complexity of the constraint graph, which represents the set of association rules in the data. It also depends on the number of exceptions associated with each rule. In chapter 7, we study these dependencies in detail and set up an experimental framework to assess the benefits of the semantic optimization techniques. We

also propose a number of changes to the proposed algorithms that are expected to reduce the cost of semantic optimization significantly.

Chapter 7

Complexity and Optimization

The constraint introduction and elimination heuristics suggest the use of association rules to optimize queries semantically (chapter 6). In this chapter we assess the complexity of the algorithm that navigates over the existing association rules and identifies indirect paths from sources to target constraints along with their exceptions. In particular, we look at its first two steps, since the cost of the third step is negligible. Recall that step 1 backtracks from a target constraint in order to find incomplete paths from sources to targets. The path annotations resulting from step 1 are then combined in step 2 to form (complete) direct or indirect paths along with their exceptions.

Studying the cost of these steps does not only tell us whether it is worth executing the algorithm. It also helps us to locate expensive parts that could possibly be optimized, deferred, or even completely avoided. In this chapter, we propose a number of optimizations that reduce the cost of the algorithm significantly.

Having implemented the optimized version of the algorithm, we set up a number of experiments to measure the benefits and the cost of exploiting semantic knowledge during query execution. In our experimental framework, we use the algorithm in order to add index constraints in the `predicate` of a `Select` algebraic operator. Consequently, we are interested in associations that relate the constraints of a query with any constraint on the existing index attribute. By running a variety of queries and changing the cluster index attribute on which data organisation is based, we manage to monitor the behaviour of the algorithm under different input conditions.

7.1 Experimental Framework

The theoretical analysis of the cost of steps 1 and 2 is supported by experimental results of running OQL queries in the following context. Queries are executed on a Pentium III processor, running at 450 MHz with 256 Mbytes of SDRAM. The OQL processor, which has been developed for the

purposes of this thesis, consists of the following components:

- a type checker that ensures the typability of a query against a database schema [BT00],
- a module that converts OQL to calculus representation (chapter 4),
- a module that normalizes calculus expressions based on a set of rules (chapter 4),
- a module that converts calculus expressions to their algebraic forms (chapter 5),
- a semantic optimizer that takes advantage of direct or indirect association rules (with their exceptions) in order to add or eliminate constraints from a predicate of a `Select` algebraic operator (chapter 6),
- a physical execution module that takes the optimized algebraic plan and executes it.

This OQL processor has been implemented on top of the storage manager of POET OSS 6.1.

The data set to which queries are addressed has been downloaded from the site of the U.S. Census Bureau (Department of Commerce). Each object represents a respondent to a census survey with 48 attributes. We have downloaded the survey data in the period 1973-1994 which consists of approximately 64,000 objects.

We also implemented a rule mining module based on the algorithms by Agrawal et al. [AIS93] and Srikant et al. [SA96]. This module identifies association rules of the form

$$\text{attr}_1 \in R_1 \wedge \dots \wedge \text{attr}_n \in R_n \Rightarrow \text{attr}_0 \in R_0 \text{ (Exc)}$$

That is, for all objects (except those in `Exc`) such that the value of `attr1` is in range `R1` and `...` and the value of `attrn` is in range `Rn`, it is expected that the value of `attr0` be in `R0`. The objects that do not satisfy the rule above, i.e. that satisfy the antecedent condition without satisfying the consequent, are included in the set of exceptions `Exc`.

Each rule has a certain confidence and support, which reflect its strength and popularity in the dataset. When running the rule mining algorithm, we may specify the minimum values of support and confidence for the derived rules. Depending on these values, the rule mining algorithm generates different sets of association rules along with their exceptions. For instance, low values of confidence and support are expected to yield numerous rules with many exceptions on average, and high connectivity of the graph of constraints that represents the association rules.

For each experiment to be presented in the remainder of this chapter, we specify the confidence and the support that were passed as parameters to the rule mining algorithm, as well as the number of resulting association rules. Both the strength and the number of existing association rules affect significantly the cost of semantic optimization of OQL queries, as will become evident in the complexity analysis.

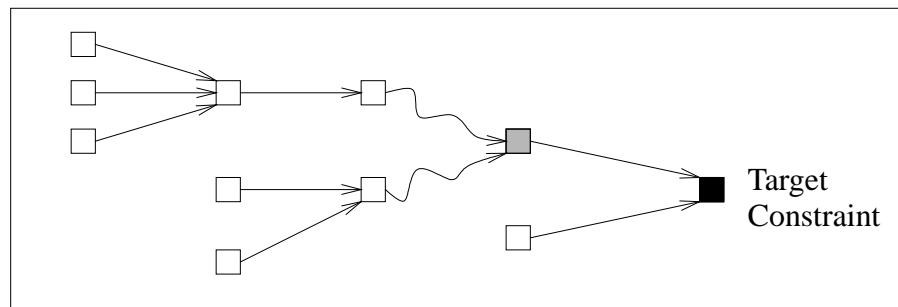


Figure 7.1: Graph of constraints forming a tree towards the target constraint

The rule mining algorithm and the OQL processor have been implemented in Java (JDK 1.2.2). In order to evaluate the cost of the semantic optimizer, we use a Java profiling tool called OptimizeIt. This tool provides the means of monitoring both the memory and the CPU usage of a Java program. We use its CPU profiler to measure the cost of running the OQL processor, paying particular attention to the semantic optimizer. The monitored times are slightly greater than the real times, because of the small overhead imposed by the CPU profiler. Time measurements are also slightly inconsistent in multiple executions of the same query. This is mainly due to the random effect of the garbage collection mechanism of the Java virtual machine and the caching mechanism of the object storage manager. In order to handle this problem, we execute each query at least five times. Hence, each time measurement in the experiments of this chapter is the average of the times monitored during multiple executions of the same query.

Despite the inaccuracy of time measurements the experimental results provide sufficient evidence to support the theoretical claims about the complexity of the semantic optimization algorithm.

7.2 Complexity of step 1

7.2.1 Initial version of step 1

This step backtracks from each target constraint, in order to find all (incomplete) paths from sources to targets. This problem is equivalent to finding all maximum-length acyclic paths leading to the target constraints.

If the graph of constraints had the form of a tree as illustrated in figure 7.1, then step 1 would access the edges of the graph (E) at most once for each target constraint. The target constraints are at most as many as the constraints of the graph (V). Hence, in the worst case, the complexity of this step would simply be $O(V * E)$.

However, in the general case, the backtracking algorithm traverses certain edges of the graph a

number of times. In order to evaluate the cost of step 1, we need to evaluate the maximum number of edge traversals given a graph with V vertices and E edges.

In appendix C we show that given a number of vertices V , the maximum number of edge traversals is given by the following formula:

$$TE = \sum_{i=1}^{V-1} \left(\prod_{j=1}^i (V - j) \right) \quad (7.1)$$

We also show that given a number of edges E , the maximum number of edge traversals is:

$$TE = (m + 1) + (m + 1)(m) + (m + 1)(m)(m - 1) + \dots + (m + 1)! \quad (7.2)$$

where $\lceil m = (-3 + \sqrt{1 + 8E})/2 \rceil$.

The maximum number of edge traversals given both V and E is derived based on formulae 7.1 and 7.2.

$$TE(V, E) \leq \min\{TE_1(V), TE_2(E)\} \quad (7.3)$$

The cost of step 1 is linear in the number of traversed edges (TE). At each edge of the graph of constraints a path annotation of average length ℓ is propagated, i.e. it is copied and extended with one more constraint. Since $\ell \ll TE$, we may deduce that the cost of step 1 is $O(TE)$.

In practice, the number of edges traversed in step 1 (TE) is much smaller than the theoretical maximum given in 7.3. TE does not depend only on the number of edges and the number of vertices; it is very strongly related to the index (target) constraints and their connectivity to the rest of constraints. For instance, consider a graph of constraints with 51 vertices and 223 edges, that corresponds to 186 association rules with minimum confidence 73% and support 31.25%. Using as index the attribute `race` led to step 1 traversing 1,670,103 edges. However, using as index the attribute `spwrks1f` resulted in traversing 2,098,510 edges. Given a set of association rules, the number of TE varies with the selected index, but remains the same for each query execution. Hence, the optimizer can estimate the complexity of step 1 by optimizing a single query, and select a different index attribute in case the cost is not acceptable.

7.2.2 Optimization of Step 1

The great number of traversed edges makes the execution of step 1 prohibitively expensive. The original version of step 1 is inefficient since it may traverse a certain edge of the graph of constraints multiple times (at most as many as the number of path annotations). Another reason is that it only stops the backtracking procedure when it finds a constraint already included in the path or a constraint

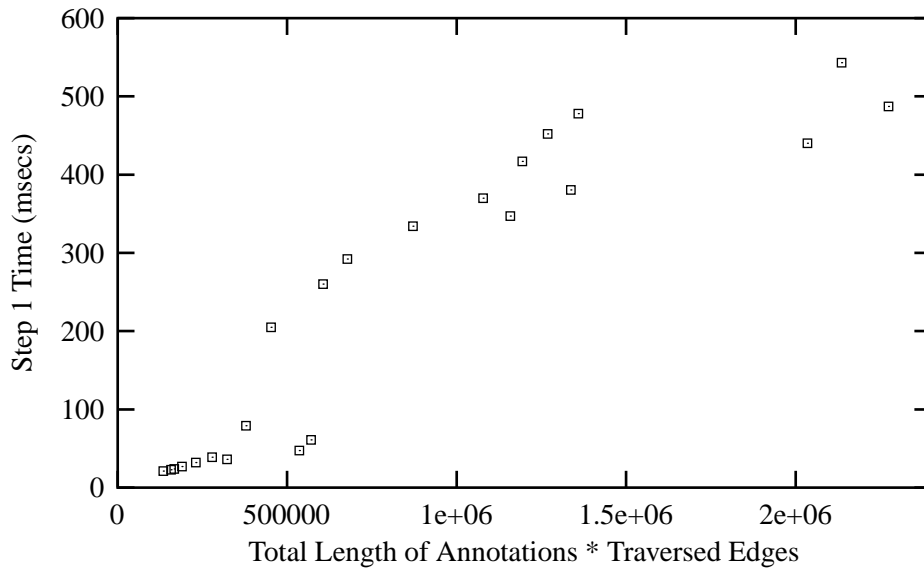


Figure 7.2: Cost of finding incomplete paths from sources to targets in step 1

without any incoming links. However, it would be sufficient to stop when a source constraint is encountered.

Based on these two remarks, we propose an optimized version of step 1. We start from a target constraint C_T and we apply the following steps recursively.

1. Annotate the current constraint C_{cur} as *visited*.
2. If C_{cur} is a source constraint, annotate it with a path consisting only of itself.
3. Otherwise, for each constraint C_i such that there is an edge from C_i to C_{cur} ,
 - if C_i is not visited apply the procedure on it recursively. This results in C_i being annotated with zero or more paths from sources to C_i . (If C_i has already been visited then it has already been annotated with these paths.)
 - Extend each one of these paths with the current constraint C_{cur} and add the new paths to the annotations of C_{cur} .

The algorithm above shows that backtracking is interrupted when a source constraint is encountered. The use of the flag *visited* prevents us from traversing the same edge more than once. The downside is that all constraints in the paths from sources to targets are annotated with paths from sources to themselves. These annotations are first created and stored in the source constraints; they are gradually extended and propagated towards the targets.

Let us consider the cost of the new version of step 1. This does not depend only on the number of traversed edges, but also on the number and the length of path annotations propagated across the graph of constraints. Propagating a number of annotations across an edge implies copying them and extending each one of them with an additional constraint. The cost of copying the annotations at each stage (across a certain edge) is bounded by the total length of path annotations derived as a result of step 1 (TLPA). If TE is the number of edges traversed during the execution of step 1, the cost of the latter should be $O(\text{TLPA} * \text{TE})$. Figure 7.2 shows the cost of the optimized version of step 1 for 23 queries, which are optimized by adding constraints on different index attributes. The graph of constraints is generated by 186 association rules with minimum confidence 73% and support 31.25%. The optimization of the same 23 queries have been used to monitor times relevant to following experiments (figures 7.2, 7.3, 7.4, 7.5, 7.6, 7.7). The optimized version of step 1 is considerably faster than the original one.

7.3 Complexity of Step 2

As a result of the first (backtracking) step the source constraints are annotated with paths leading to the target constraints. Based on these path annotations, the second step of the algorithm performs two operations iteratively. It first absorbs all simple constraints in the paths until a pair of simple-composite constraints is reached. It then finds combinations of annotations having the same composite constraint, such that the preceding simple constraints are all the conjuncts of the composite one. If new annotations result from the combinations they are added to the existing annotations and the two substeps are repeated until no more combinations are generated.

An important issue that should be stressed is the handling of exceptions during the iteration of the two substeps. Substep 1 absorbs the simple constraints in the path annotations, by going through the edges (i.e. the direct associations) that link them. The exception sets associated with each edge are taken into consideration whilst going through the edges. If constraint A leads to constraint B based on a rule with exceptions $E_{A \Rightarrow B}$, and B leads to C based on a rule with exceptions $E_{B \Rightarrow C}$, then the exception set of the indirect path $A \rightarrow C$ (or the indirect rule $A \Rightarrow C$) is:

$$E = \{exc | exc \leftarrow E_{A \Rightarrow B}, \text{not } C(exc)\} \cup \{exc | exc \leftarrow E_{B \Rightarrow C}, A(exc)\} \quad (7.4)$$

Substep 2 combines path annotations that pass from the same composite constraint through different conjuncts. Doing so, it also combines the exception sets evaluated so far for each of the annotations; that is, it creates the union of the exceptions evaluated so far for each combined path that

satisfy the source constraints of the other combined paths. If $C = C_1 \wedge C_2$,

$$\begin{aligned} E(C_A \wedge C_B \Rightarrow C) = \\ \{e \mid e \leftarrow E(C_A \Rightarrow C_1), C_B(e)\} \cup \\ \{e \mid e \leftarrow E(C_B \Rightarrow C_2), C_A(e)\} \end{aligned} \quad (7.5)$$

The original algorithm implies the application of the two formulae (7.4, 7.5) at each iteration of the substeps 2.1 and 2.2.. This results in filtering many exception sets against various constraints, even if the resulting path annotations do not eventually lead to complete paths from a source to a target constraint. It is indicative to say that measuring the time of step 2 in the context of rules with 90 percent confidence, the prevailing cost of both substeps was the filtering of exceptions. This is because filtering implies resolving the references of the exceptions, fetching the actual instances from the database and checking their values against certain constraints.

To avoid this overhead we optimize step 2 by deferring the filtering of exceptions as much as possible. Instead of applying formulae 7.4 and 7.5, we just accumulate the sets of exceptions both while absorbing simple constraints (substep 2.1) and while combining path annotations (substep 2.2). This is a relatively cheap operation since it does not involve resolving exceptions, and therefore accessing the database.

The resolution and filtering of exceptions will take place only for those exceptions that are relevant to the complete paths from sources to target constraints. Hence, annotations that do not form complete paths should not burden the optimization process with exception filtering. Further, the exceptions do not need to be tested against the intermediate constraints of the path annotations, as is implied by the formulae 7.4 and 7.5, but only against the sources and the target constraints.

The optimization of step 2 - by deferring exception filtering and actually doing it only against the sources and the target constraints - led to significant performance improvements. In the remainder of this section, we reason about the complexity of step 2 (in its optimized form) and test its performance by running a series of queries.

7.3.1 Substep 2.1 - Absorbing simple constraints in path annotations

The nature of substep 2.1 denotes that its cost should be linear to the number of absorbed constraints. This is explained by the fact that absorbing each constraint, by traversing a simple edge, involves a constant cost as is explained below. Each annotation is associated with a set of exception sets which is initially empty. Absorbing a constraint involves i) traversing a simple edge (i.e. a direct association rule) and ii) adding the exception set of this edge to the set (of exception sets) associated with the annotation.

The linearity of the relation between the number of absorbed constraints and the cost of step 2.1 is demonstrated in figure 7.3.

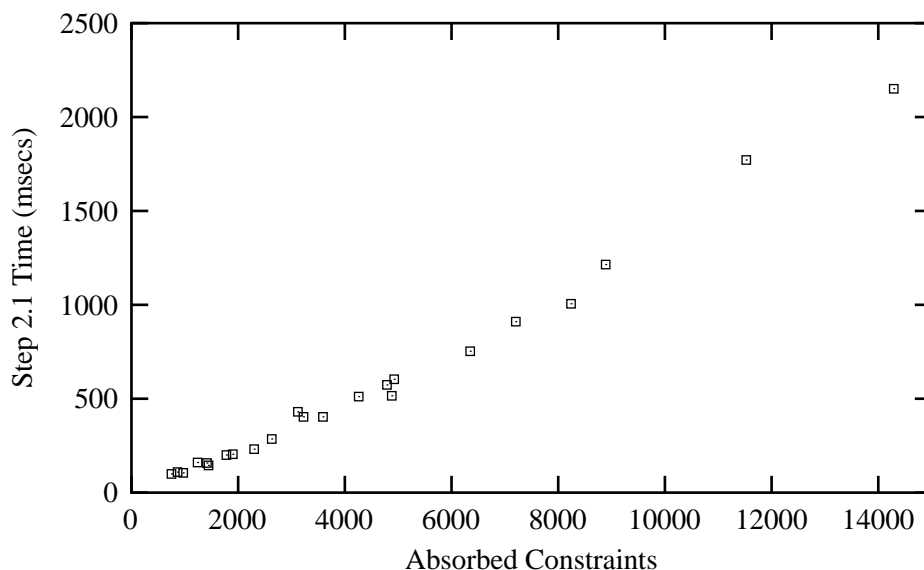


Figure 7.3: Relating the cost of step 2.1 to the number of absorbed constraints

Proposition 1 *The number of absorbed constraints is not bounded by the total length of path annotations.*

Proof. This proposition is easy to prove by the use of a simple example. Consider six path annotations containing the same composite constraint. In two of them the composite constraint is preceded by its first conjunct; the other four annotations reach the composite constraint through its second conjunct.

$$\begin{aligned}
 C_1 &\rightarrow \dots C_{\text{conj}_1} \rightsquigarrow C_{\text{comp}} \rightarrow C_7 \rightarrow C_8 \rightarrow C_9 \\
 C_2 &\rightarrow \dots C_{\text{conj}_1} \rightsquigarrow C_{\text{comp}} \rightarrow C_7 \rightarrow C_8 \rightarrow C_9 \\
 C_3 &\rightarrow \dots C_{\text{conj}_2} \rightsquigarrow C_{\text{comp}} \rightarrow C_7 \rightarrow C_8 \rightarrow C_9 \\
 C_4 &\rightarrow \dots C_{\text{conj}_2} \rightsquigarrow C_{\text{comp}} \rightarrow C_7 \rightarrow C_8 \rightarrow C_9 \\
 C_5 &\rightarrow \dots C_{\text{conj}_2} \rightsquigarrow C_{\text{comp}} \rightarrow C_7 \rightarrow C_8 \rightarrow C_9 \\
 C_6 &\rightarrow \dots C_{\text{conj}_2} \rightsquigarrow C_{\text{comp}} \rightarrow C_7 \rightarrow C_8 \rightarrow C_9
 \end{aligned}$$

Combining these annotation results in eight new annotations. Hence, the edges following the composite constraint will be traversed eight instead of six times. The total traversed edges will obviously be more than the total length of path annotations.

Proposition 2 *The number of absorbed constraints is bounded by the following formula:*

$$(\text{New_Annotations} + \text{Initial_Annotations}) \times \text{Composite_Distance},$$

where *Composite_Distance* is the maximum number of consecutive simple constraints preceding or following a composite constraint.

Proof. `New_Annotations` is the number of paths generated by combining existing annotations in step 2.2. These along with the `Initial_Annotations`, generated by the backtracking algorithm in step 1, are all the annotations to which the absorbing algorithm may apply. The maximum number of constraints absorbed in each annotation is `Composite_Distance`. Hence, the number of absorbed constraints never exceeds

$$(\text{Combined_Annotations} + \text{Initial_Annotations}) \times \text{Composite_Distance}$$

7.3.2 Substep 2.2 - Combining path annotations

Substep 2.2 involves the following costs:

Sorting Cost: The path annotations resulting from substep 2.1 have a simple constraint at their first position and a composite one at their second. It would be convenient to sort these annotations in order to facilitate their combination. The initial implementation of step 2.2 sorted paths based first on the composite and then on the simple constraints. It then found combinations of annotations satisfying two conditions: i) their initial composite constraint should be the same and ii) the simple constraints preceding the common composite constraint should form all its conjuncts. This method had two basic disadvantages. First it resulted in great numbers of new annotations that needed to be sorted and then combined. Second, the new annotations contained many duplicates that needed to be removed. In short, this method made the cost of step 2.2 prohibitively expensive.

A simple observation optimized the implementation of step 2.2 significantly. *A number of paths should be combined only if they are identical from their first composite constraint onwards.* For instance, say we have the following path annotations:

$$\begin{aligned} C_1 &\rightsquigarrow C_{\text{comp}} \rightarrow C_3 \rightarrow C_4 \rightarrow C_5 \\ C_2 &\rightsquigarrow C_{\text{comp}} \rightarrow C_3 \rightarrow C_4 \rightarrow C_5 \\ C_3 &\rightsquigarrow C_{\text{comp}} \rightarrow C_5 \\ C_4 &\rightsquigarrow C_{\text{comp}} \rightarrow C_5 \\ C_1 &\rightsquigarrow C_{\text{comp}} \rightarrow C_6 \rightsquigarrow C'_{\text{comp}} \rightarrow C_5 \\ C_2 &\rightsquigarrow C_{\text{comp}} \rightarrow C_6 \rightsquigarrow C'_{\text{comp}} \rightarrow C_5 \end{aligned}$$

In the original implementation the first annotation would be combined with the second, the third and the sixth, since they have the same first composite constraint and different simple constraints (conjuncts) preceding the composite one. However in the current implementation, the first annotation is combined only with the second, because it has the same body from the first composite constraint onwards. This method of combining annotations yields a smaller number of new annotations, without limiting the ability to find indirect association rules. It also avoids the generation of duplicate new annotations.

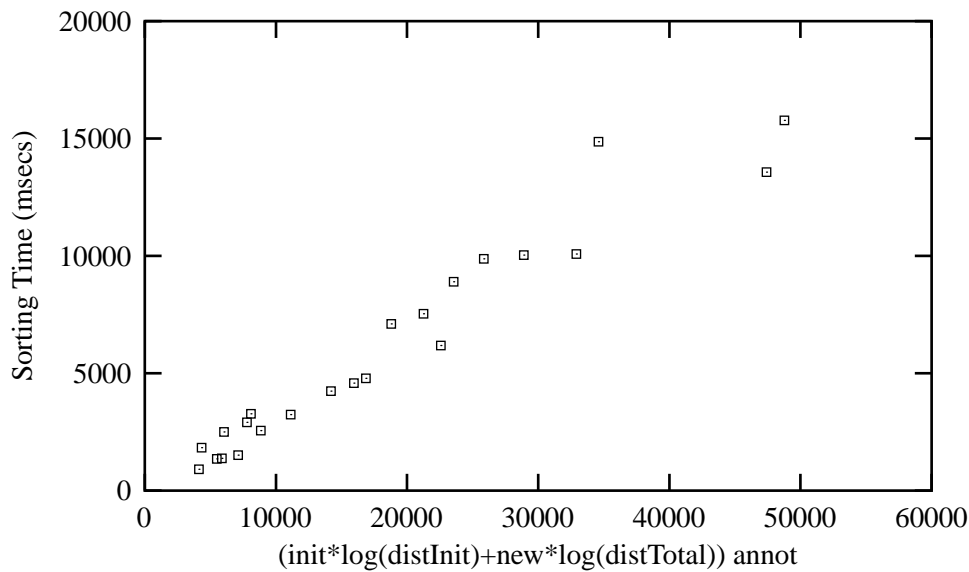


Figure 7.4: Cost of sorting initial and new path annotations

In order to facilitate combining paths, we sort them based on:

1. the number of composite constraints in each of them (desc)
2. their body from the first composite constraint onwards (asc)
3. the simple constraint preceding the first composite constraint (asc)

The paths of the previous example would be sorted as follows:

$$\begin{aligned}
 C_1 &\rightsquigarrow C_{\text{comp}} \rightarrow C_6 \rightsquigarrow C'_{\text{comp}} \rightarrow C_5 \\
 C_2 &\rightsquigarrow C_{\text{comp}} \rightarrow C_6 \rightsquigarrow C'_{\text{comp}} \rightarrow C_5 \\
 C_1 &\rightsquigarrow C_{\text{comp}} \rightarrow C_3 \rightarrow C_4 \rightarrow C_5 \\
 C_2 &\rightsquigarrow C_{\text{comp}} \rightarrow C_3 \rightarrow C_4 \rightarrow C_5 \\
 C_1 &\rightsquigarrow C_{\text{comp}} \rightarrow C_5 \\
 C_2 &\rightsquigarrow C_{\text{comp}} \rightarrow C_5
 \end{aligned}$$

Consider the general case in which Step 1 yields IA initial annotations. At the first iteration, step 2.1 absorbs the single constraints preceding the first composite constraint of each path. The resulting annotations are inserted into a sorted vector and duplicate paths are eliminated. This results in a vector of DA distinct ordered annotations. The cost of inserting/sorting the initial annotations is $O(\text{IA} * \log(\text{DA}))$. Step 2.2 then combines these annotations into new ones that are also inserted in the

ordered vector. The total number of annotations in the final sorted vector is $TA = DA + NA$, where NA are the new annotations generated during all iterations of step 2.2. The cost of inserting/sorting the new annotations is therefore $O(NA * \log(TA))$.

The total sorting cost of step 2.2 is $O(IA * \log(DA) + NA * \log(TA))$ (fig 7.4). We must take into account that comparing path annotations, in order to sort them, resembles comparing strings. The maximum number of operations per comparison is bounded by the maximum length of path annotations. If this is a constant $k \ll IA * \log(DA) + NA * \log(TA)$ then we may still assume that

$$\text{Sorting_Time} = O(IA * \log(DA) + NA * \log(TA)).$$

In order to predict the sorting cost of step 2, we need to evaluate IA, DA, NA and TA . The number of initial annotations (IA) is determined at the end of step 1, whilst the number of distinct annotations (DA) is known by the end of the first iteration of step 2.1. If we managed to predict the number of new annotations (NA), then the total annotations would be $TA = DA + NA$.

Assume that after sorting the initial path annotations, we have

- p_1 annotations containing ℓ_1 composite constraints,
- p_2 annotations containing ℓ_2 composite constraints,
- ...
- p_n annotations containing ℓ_n composite constraints,

where $DA = p_1 + \dots + p_n$.

We assume for simplicity that composite constraints consist of two conjuncts. If s is the number of source constraints, then each composite constraint can be reached by at most $2^s - 1$ combinations of sources. The minimum number of paths that can generate $2^s - 1$ combinations is $2 * \sqrt{2^s - 1}$ ($\sqrt{2^s - 1}$ paths reach the composite constraint through each one of its conjuncts). Hence, by combining p_1 annotations, we derive at most $(p_1 * \sqrt{2^s - 1})/2$ new paths.

If $\ell_2 = \ell_1 - 1$ then the new paths $(p_1 * \sqrt{2^s - 1})/2$ are added to p_2 . Combining the annotations with ℓ_2 composite constraints ($p'_2 = p_2 + (p_1 * \sqrt{2^s - 1})/2$) results in at most $(p'_2 * \sqrt{2^s - 1})/2$ new paths.

If $\ell_2 < \ell_1 - 1$ then only $p'_1 = (p_1 * \sqrt{2^s - 1})/2$ paths are combined into $(p'_1 * \sqrt{2^s - 1})/2$ new annotations. These have $\ell_1 - 2$ composite constraints.

This procedure continues in a similar way. If $\ell_1 = \ell_2 + 1 = \dots = \ell_n + n - 1$ then

$$NA = \sum_{i=1}^n \text{New_Annot}_i \tag{7.6}$$

where $\text{New_Annot}_i = (p'_i * \sqrt{2^s - 1})/2$ and $p'_i = p_i + \text{New_Annot}_{i-1}$.

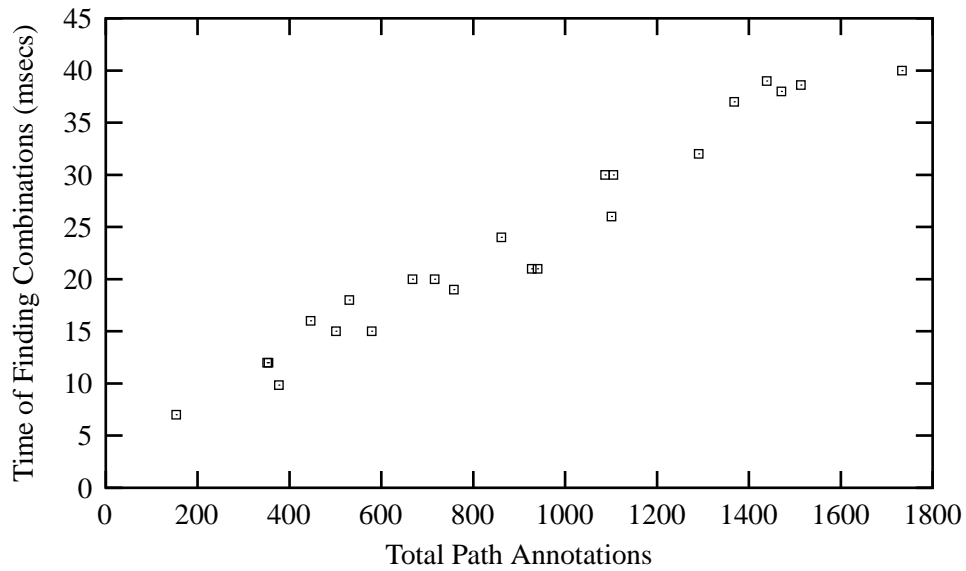


Figure 7.5: Cost of finding new combinations in step 2.2

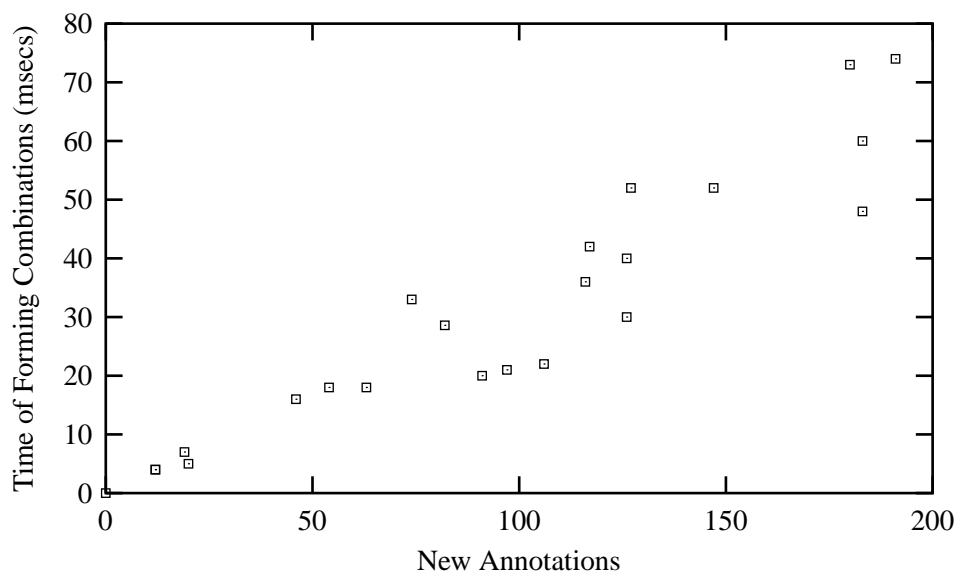


Figure 7.6: Cost of forming combinations in step 2.2

Cost of Combining Annotations: This cost involves iterating over the ordered vector of paths in order to combine annotations with the same initial composite constraint, different conjuncts preceding it and the same body from the composite constraint onwards. This is easy to do since the vector is already sorted based on these three criteria.

Consider a composite constraint as having n positions, one for each conjunct. By browsing the sorted vector of annotations, we can linearly determine which paths can fill the first, second, or n^{th} position of the composite constraint. Groups of n annotations that can fill all different positions of a composite constraint are combined to form a new path annotation. Hence the cost of finding combinations is linear in the total number of path annotations TA (figure 7.5).

The total number of combinations equals the number of new annotations. For each combination the sources of the constituent paths are checked that they have compatible ranges. The exceptions of these paths are then merged to form the exceptions of the new path. The cost of forming combinations is therefore linear in the number of new annotations NA (figure 7.6).

The total cost of combining sorted annotations is given by a function of the form:

$$\text{Combining_Cost} = a * \text{TA} + b * \text{NA} + c, \quad a, b, c > 0$$

Based on this equation, we can predict the worst-case cost of combining annotations by predicting NA and setting $\text{TA} = \text{DA} + \text{NA}$.

7.4 Combining steps 1 and 2 optimizes the algorithm significantly

Extending and propagating path annotations across the graph of constraints is the most expensive operation of step 1. Likewise, sorting path annotations in order to facilitate their combination is the dominant cost of step 2. It is possible to avoid both of these expensive operations by combining the functionality of steps 1 and 2.

In the new optimized algorithm, we avoid storing in each constraint C_{cur} (step 1) long paths from source constraints. Instead we store pairs of the form (C_s, SE) , which means that there is a complete (indirect) path from the source constraint C_s to C_{cur} that involves the sets of exceptions $E_i \in \text{SE}$. Here $E_i, i = 1, \dots, n$ correspond to the sets of exceptions of the direct edges (rules) that form the path $C_s \rightarrow \dots \rightarrow C_{\text{cur}}$. Based on this remark, we present the combined version of steps 1 and 2.

We start from a target constraint C_T and we apply the following steps recursively:

1. Annotate the current constraint C_{cur} as *visited*.
2. If C_{cur} is a source constraint, annotate it with the pair $(C_{\text{cur}}, \{\})$.
3. Otherwise, if C_{cur} is a simple constraint (not a conjunction of constraints)

- for each constraint C_i such that there is an edge from C_i to C_{cur} ,
 - If C_i is not visited apply the procedure to it recursively. This results in C_i being annotated with zero or more pairs of the form C_s, SE (if C_i has already been visited then it has already been annotated with these pairs).
 - Let E_i be the set of exceptions of the association rule $C_i \rightarrow C_{\text{cur}}$.
 - Take copies of the annotations $\{(C_j, SE_j)\}$, $j = 1, \dots, m$ of C_i and update their sets of exceptions with E_i : $SE'_j = SE_j \cup \{E_i\}$, $j = 1, \dots, m$. Add the updated copies $\{(C_j, SE'_j)\}$ to the annotations of C_{cur} .
- Remove duplicate annotations of C_{cur}

4. Otherwise, (if C_{cur} is a composite constraint)

- for each conjunct C_i , $i = 1, \dots, n$ such that $C_i \rightsquigarrow C_{\text{cur}}$,
 - If C_i is not visited apply the procedure to it recursively. This results in C_i being annotated with zero or more pairs of the form C_s, SE (if C_i has already been visited then it has already been annotated with these pairs).
- Combine the annotations of conjuncts C_i to form new annotations $\{(C'_j, SE'_j)\}$ that reflect complete paths that pass through all conjuncts of C_{cur} . For instance, if each conjunct C_i is annotated with a single pair (C_{s_i}, SE_i) , these pairs are combined to the new annotation $(C_{s_1} \wedge \dots \wedge C_{s_n}, SE_1 \cup \dots \cup SE_n)$.
- Remove duplicate annotations of C_{cur}

The algorithm above propagates pairs of the form $(\text{source}, \text{exception_sets})$ instead of long path annotations that include all intermediate constraints. There is also no need to sort path annotations before combining them, since combinations are formed on the fly at the composite constraint, as soon as the indirect paths from sources to its conjuncts have been identified. The cost of the algorithm depends on the number of edges traversed in the graph of constraints (TE). Note that TE is bounded by the total number of edges E in the graph. The length of path annotations no longer plays an important role (as in the previous version of step 1), since for each indirect path we store only its source and target constraint. The number of tuples $(\text{source}, \text{exception_sets})$ propagated across each edge (or combined in composite constraints) is very small compared to the number of traversed edges. These tuples represent complete paths from sources to a particular constraint and are considerably fewer than the incomplete paths propagated across the graph in the previous version of the algorithm.

To be more specific, if s is the number of source constraints, the maximum number of annotations propagated across an edge (or stored in a constraint) is $2^s - 1$. For instance, given two sources C_1

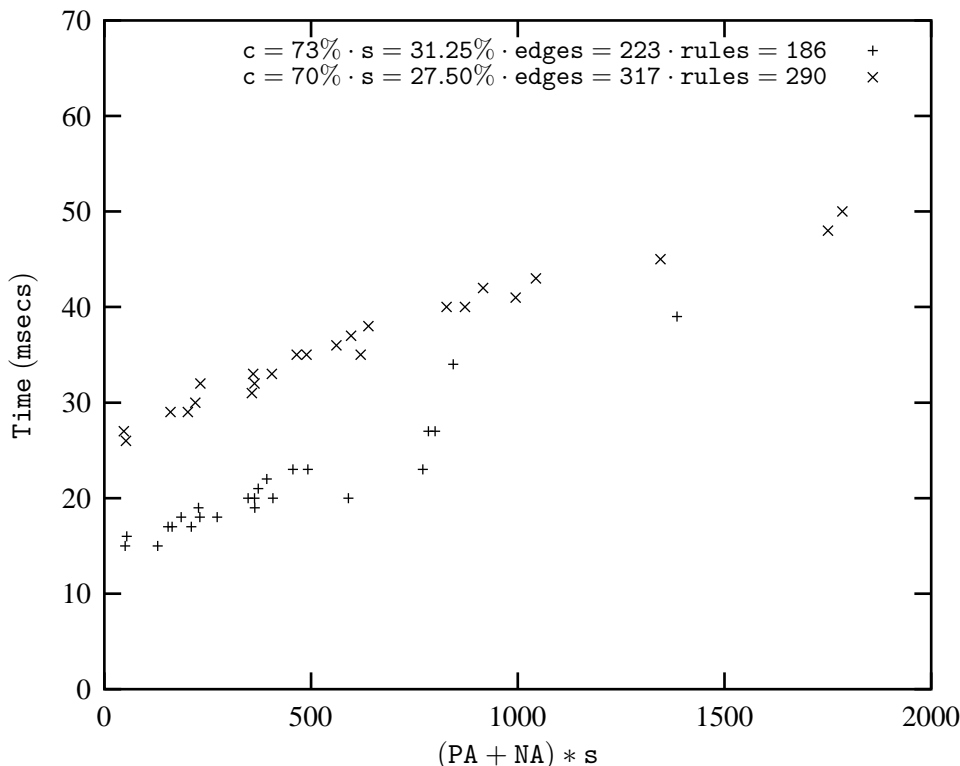


Figure 7.7: Cost of optimized algorithm

and C_2 , the annotations of any other constraint can be at most $2^2 - 1 = 3$:

$$(C_1, E_1) \quad (C_2, E_2) \quad (C_1 \wedge C_2, E_3)$$

Consider PA as the total number of annotations propagated towards simple (non-composite) constraints. Assume that NA is the number of new annotations that are formed in composite constraints by combining the annotations of the conjuncts. For each one of the PA + NA annotations, we evaluate the corresponding sets of exceptions. Then we add each tuple (C, E) to the annotations of a constraint ensuring that we eliminate duplicates. Since the maximum number of annotations per constraint is $2^s - 1$, the cost of forming annotations and removing duplicate ones is $O((PA + NA) * \log(2^s - 1)) = O((PA + NA) * s)$ (figure 7.7).

We monitored the cost of the algorithm by optimizing 23 queries (the same ones used in previous experiments) using first 186 association rules and then 290 rules (figure 7.7). The first set of rules was generated by running the rule mining algorithm with confidence 73% and support 31.25%. The second set was generated based on 70% confidence and 27.5% support. Optimizing queries using 186 association rules is cheaper than using 290 rules, since the generated graph of constraints has

fewer edges in the first case (223 edges) than in the second (317 edges). This implies that there are more complete paths from source constraints to a certain constraint C , i.e. C has more annotations. The cost of removing duplicate annotations, and hence the total cost of the algorithm, is increased. This explains why optimizing queries with the same value for $(PA + NA) * s$ is more expensive when the number of association rules, and thus the number of edges, is increased. The factor s provides an upper bound for the cost of removing duplicates at each constraint, but in practice this factor depends on the connectivity of the graph.

If E is the number of edges in the graph of constraints then the total number of annotations propagated towards simple constraints is:

$$PA \leq E * (2^s - 1)$$

If the graph contains $comp$ composite constraints, each one having at most $conj$ conjuncts, then the total number of annotations generated in composite constraints is:

$$NA \leq comp * (2^s - 1)^{conj}$$

Hence, the complexity of the new optimized algorithm is the following:

$$Cost = O((E * (2^s - 1) + comp * (2^s - 1)^{conj}) * s) \quad (7.7)$$

The number of sources s reflects the number of predicates in the where clause of the initial query. Hence, formula 7.7 shows that the complexity of the query affects the cost of optimization exponentially. The number of edges E , composite constraints $comp$ and conjuncts per composite constraint $conj$ reflect the connectivity of the graph of constraints. Hence, the optimization cost also depends on the number and the nature of association rules which are represented in the graph of constraints. For a particular query, we can bound the cost of the algorithm by selecting an appropriate set of association rules as the basis for the optimization.

7.5 Overall performance of optimizing queries semantically

In this chapter, we studied the complexity of the algorithm that navigates over a graph of constraints and finds direct or indirect paths along with their exceptions. Having implemented various versions of this algorithm, we used it in order to optimize queries semantically. Our intention was to identify association rules that would help us add index constraints in the where clauses of the queries. A series of time measurements on the various (sub)steps of the algorithm were carried out in order to shed light on the analysis of complexity.

Studying the complexity of the algorithm led us to a number of optimizations that improved its performance significantly. In the final version of the algorithm, the cost of steps 1 and 2 is negligible,

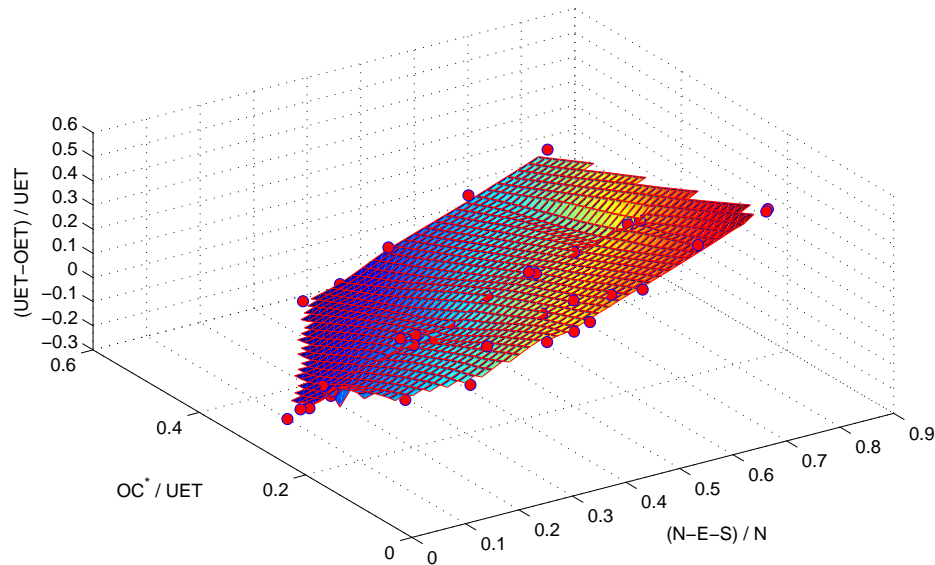


Figure 7.8: Relating the query execution benefits to the optimization cost and the evaluation benefits (top view)

compared to the operations that fetch objects from the database. These operations include i) fetching the association rules in order to generate the graph of constraints and ii) filtering the exceptions involved in the (indirect) paths that contribute to the query optimization. The cost of the first operation is ignored, since the OQL processor can load the association rules into a graph of constraints once in order to avoid burdening the execution of each individual query. Since the cost of resolving and filtering exceptions is linear in the number of exceptions, the optimizer decides whether to optimize a query semantically, as soon as it evaluates the number of exceptions E_i corresponding to each useful indirect association i . The consequent (index) constraint of an association i can be used to retrieve S_i instead of N objects from the extent. The optimal association is the one that minimizes $E_i + S_i$. The optimizer applies this association to transform a query semantically only provided that $N > E_{opt} + S_{opt}$. Otherwise, the optimization algorithm is interrupted and the query is executed in its initial form.

The surface in figure 7.8 demonstrates the benefits of the algorithm in the optimization of OQL queries. The x-axis $((N - E - S)/N)$ represents the percentage of objects that we avoided retrieving from the extent by introducing the index constraint. It shows the benefits of the optimization solution at evaluation time. The y-axis represents the fraction of the optimization cost to the unoptimized total execution time (OC^*/UET). Note that the optimization cost (OC^*) does not include the cost of resolving exceptions from the database. It only involves loading the association rules into

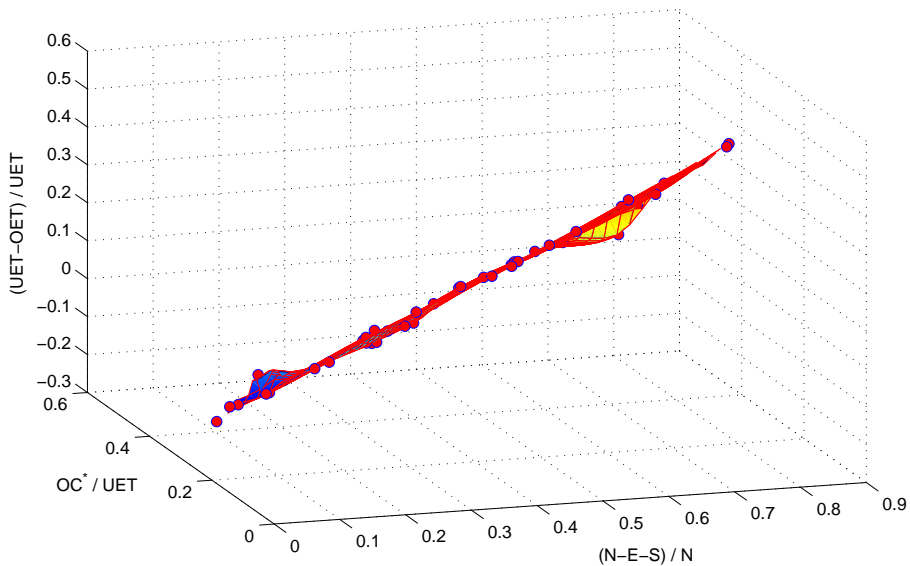


Figure 7.9: Relating the query execution benefits to the optimization cost and the evaluation benefits (1st side view)

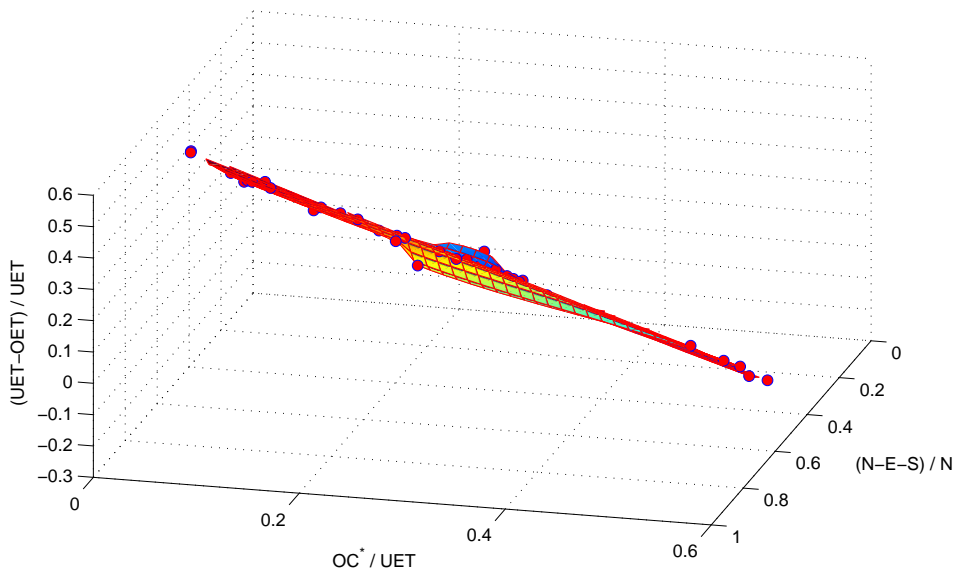


Figure 7.10: Relating the query execution benefits to the optimization cost and the evaluation benefits (2nd side view)

a graph of constraints and finding the constraint introduction solution. OC^* should be distinguished from the optimization cost measured in figure 7.7, which does not include the cost of loading rules in memory. The z-axis represents the percentage of the total execution time saved by the optimization $(UET - OET)/UET$. It is evident that as we decrease the cost factor (OC^*/UET) and increase the evaluation benefit $(N - E - S)/N$, the total benefit of the optimization $(UET - OET)/UET$ is increased.

We optimized 40 queries using a variety of index constraints with different selectivities. We also used different sets of association rules by running the rule mining algorithm with the following values of confidence and support $\{(c = 50\%, s = 10\%), (c = 60\%, s = 10\%), (c = 70\%, s = 18.75\%), (c = 73\%, s = 25\%), (c = 73\%, s = 31.25\%)\}$. Although queries were evaluated over a small extent (4000 objects) the benefits of applying the constraint introduction heuristic were often significant. For large extents the cost factor OC^*/UET is expected to be negligible and the total execution benefit to be linear in the evaluation benefit $((N - E - S)/N)$. This is also the case for queries that are optimized once at compile time, and are then executed multiple times.

The side views of the same surface (figures 7.9 and 7.10) show that its shape is almost planar, as is expected theoretically:

$$\begin{aligned} \text{Total_Optimization_Benefit} &= \text{Evaluation_Benefit} - \text{Semantic_Optimization_Cost} \\ \Rightarrow UET - OET &= a * (N - E - S) - OC^* \\ \Rightarrow (UET - OET)/UET &= a * (N - E - S)/UET - OC^*/UET \\ \Rightarrow (UET - OET)/UET &= a * (N - E - S)/(a * N + b) - OC^*/UET \end{aligned}$$

where a reflects the time of fetching and processing a single object and b reflects the cost of checking the type of a query, translating it into its intermediate representations and applying syntactic transformations. If b is much smaller than the evaluation time ($a * N$), we may assume that $UET \approx a * N$. Hence,

$$(UET - OET)/UET \approx (N - E - S)/N - OC^*/UET$$

The anomalies of the surface in figures 7.8, 7.9 and 7.10 are partly caused by the inaccuracy of time measurements, as well as by the assumption $UET \approx a * N$. The surface is the result of interpolating the irregularly sampled data to a regularly spaced grid and then plotting the interpolated and the non-uniform data (using MATLAB).

There are cases when the optimization has a negative effect on the query execution time. This may be due to either of the following reasons: i) the optimizer may have decided that the optimization is not worth applying, but cannot undo the optimization cost incurred so far; ii) the optimization was worth applying, but its benefits were not greater than the optimization cost. In both cases the negative effect is negligible, provided that the optimization cost does not include fetching the association rules from the database at each query execution.

7.6 Discussion

In chapter 6, we proposed an algorithm that navigates through a set of existing association rules, and identifies indirect associations between certain source and target constraints. This algorithm was proposed as part of the overall algorithm for optimizing queries using the constraint introduction (or elimination) heuristic. In this chapter, we reasoned about the complexity of the algorithm by looking at its functionality step by step. The complexity analysis highlighted aspects of the algorithm that needed to be improved and hence contributed to its optimization. Having implemented various versions of the algorithm we tested its performance by running a series of OQL queries. The benefits of applying the constraint introduction heuristic depend on the selectivity of the index constraint and the exceptions involved in the useful associations. Our experience is that this heuristic often reduces the time of query execution significantly. The constraint elimination heuristic is worth applying only when the associations involved have practically no exceptions.

Chapter 8

Conclusion and Future Work

In this dissertation we developed an OQL processing framework that encompasses all features of the object-oriented paradigm and is flexible enough to benefit from semantic knowledge in the form of association rules. This was achieved by extending previous work on type inference, calculus and algebraic intermediate representations, as well as on the use of semantic optimization heuristics. The proposed OQL framework is useful for a number of commercial or scientific applications that handle large amounts of complex data rich in semantic correlations (chapter 2). We give a summary of the contributions of this dissertation in the various stages of OQL processing.

8.1 Summary

We first proposed a type inference model (chapter 3) in order to ensure that an OQL query is well-typed. The inference algorithm, which is the core of this model, identifies the most general type of an OQL query in the absence of schema information; it also derives the minimum type requirements that a schema should satisfy to be compatible with this query. The expressiveness and the complex type system of OQL raised issues that could not be handled by previous work on type inference. Hence, by developing a type inference model for OQL, we gained a lot of experience in dealing with polymorphic collection types, structures and classes all being related through a complex subtype relationship and being manipulated using a variety of constructs. It was shown that the unification mechanism used to resolve constraints in previous models is not sufficient in our context. The inference rules defined in this thesis generate a variety of constraints specific to OQL and a mechanism is provided to resolve them. The proposed inference model should be useful in any distributed or even semi-structured database environment.

Having ensured the typability of an OQL query, we then investigated its translation into the monoid calculus intermediate representation (chapter 4). In the translation process, we identified

several weak points of ODMG OQL [CBB⁺00], e.g. the nondeterministic behaviour of the `order by` clause, and proposed ways of dealing with them. We pointed out that not all OQL queries reflect homomorphic mappings, and hence, their translation using monoid comprehensions is not always possible. Nonhomomorphic functions like `bagof`, `listof`, `arrayof` and `sortof` were introduced in order to resolve this problem. The resulting translation rules, which cover all OQL constructs, differ from those defined in related contexts [Afs98, FM00]. In [FM00], a set of normalization rules is defined in order to simplify calculus expressions. In this dissertation, we extended this set with additional rules and we proposed a normalization algorithm that determines the order of their application.

We then investigated the translation of calculus expressions to the monoid algebra (chapter 5). Following the approach of Fegaras and Maier [FM00], we adopted a set of algebraic operators (e.g. `Select`, `Join`) similar to the ones used in relational contexts. However, we realised that new operators should be introduced in order to translate calculus expressions containing nonhomomorphic functions. We also pointed out that nested comprehensions should be processed locally before being merged with the algebraic expressions generated in outer comprehensions. Our approach provides opportunities for further optimization compared to the approach in [FM00]. New algebraic operators were added in order to handle merge operations of collections. We highlighted the need for the `OuterSelect` operator, which enables the OQL processor to filter the results of a nested comprehension before processing its head. At the physical level, this implies reducing the volume of data to be processed at an early stage. Finally, by annotating all algebraic operators with monoid information, we provided opportunities for optimizing the physical execution plan. The translation rules proposed in this thesis provide a complete mapping from the monoid calculus to the algebraic query representation.

In chapter 6, we proposed an algorithm that uses association rules in order to optimize queries semantically. This algorithm is the main step towards applying the constraint introduction and elimination heuristics presented in [SSS92]. It is assumed that a rule mining algorithm is used in order to generate a set of association rules along with their exceptions. The number of derived rules depends on the minimum support and confidence that are passed as parameters. The exceptions to the rules, i.e. the objects that satisfy the antecedent but not the consequent, are maintained properly in the presence of database updates. Given a particular query, we showed how query predicates are mapped into source or target constraints found in association rules. We then proposed an algorithm that identifies direct or indirect correlations between source and target constraints along with their exceptions. Using these correlations, a set of optimization solutions is identified, and the optimal one is selected for the query transformation. Previous work on semantic optimization has mainly considered integrity rules rather than association rules with exceptions [CGM90, GGMR97]. There has also been no previous effort to identify indirect correlations that could be useful for semantic transformations. The

algorithm proposed in this thesis has two valuable properties: i) it is flexible enough to take advantage of weak or indirect correlations and ii) it can be used to apply a variety of heuristics besides the ones that suggest the introduction or elimination of constraints. For instance, the proposed algorithm is equally useful for the following two heuristics proposed by Siegel et al. [SSS92]:

- If a query contains a restriction on a range attribute, then try to locate a constraint on that attribute that is a contradiction to the existing restriction (contradiction heuristic).
- If a restriction on a range attribute implies the answer to a query then the latter can be answered without accessing the database (tautology heuristic).

In order to assess the benefits of semantic optimization, we evaluated the cost of the algorithm that finds indirect correlations between source and target constraints. The complexity of the algorithm was estimated from a theoretical point of view and was confirmed by optimizing a number of queries using the constraint introduction heuristic. The complexity analysis led us to identify expensive tasks of the algorithm that could be changed, deferred or avoided. A number of optimizations were proposed that reduced the cost of the algorithm considerably.

Our experience from implementing the algorithm showed that a significant amount of time could be saved if the association rules are loaded in memory (into a graph of constraints) once and not every time a query is optimized. Given a graph of constraints, the semantic optimizer navigates over its edges and identifies indirect associations that could be used to apply the constraint introduction heuristic. We decided that the exceptions involved in the indirect associations should be accumulated but not resolved, until the OQL processor decides on whether to apply the optimization. A query is optimized, if the sum of exceptions and objects retrieved using the index constraint is smaller than the cardinality of the extent. Hence, instead of resolving all objects from the original extent, the OQL processor only fetches a subset of them, and thus saves significant data-access time. If the optimization is not applicable, the query execution is burdened with a negligible CPU-related cost. It is shown that this cost is bounded by the number of edges in the graph of constraints.

The constraint elimination heuristic is worth applying when the number of exceptions involved in the optimization solution is very small. In a relational context, the cost of resolving these exceptions may easily outbalance the CPU-related benefits of applying the heuristic. However, in an object-oriented context, where constraints may potentially include expensive method invocations, the optimizer may decide to apply the heuristic despite the great number of exceptions involved.

In the case of the tautology and contradiction heuristics, a potential optimization solution is beneficial independent of the number of exceptions involved. The optimized queries can be executed by resolving only the exceptions of the solution, which are always fewer than the cardinality of the extent to which the query is addressed.

8.2 Future work

The areas studied in this dissertation include type inference of OQL, calculus representation, normalization of calculus expressions, algebraic representation and semantic optimization using association rules. This work can be extended in the last two phases of OQL processing, namely the algebraic level and the semantic optimization techniques.

In particular, in this dissertation we defined a complete mapping from calculus expressions to the monoid algebra. However, we did not consider syntactic transformations at the algebraic level. It would be interesting to identify a set of rules that optimize the algebraic plan, taking into account the new algebraic operators introduced in this thesis. Further work could focus on optimizing the physical execution plan based on the monoid information associated with the algebraic operators.

Future research might also concern the semantic optimization of OQL queries. In this dissertation we presented in detail all the steps of applying the constraint introduction and elimination heuristics (chapter 6). The most important step is an algorithm that navigates over a graph of constraints and finds direct or indirect associations along with their exceptions. The heuristics considered so far are applicable to both relational and object-oriented systems. Future work could investigate the use of the proposed algorithm to apply heuristics specific to the object-oriented paradigm. For instance, the query optimizer could benefit from association rules with the following kinds of constraints:

$$\begin{array}{l}
 C \equiv x_1 \in \text{extent} \\
 | \quad x_1.\text{relationship} = x_2 \\
 | \quad x_1 \in x_2.\text{relationship} \\
 | \quad x_1.\text{attribute}_1 < \text{comp_op} > x_2.\text{attribute}_2 \\
 | \quad x_1.\text{method}(y_1, \dots, y_k) < \text{comp_op} > x_2.\text{attr} \\
 | \quad x_1.\text{method}(y_1, \dots, y_k) < \text{comp_op} > \text{value}
 \end{array}$$

where x_1, x_2 represent objects and y_1, \dots, y_k may be instances of any OQL type (chapter 3). It is assumed that in the beginning of a rule, each variable x_i (or y_i) is annotated with the most general type for which the rule is applicable.

Examples of rules containing constraints on the attributes, methods, relationships or extents of one or more objects are given below:

$$\lambda x : \text{Employee.}$$

$$x.\text{salary} > 50,000 \Rightarrow x \in \text{Managers}$$

$$\lambda x : \text{Company.}$$

$$x.\text{revenues} > x.\text{expenses} \Rightarrow x.\text{profits}() > 0$$

$$\lambda x : \text{Person}, y : \text{Person.}$$

$$x \in \text{Unemployed} \wedge y \in \text{Employees} \Rightarrow x.\text{expenses} < y.\text{expenses}$$

$$\lambda x : \text{Hotel}, \text{month} : \text{integer.}$$

$$x.\text{stars} = 3 \wedge \text{month} < 9 \wedge \text{month} > 5 \Rightarrow x.\text{revenues}(\text{month}) > 200,000$$

In order to benefit from such semantic knowledge, new data mining algorithms should be devised to identify correlations between the values of attributes, relationships and methods of objects belonging to the same or different classes. Mechanisms like the use of join indices should be considered in order to optimize the execution of mining algorithms [ST95, JS97]. Standard procedures for defining semantics about method invocations (e.g. relating parameters to their results or restricting the range of the results) are also worth studying.

Assuming that we have the means of deriving *complex* association rules, the next issue that arises is how we could possibly generate a graph of constraints that reflects these associations. For instance, information about the class hierarchy of a particular schema is necessary in order to add a link from a constraint $x \in \text{extent}_1$ towards $x \in \text{extent}_2$, where $\text{extent}_1 \subseteq \text{extent}_2$ (class_1 corresponding to extent_1 is a subclass of class_2 corresponding to extent_2).

Association rules have traditionally been used to define *positive* correlations between the values of data. For instance, if the values of an attribute A are in a range R_1 the values of attribute B are in range R_2 . However, query optimization could take advantage of association rules that contain negations and thus express the fact that certain ranges of attributes never appear together [EGLM01]. It would be interesting to study how we could identify indirect associations with exceptions in the presence of such rules.

We might wish to define indirect relationships between objects, for example, if $x.\text{father} = y$ and $y.\text{father} = z$ then $x.\text{grandfather} = z$. This is a typical rule that expresses intensional information in a deductive database. It could be implemented as a deductive method attached to a class in an object-oriented system [Sam00, SP00]. Hence, a node which represents the constraint $x.\text{grandfather} = z$ might be an abstraction of a subgraph which needs to be traversed. The nodes of the subgraph would be variables instead of constraints, whilst the edges would denote the relationships traversed. Combining known traversal techniques for the subgraphs [Pul96] with the traversal

techniques proposed in this thesis, we could take advantage of associations involving intentional and extensional information.

Given that the previous problems are solved, the proposed algorithm navigates over the graph of constraints and finds paths from source to target constraints. Identifying which constraints are used as sources and targets is another open problem, the answer to which depends on the optimization heuristic.

Finally, future research could consider combining data mining with warehousing techniques. If we were able to identify association rules that involve aggregate values (stored in a warehouse), then it would be possible to optimize expensive `select . . . group by` queries [AHLS00].

8.3 Conclusion

This dissertation presents a framework for the efficient execution of OQL queries. Starting from inferring the type of an OQL query, we proceed to translate it into a calculus form, normalize it at this level and translate it into an algebraic representation. We then demonstrate how OQL queries can be optimized semantically with the aid of association rules. Many optimization heuristics, e.g. the constraint introduction and elimination heuristics, suggest the use of semantic knowledge expressed in (direct) integrity rules. We point out that optimization opportunities would be increased if we used indirect associations that are not necessarily 100% valid. An algorithm is proposed that navigates over existing association rules and identifies indirect associations along with their exceptions. We assess the cost and the benefits of this algorithm in the context of two optimization heuristics applicable in both object-oriented and relational systems. It is expected that the algorithm be equally useful for applying heuristics specific to the object-oriented paradigm. We implemented and tested the ideas reported in this thesis by developing an OQL engine and enhancing it with a semantic optimizer.

Appendix A

Inference Rules

This appendix presents the inference rules for OQL that have not been discussed in chapter 3 of the dissertation:

$$\frac{\mathcal{H}; \mathcal{C} \vdash q_1: \sigma \Rightarrow \mathcal{C}_1 \quad \mathcal{H}; \mathcal{C}_1 \vdash q_2: \tau \Rightarrow \mathcal{C}_2}{\mathcal{H}; \mathcal{C} \vdash q_1 \text{ and } q_2: \text{bool} \Rightarrow \mathcal{C}_2 \wedge \{\sigma = \text{bool}\} \wedge \{\tau = \text{bool}\}}$$
$$\frac{\mathcal{H}; \mathcal{C} \vdash q_1: \sigma_1 \Rightarrow \mathcal{C}_1 \quad \mathcal{H}; \mathcal{C}_1 \vdash q_2: \sigma_2 \Rightarrow \mathcal{C}_2}{\mathcal{H}; \mathcal{C} \vdash q_1 \text{ union } q_2: \sigma \Rightarrow \mathcal{C}_2 \wedge \{\sigma = \text{Merge_Result}(\sigma_1, \sigma_2)\} \wedge \{\text{Equality_Compatible}(\text{Member_Type}(\sigma_1), \text{Member_Type}(\sigma_2))\}}$$
$$\frac{\mathcal{H}; \mathcal{C} \vdash q_1: \sigma_1 \Rightarrow \mathcal{C}_1 \quad \mathcal{H}; \mathcal{C}_1 \vdash q_2: \sigma_2 \Rightarrow \mathcal{C}_2}{\mathcal{H}; \mathcal{C} \vdash q_1 = q_2: \text{bool} \Rightarrow \mathcal{C}_2 \wedge \{\text{Equality_Compatible}(\sigma_1, \sigma_2)\}}$$
$$\frac{\mathcal{H}; \mathcal{C} \vdash q_1: \sigma_1 \Rightarrow \mathcal{C}_1 \quad \mathcal{H}; \mathcal{C}_1 \vdash q_2: \sigma_2 \Rightarrow \mathcal{C}_2}{\mathcal{H}; \mathcal{C} \vdash q_1 < q_2: \text{bool} \Rightarrow \mathcal{C}_2 \wedge \{\text{Greater_Less_Than_Compatible}(\sigma_1, \sigma_2)\}}$$
$$\frac{\mathcal{H}; \mathcal{C} \vdash q_1: \sigma_1 \Rightarrow \mathcal{C}_1 \quad \mathcal{H}; \mathcal{C}_1 \vdash q_2: \sigma_2 \Rightarrow \mathcal{C}_2}{\mathcal{H}; \mathcal{C} \vdash q_1 + q_2: \sigma_0 \Rightarrow \mathcal{C}_2 \wedge \{\sigma_0 = \text{Arith_Result}(\sigma_1, \sigma_2)\}}$$
$$\frac{\mathcal{H}; \mathcal{C} \vdash q_1: \sigma_1 \Rightarrow \mathcal{C}_1 \quad \mathcal{H}; \mathcal{C}_1 \vdash q_2: \sigma_2 \Rightarrow \mathcal{C}_2}{\mathcal{H}; \mathcal{C} \vdash q_1 \text{ mod } q_2: \text{int} \Rightarrow \mathcal{C}_2 \wedge \{\sigma_1 = \text{int}\} \wedge \{\sigma_2 = \text{int}\}}$$

$$\begin{array}{c}
\frac{\mathcal{H}; \mathcal{C} \vdash q: \sigma \Rightarrow \mathcal{C}_1}{\mathcal{H}; \mathcal{C} \vdash \neg q: \sigma \Rightarrow \mathcal{C}_1 \wedge \{\sigma \subseteq \text{int/float}\}} \\
\frac{\mathcal{H}; \mathcal{C} \vdash q: \sigma \Rightarrow \mathcal{C}_1}{\mathcal{H}; \mathcal{C} \vdash \text{not } q: \text{bool} \Rightarrow \mathcal{C}_1 \wedge \{\sigma = \text{bool}\}} \\
\frac{\mathcal{H}; \mathcal{C} \vdash q: \sigma \Rightarrow \mathcal{C}_1}{\mathcal{H}; \mathcal{C} \vdash \text{abs}(q): \sigma \Rightarrow \mathcal{C}_1 \wedge \{\sigma \subseteq \text{int/float}\}} \\
\frac{\mathcal{H}; \mathcal{C} \vdash q: \sigma \Rightarrow \mathcal{C}_1}{\mathcal{H}; \mathcal{C} \vdash \text{min}(q): \phi \Rightarrow \mathcal{C}_1 \wedge \{\phi \subseteq \text{orderable}\} \wedge \{\text{Member.Type}(\sigma) = \phi\}} \\
\frac{\mathcal{H}; \mathcal{C} \vdash q: \sigma \Rightarrow \mathcal{C}_1}{\mathcal{H}; \mathcal{C} \vdash \text{sum}(q): \phi \Rightarrow \mathcal{C}_1 \wedge \{\phi \subseteq \text{int/float}\} \wedge \{\text{Member.Type}(\sigma) = \phi\}} \\
\frac{\mathcal{H}; \mathcal{C} \vdash q: \sigma \Rightarrow \mathcal{C}_1}{\mathcal{H}; \mathcal{C} \vdash \text{avg}(q): \text{float} \Rightarrow \mathcal{C}_1 \wedge \{\sigma \subseteq \text{collection}(\text{int/float})\}} \\
\frac{\mathcal{H}; \mathcal{C} \vdash q: \sigma \Rightarrow \mathcal{C}_1}{\mathcal{H}; \mathcal{C} \vdash \text{count}(q): \text{int} \Rightarrow \mathcal{C}_1 \wedge \{\sigma \subseteq \text{collection}(\text{nonfunctional})\}} \\
\frac{\mathcal{H}; \mathcal{C} \vdash q: \sigma \Rightarrow \mathcal{C}_1}{\mathcal{H}; \mathcal{C} \vdash \text{distinct}(q): \phi \Rightarrow \mathcal{C}_1 \wedge \{\phi = \text{Distinct.Result}(\sigma)\}} \\
\frac{\mathcal{H}; \mathcal{C} \vdash q: \sigma \Rightarrow \mathcal{C}_1}{\mathcal{H}; \mathcal{C} \vdash \text{listtoiset}(q): \text{set}(\phi) \Rightarrow \mathcal{C}_1 \wedge \{\sigma = \text{list}(\phi)\}} \\
\frac{\mathcal{H}; \mathcal{C} \vdash q: \sigma \Rightarrow \mathcal{C}_1}{\mathcal{H}; \mathcal{C} \vdash \text{element}(q): \phi \Rightarrow \mathcal{C}_1 \wedge \{\text{Member.Type}(\sigma) = \phi\}}
\end{array}$$

Appendix B

Translating OQL to the monoid calculus

This appendix presents the full mapping from OQL constructs to their calculus forms. The notation $|q|^*$ stands for the calculus equivalent of an OQL query q ($| - |^* : \text{OQL} \rightarrow \text{Calculus}$).

OQL Expressions	Calculus Expressions
id	id
literal	literal
q.id	$ q ^*.id$
$q(q_1, \dots, q_n)$	$ q ^*(q_1 ^*, \dots, q_n ^*)$
$q_1 \text{ union } q_2$	$\text{merge}[M](q_1 ^*, q_2 ^*)$, (to determine M see 4.2.3)
$q_1 \text{ intersect } q_2,$ ($q_1 \text{ or } q_2 : \text{bag}(-)$)	$\text{bag}\{x \parallel x \leftarrow \text{bagof}(q_1 ^*), \text{some}\{x = y \parallel$ $y \leftarrow \text{bagof}(q_2 ^*)\}\}$
$q_1 \text{ intersect } q_2,$ ($q_1 \text{ and } q_2 : \text{set}(-)$)	$\text{set}\{x \parallel x \leftarrow q_1 ^*, \text{some}\{x = y \parallel$ $y \leftarrow q_2 ^*\}\}$
$q_1 \text{ except } q_2$ ($q_1 \text{ or } q_2 : \text{bag}(-)$)	$\text{bag}\{x \parallel x \leftarrow \text{bagof}(q_1 ^*), \text{all}\{x \neq y \parallel$ $y \leftarrow \text{bagof}(q_2 ^*)\}\}$
$q_1 \text{ except } q_2$ ($q_1 \text{ and } q_2 : \text{set}(-)$)	$\text{set}\{x \parallel x \leftarrow q_1 ^*, \text{all}\{x \neq y \parallel y \leftarrow q_2 ^*\}\}$
first(q)	$ q ^*[0]$
last(q)	$ q ^*[\text{sum}\{1 \parallel x \leftarrow \text{bagof}(q ^*)\} - 1]$
$q_1[q_2]$	$ q_1 ^* [q_2 ^*]$

OQL Expressions	Calculus Expressions
$\text{set}(q_1, \dots, q_n)$	$\text{merge}[\text{set}](\text{unit}[\text{set}](q_1 ^*), \dots, \text{unit}[\text{set}](q_n ^*))$
$\text{bag}(q_1, \dots, q_n)$	$\text{merge}[\text{bag}](\text{unit}[\text{bag}](q_1 ^*), \dots, \text{unit}[\text{bag}](q_n ^*))$
$\text{list}(q_1, \dots, q_n)$	$\text{merge}[\text{list}](\text{unit}[\text{list}](q_1 ^*), \dots, \text{unit}[\text{list}](q_n ^*))$
$\text{array}(q_1, \dots, q_n)$	$\text{arrayof}(\text{merge}[\text{list}](\text{unit}[\text{list}](q_1 ^*), \dots, \text{unit}[\text{list}](q_n ^*)))$
$\text{dictionary}((q_1, q'_1), \dots, (q_n, q'_n))$	$\text{merge}[\text{dict}](\text{unit}[\text{dict}](q_1 ^*, q'_1 ^*), \dots, \text{unit}[\text{dict}](q_n ^*, q'_n ^*))$
$\text{flatten}(q, (q : \text{col}_1(\text{set}(-)))$	$\text{set}\{x \parallel s \leftarrow \text{bagof}(q ^*), x \leftarrow s\}$
$\text{flatten}(q, (q : \text{col}_1(\text{bag}(-)))$	$\text{bag}\{x \parallel s \leftarrow \text{bagof}(q ^*), x \leftarrow s\}$
$\text{flatten}(q, (q : \text{set}(\text{list}/\text{oset}/\text{array}(-)))$	$\text{set}\{x \parallel s \leftarrow q ^*, x \leftarrow \text{bagof}(s)\}$
$\text{flatten}(q, (q : \text{bag}(\text{list}/\text{oset}/\text{array}(-)))$	$\text{bag}\{x \parallel s \leftarrow q ^*, x \leftarrow \text{bagof}(s)\}$
$\text{flatten}(q, (q : \text{list}/\text{oset}(\text{oset}(-)))$	$\text{oset}\{x \parallel s \leftarrow q ^*, x \leftarrow s\}$
$\text{flatten}(q, (q : \text{array}(\text{oset}(-)))$	$\text{oset}\{x \parallel s \leftarrow \text{listof}(q ^*), x \leftarrow s\}$
$\text{flatten}(q, (q : \text{list}(\text{list}(-)))$	$\text{list}\{x \parallel s \leftarrow q ^*, x \leftarrow s\}$
$\text{flatten}(q, (q : \text{oset}/\text{array}(\text{list}(-)))$	$\text{bag}\{x \parallel s \leftarrow \text{list}(q ^*), x \leftarrow s\}$
$\text{flatten}(q, (q : \text{list}/\text{oset}(\text{array}(-)))$	$\text{array}\{x \parallel s \leftarrow \text{arrayof}(q ^*), x \leftarrow \text{bagof}(s)\}$
$\text{flatten}(q, (q : \text{array}(\text{array}(-)))$	$\{x \parallel s \leftarrow q ^*, x \leftarrow s\}$
$\text{order_flatten}(q, q : \text{list}(\text{bag}(-)))$	$\text{list}\{x \parallel s \leftarrow q ^*, x \leftarrow \text{listof}(s)\}$
$\text{order_flatten}(q, q : \text{oset}(\text{bag}(-)))$	$\text{oset}\{x \parallel s \leftarrow q ^*, x \leftarrow \text{listof}(s)\}$
$\text{distinct}(q) (q : \text{set}/\text{bag}(-))$	$\text{set}\{x \parallel x \leftarrow q ^*\}$
$\text{distinct}(q) (q : \text{oset}/\text{list}(-))$	$\text{oset}\{x \parallel x \leftarrow q ^*\}$
$\text{distinct}(q) (q : \text{array}(-))$	$\text{oset}\{x \parallel x \leftarrow \text{listof}(q ^*)\}$
$q_1 < \text{comp_op} > q_2$ ($< \text{comp_op} > \in \{<, >, \leq, \geq, =, \neq\}$)	$ q_1 ^* < \text{comp_op} > q_2 ^*$
$q_1 < \text{comp_op} > \text{all } q_2$ ($< \text{comp_op} > \in \{<, >, \leq, \geq, =, \neq\}$)	$\text{all}\{ q_1 ^* < \text{comp_op} > x_2 \parallel x_2 \leftarrow \text{bagof}(q_2 ^*)\}$
$q_1 < \text{comp_op} > \text{any/some } q_2$ ($< \text{comp_op} > \in \{<, >, \leq, \geq, =, \neq\}$)	$\text{some}\{ q_1 ^* < \text{comp_op} > x_2 \parallel x_2 \leftarrow \text{bagof}(q_2 ^*)\}$
$q_1 < \text{arithm_op} > q_2$ ($< \text{arithm_op} > \in \{+, -, *, /, \text{mod}\}$)	$ q_1 ^* < \text{arithm_op} > q_2 ^*$

OQL Expressions	Calculus Expressions
$\neg q$	$\neg q ^*$
not (q)	not ($ q ^*$)
q_1 and q_2	merge[all]($ q_1 ^*$, $ q_2 ^*$)
q_1 or q_2	merge[some]($ q_1 ^*$, $ q_2 ^*$)
element(q)	elementof($ q ^*$)
listtoset(q)	set{x x \leftarrow q}
max(q)	max{x x \leftarrow bagof($ q ^*$)}
min(q)	min{x x \leftarrow bagof($ q ^*$)}
avg(q)	unit[avg](< avg_sum = sum{x x \leftarrow bagof($ q ^*$)}, avg_count = sum{1 x \leftarrow bagof($ q ^*$)} >}
count(q)	sum{1 x \leftarrow bagof($ q ^*$)}
sum(q)	sum{x x \leftarrow bagof($ q ^*$)}
q_1 in q_2	some{x = $ q_1 ^*$ x \leftarrow bagof($ q_2 ^*$)}
for all x in $q_1 : q_2$	all{ $q_2 ^*$ x \leftarrow bagof($ q_1 ^*$)}
exists x in $q_1 : q_2$	some{ $q_2 ^*$ x \leftarrow bagof($ q_1 ^*$)}
exists(q)	some{true x \leftarrow $ q ^*$ }
unique(q)	sum{1 x \leftarrow bagof($ q ^*$)} = 1
struct($\ell_1 : q_1, \dots, \ell_n : q_n$)	struct($\ell_1 = q_1 ^*, \dots, \ell_n = q_n ^*$)
C($\ell_1 : q_1, \dots, \ell_n : q_n$)	C($\ell_1 = q_1 ^*, \dots, \ell_n = q_n ^*$)
select [distinct] q' from q_1 as x_1, \dots, q_k as x_k where q''	bag [set]{ $ q' ^*$ $x_1 \leftarrow$ bagof($ q_1 ^*$), ..., $x_k \leftarrow$ bagof($ q_k ^*$), $ q'' ^*$ }
select [distinct] q' from q_1 as x_1, \dots, q_k as x_k where q'' order by q_1''' [asc/desc], ..., q_j''' [asc/desc]	list [oset]{ v_1 .partition $v_1 \leftarrow$ orderof(group{ struct(asc/desc $_1 = q_1''' ^*, \dots, \text{asc/desc}_j = q_j''' ^*$, partition = unit[bag]($ q' ^*$) $x_1 \leftarrow$ bagof($ q_1 ^*$), ..., $x_k \leftarrow$ bagof($ q_k ^*$), $ q'' ^*$)}}
select [distinct] q' from q_1 as x_1, \dots, q_k as x_k where q'' group by $\ell_1 : q_1''', \dots, \ell_j : q_j'''$ having q'''' order by q_1'''' [asc/desc], ..., q_n'''' [asc/desc]	list [oset]{ v_1 .partition $v_1 \leftarrow$ orderof(group{ struct(asc/desc $_1 = q_1'''' ^*, \dots, \text{asc/desc}_n = q_n'''' ^*$, partition = unit[bag](v_2) $v_2 \leftarrow$ bagof(group{ struct($\ell_1 = q_1''' ^*, \dots, \ell_j = q_j''' ^*$, partition = unit[bag](struct($x_1 = x_1, \dots, x_k = x_k$) $x_1 \leftarrow$ bagof($ q_1 ^*$), ..., $x_k \leftarrow$ bagof($ q_k ^*$), $ q'' ^*$)), $\ell_1 \equiv v_2.\ell_1$, ..., $\ell_j \equiv v_2.\ell_j$, partition $\equiv v_2$.partition, $ q'''' ^*$)}}

Appendix C

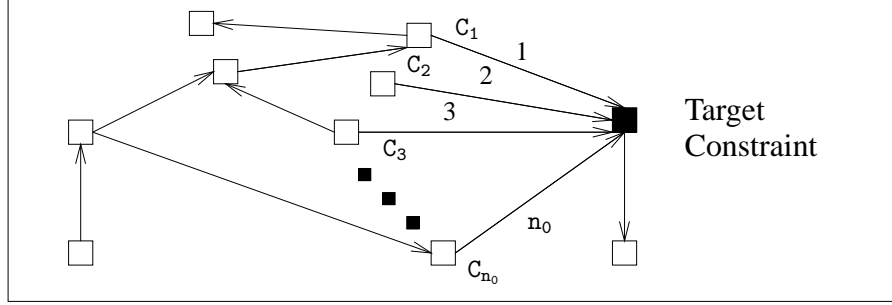
Edge traversals in the backtracking algorithm

C.1 Edge traversals given a graph with V vertices

Given a number of vertices V , the more connected the graph, the more the edge traversals in step 1. Hence, the worst case scenario occurs when the graph is fully connected. For a directed graph, this implies that $E = V * (V - 1)$. Each constraint has $V - 1$ incoming links from all other constraints and $V - 1$ outgoing edges towards them.

Assume that backtracking starts from a certain target constraint. Backtracking can be implemented both as a depth-first-search (DFS) and as a breadth-first-search (BFS) method, both having the same complexity. Starting from the target constraint, the first level of BFS traverses $V - 1$ edges. The paths traversed so far consist of two constraints or vertices. The second time, each one of the $V - 1$ vertices traverses $V - 2$ edges, thus forming acyclic paths containing three vertices. This procedure continues until the $V - 1$ level of BFS, that is, until the traversed paths have V constraints. The table below shows the number of edges traversed at each level.

Level	TraversedEdges
1	$(V - 1)$
2	$(V - 1)(V - 2)$
3	$(V - 1)(V - 2)(V - 3)$
4	$(V - 1)(V - 2)(V - 3)(V - 4)$
...	...
$n = V - 1$	$(V - 1)(V - 2)(V - 3) \dots (V - (V - 1))$

Figure C.1: The target constraint has n_0 incoming links

The total number of traversed edges is:

$$\begin{aligned} \text{TE} &= \sum_{i=1}^{v-1} [(v-1)(v-2) \dots (v-i)] \\ &\Rightarrow \text{TE} = \sum_{i=1}^{v-1} \left(\prod_{j=1}^i (v-j) \right) \end{aligned} \quad (\text{C.1})$$

C.2 Edge traversals given a graph with E edges

Let us consider the backtracking algorithm for a target constraint C_0 . Assume that C_0 has n_0 incoming links, as illustrated in figure C.1. The maximum number of edges traversed backwards starting from C_0 , in a graph of E edges, is denoted as $\text{TE}(C_0, E)$, or simply $\text{TE}(E)$. It is the sum of the following edges:

- n_0 (the incoming links to C_0)
- $\sum_{i=1}^{n_0} \text{TE}(E - n_0)$ (the traversed edges from the C_i , $i = 1, \dots, n_0$ backwards).

Notice that the traversed edges from the constraints C_i , $i = 1, \dots, n_0$ do not consider the initial n_0 edges. The latter will never be traversed again, since they are incoming links to a vertex that has already been encountered. We refer to C_i , $i = 1, \dots, n_0$ as first-level constraints, because they would be encountered at the first level in a breadth-first-search implementation of the algorithm.

Assume that the first-level constraints have on average n_1 incoming links. It can be easily inferred that

$$\text{TE}(E) = n_0 + \sum_{i_0=1}^{n_0} (\text{TE}(E - n_0)) = n_0 + \sum_{i_0=1}^{n_0} (n_1 + \sum_{i_1=1}^{n_1} (\text{TE}(E - n_0 - n_1))) \Rightarrow$$

$$\begin{aligned} \text{TE} &= n_0 + \sum_{i_0=1}^{n_0} (n_1) + \dots + \sum_{i_0=1}^{n_0} \sum_{i_1=1}^{n_1} \dots \sum_{i_{m-1}=1}^{n_{m-1}} (n_m) + \\ &\quad + \sum_{i_0=1}^{n_0} \sum_{i_1=1}^{n_1} \dots \sum_{i_m=1}^{n_m} (\text{TE}(E - n_0 - \dots - n_m)) \end{aligned}$$

where the last term $\sum_{i_0=1}^{n_0} \sum_{i_1=1}^{n_1} \dots \sum_{i_m=1}^{n_m} (\text{TE}(E - n_0 - \dots - n_m)) = 0$. Hence,

$$\text{TE} = n_0 + n_0 * n_1 + n_0 * n_1 * n_2 + \dots + n_0 * \dots * n_m \tag{C.2}$$

$$E = n_0 + n_1 + \dots + n_m \tag{C.3}$$

The value of TE in formula C.2 is maximized when $n_0 > n_1 > \dots > n_m$ and in particular when $n_i = n_{i-1} + 1$. If we set $n = n_0$, then the formulae C.2, C.3 take the following forms:

$$\text{TE} = n + n(n - 1) + \dots + n(n - 1) \dots (n - m) \tag{C.4}$$

$$E = n(m + 1) - m(m + 1)/2 \tag{C.5}$$

The question that arises at this point is which is the combination of n and m that maximizes the bound TE on the number of traversed edges. The example in the following table shows that TE is maximized when m is maximized.

Edges = 15						
n_0	n_1	n_2	n_3	n_4	$\sum n_i$	TE
5	4	3	2	1	15	325
6	5	4			15	156
8	7				15	64

Hence, TE is maximized if $n = m + 1$. This enables the following transformations:

$$\begin{aligned} E &= (m + 1)(m + 1) - m(m + 1)/2 \Rightarrow \\ m^2 + 3m + 2(1 - E) &= 0 \Rightarrow \\ m &= (-3 + \sqrt{1 + 8E})/2 \end{aligned}$$

The formula that calculates the traversed edges becomes:

$$\text{TE} = (m + 1) + (m + 1)(m) + (m + 1)(m)(m - 1) + \dots + (m + 1)! \tag{C.6}$$

where $m = \lceil (-3 + \sqrt{1 + 8E})/2 \rceil$.

For instance, when $E = 3$, $m = (-3 + \sqrt{1 + 24})/2 = 1$ and $\text{TE} = 2 + 2 * 1 = 4$. In fact, when $E = 3$ the maximum number of traversed edges is $3 \leq 4$.

When $E = 6$, $m = (-3 + \sqrt{1 + 48})/2 = 2$ and $\text{TE} = 3 + 3 * 2 + 3 * 2 * 1 = 15$. Indeed, when $E = 6$ the maximum number of traversed edges is $9 \leq 15$.

Notice that it is wrong to assume that the maximum number of traversed edges occurs when the edges form a fully connected graph. If this was true, then we would calculate the number of vertices using the formula $V(V - 1) = E$; given V , the maximum traversed edges would be evaluated as in section C.1.

Appendix D

Implementation of OQL Processor

This appendix gives an overview of the OQL processor developed to test the ideas and to perform the experiments discussed in this thesis.

An OQL interpreter was developed that performs the following tasks:

- lexical analysis
- parsing
- type checking and type inference
- translation of OQL into calculus representation
- normalization at the calculus level
- translation of calculus into algebraic representation
- semantic optimization
- query evaluation

The query interpreter uses the object storage manager and the Java binding of POET OSS 6.1.

The first step towards developing an OQL interpreter is to write a specification file for OQL terminals (OQL keywords, atomic literals and variables). This file is used by JLex, a lexical analyzer generator for Java, in order to produce a class that takes a query and whenever invoked, returns the next terminal encountered.

A set of Java classes are then defined corresponding to the various OQL constructs. For instance, the classes `Select_Query`, `Ref_Access_Query`, `Index_Access_Query` correspond to the `select`

clause, `e.label` and `e1[e2]` respectively. The class `Main_Query` corresponds to the outer OQL expression, which potentially includes query definitions. All classes inherit from the general class `Query`.

A second specification file includes the grammar of OQL, i.e. the query expressions that are syntactically acceptable. Each OQL expression is paired with the action that should be taken when the expression is encountered. CUP, an LALR Parser Generator for Java, takes as input the OQL grammar specification and produces a parser for it. The generated parser takes a query string and, with the aid of the lexical analyzer, it generates a `Query` object for each query subexpression encountered. When the whole query is parsed a `Main_Query` object is created and its method `process` is invoked. This method checks the type of the query, translates it to calculus form, normalizes it, translates it to algebraic form, optimizes it semantically and evaluates its result.

All subclasses of `Query`, including `Main_Query`, define the methods `checkType`, `inferTypes` and `toMonoidCalculus`. Their implementation varies depending on the specific OQL construct that they represent.

In the `process` method of `Main_Query`, we first invoke the method `checkType`. In order to derive the type of a query, we need to identify the types of its operands; hence, `checkType` is recursively called on its subqueries. Access to the database schema is eventually required when `checkType` is invoked on objects that represent identifiers that are not range variables. Since the Java binding of POET OSS 6.1 does not provide an interface to the schema information, we needed to define and use a method `getType` for all classes in the schema. The method `inferTypes` is also called recursively, but it is used only when the input OQL expression includes query definitions with untyped parameters. Only the types of the parameters are inferred, whilst the remaining identifiers are type-checked in the usual way.

After ensuring that the input query is well-typed, the method `toMonoidCalculus` is invoked in order to translate the query into its monoid calculus form. The latter is invoked recursively on the subqueries (operands or attributes), in a similar way as in the case of `checkType`. For each monoid calculus expression, we define a corresponding class, whose name starts with `MC_`, e.g. `MC_Comprehension`, `MC_Ref_Access`, `MC_Identifier`. All `MC_` classes inherit from the general monoid calculus class `MC_Query`. Hence, the method `toMonoidCalculus` returns an object of type `MC_Query`, referred to hereafter as the calculus object.

All subclasses of `MC_Query` implement the method `normalize`, based on the normalization algorithm defined in chapter 4. They also define a method `toALG` for the translation of calculus constructs to algebraic representation (chapter 5). The implementation of both methods depends on the calculus construct (subclass).

We first invoke the method `normalize` on the calculus object and then the method `toALG` on the normalized calculus object. For each algebraic expression we define a corresponding class whose

name starts with ALG_, e.g. ALG_Selection, ALG_Nest. All algebraic classes inherit from the general class ALG_Query. Hence, the method `toALG` returns an object of type ALG_Query, referred to as the algebraic object.

A method that performs semantic optimization is defined in the class ALG_Query and has been implemented to apply the constraint introduction or elimination heuristic whenever possible. Hence, it is only implemented in subclasses of ALG_Query that represent algebraic expressions with predicates, e.g. ALG_Selection. This method uses the association rules that refer to the operand of the algebraic object and finds direct or indirect associations (with exceptions) in order to add a constraint on an index attribute, or to eliminate any redundant non-index constraints.

We have implemented a rule mining module based on the algorithms by Agrawal et al. [AIS93] and Srikant et al. [SA96]. This module identifies association rules corresponding to a specific extent, which are stored in the database along with their exceptions. The input values of the rule mining algorithm are the minimum support and the minimum confidence required for the resulting rules. In order to find rules between quantitative attributes of an extent, we first apply a clustering algorithm on their values. The Birch algorithm [ZRL96] has been implemented to distribute the values of an attribute into a number of clusters. For each quantitative attribute we define a mirror attribute in order to store cluster identifiers. Hence, instead of finding association rules between individual values by looking at the quantitative attributes, we discover rules relating clusters. That is, we run the rule mining algorithm over the corresponding mirror attributes that hold cluster information.

After an algebraic query has been optimized semantically, it is evaluated using the Java binding to access the POET object database. All algebraic operators (queries) implement the following method.

```
public void exec (ALG_Basket initiator_basket, Thread initiator_exec_thread)
{
    this.start_resource_input ();
    this.consume_resources (initiator_basket, initiator_exec_thread);
}
```

The `initiator_basket` is the structure where the current operator is expected to put its results. This structure is initialized by the parent operator. The first method (`start_resource_input`) initializes the baskets of the operands; the latter may then start populating these baskets with objects that the current operator will use as input. The next method (`consume_resources`) is invoked on the current operator in order to start processing the objects retrieved from the baskets of the operands. The output of this processing is accumulated into the `initiator_basket`.

If the operand of an algebraic query is an extent, a database name or an index constraint on an extent, the method `consume_resources` accesses the database using the Java binding of POET. In case the operand is an extent, an iterator is started and the method `Iterator.next()` is used to fetch the next object from the database. Indices are implemented as objects containing a list of

values and a list of objects. Hence, accessing an extent using an index implies resolving the index object using the method `Transaction.lookup(index_name)` and then invoking on it the method `Index.select(range)`, that fetches from the database all objects for which the value of the index attribute is in the given range. Named objects are also resolved from the database using the built-in method `Transaction.lookup(object_name)`.

The implementation of the method `consume_resources` depends on the algebraic construct (select, join, project, nest etc). Although the focus of this dissertation was not on syntactic algebraic optimization and efficient physical execution plans, this framework provides the opportunity to implement the method `consume_resources` efficiently. For instance, the algebraic operator `Join` could be implemented as a `SemiJoin`, i.e. it could store only the instances of the left operand in the result; this would be possible if the right operand is only referred to only in the predicate of `Join` and in none of the parent operators.

If algebraic operators (objects) are viewed as nodes and their relationships to operand objects as edges, the outer algebraic object can be viewed as the root of an algebraic tree that can potentially be rearranged and optimized using well-known techniques from relational databases. However, the use of new and monoid-aware algebraic operators in our context provides opportunities for further research in order to identify algebraic syntactic transformations and efficient execution plans for OQL.

Bibliography

- [AFP00] M. Akhtar Ali, A.A.A. Fernandes, and N.W. Paton. Incremental maintenance of materialized OQL views. In *ACM Intl Workshop on Data Warehousing and OLAP*, pages 41–48, 2000.
- [Afs98] M. Afshar. *An Open Parallel Architecture for Data-intensive Applications*. PhD thesis, University of Cambridge, Computer Laboratory, 1998.
- [AH95] O. Agesen and U. Holzle. Type feedback vs. concrete type inference: a comparison of optimization techniques for object-oriented languages. In *ACM Conf on Object-Oriented Programming, Systems, Languages and Applications (OOPSLA)*, pages 91–107, 1995.
- [AHLS00] J. Albrecht, W. Hümmer, W. Lehner, and L. Schlesinger. Query optimization by using derivability in a data warehouse environment. In *ACM Intl Workshop on Data Warehousing and OLAP*, pages 49–56, 2000.
- [AIS93] R. Agrawal, T. Imielinski, and A. Swami. Mining association rules between sets of items in large databases. In *ACM SIGMOD Intl Conf on Management of Data*, pages 207–216, 1993.
- [Ala97] S. Alagić. The ODMG object model: Does it make sense? In *ACM Conf on Object-Oriented Programming, Systems, Languages and Applications (OOPSLA)*, pages 253–270, 1997.
- [Ala99] S. Alagić. Type checking OQL queries in the ODMG type systems. *ACM Transactions on Database Systems*, 24(3):319–360, September 1999.
- [ASY98] C.C. Aggarwal, Z. Sun, and P.S. Yu. Online algorithms for finding profile association rules. In *ACM Intl Conf on Information and Knowledge Management (CIKM)*, pages 86–95, 1998.

- [BBSV97] D. Beneventano, S. Bergamaschi, C. Sartori, and M. Vincini. ODB-QOPTIMIZER: A tool for semantic query optimization in OODB. In *Intl Conf on Data Engineering (ICDE)*, 1997.
- [BDHS96] P. Buneman, S. Davidson, G. Hillebrand, and D. Suciu. A query language and optimization techniques for unstructured data. In *ACM SIGMOD Intl Conf on Management of Data*, pages 505–516, 1996.
- [Bel96] S. Bell. Deciding distinctness of query results by discovered constraints. In *Intl Conf on the Practical Application of Constraint Technology (PACT)*, pages 399–417, 1996.
- [BHP94] M.W. Bright, A.R. Hurson, and S. Pakzad. Automated resolution of semantic heterogeneity in multidatabases. *ACM Transactions on Database Systems*, 19(2), 1994.
- [BLS⁺94] P. Buneman, L. Libkin, D. Suciu, V. Tannen, and L. Wong. Comprehension syntax. *SIGMOD Record*, 23(1):87–96, 1994.
- [BO96] P. Buneman and A. Ohori. Polymorphism and type inference in database programming. *ACM Transactions on Database Systems*, 21(1):30–76, March 1996.
- [BT00] G.M. Bierman and A. Trigoni. Towards a formal type system for ODMG OQL. Technical Report 497, University of Cambridge, Computer Laboratory, October 2000.
- [CA97] K.C.C. Chan and W-H Au. Mining fuzzy association rules. In *ACM Intl Conf on Information and Knowledge Management (CIKM)*, pages 209–215, 1997.
- [Car87] L. Cardelli. Basic polymorphic type checking. *Science of Computer Programming*, 8(2):147–172, April 1987.
- [CBB⁺00] R.G.G. Cattell, D. Barry, M. Berler, J. Eastman, D. Jordan, C. Russell, O. Schadow, T. Stanienda, and F. Velez. *The Object Data Standard: ODMG 3.0*. Morgan Kaufmann, 2000.
- [CGK⁺99] Q. Cheng, J. Gryz, F. Koo, C. Leung, L. Liu, X. Qian, and B. Schiefer. Implementation of two semantic query optimization techniques in DB2 UDB. In *Intl Conf on Very Large Databases (VLDB)*, pages 687–698, 1999.
- [CGM88] U.S. Chakravarthy, J. Grant, and J. Minker. Foundations of semantic query optimization for deductive databases. In *Foundations of Deductive Databases and Logic Programming (DDL P)*, pages 243–273, 1988.

- [CGM90] U. Chakravarthy, J. Grant, and J. Minker. Logic-based approach to semantic query optimization. *ACM Transactions on Database Systems*, 15(2):162–207, 1990.
- [Cha98] S. Chaudhuri. An overview of query optimization in relational systems. In *ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS)*, pages 34–43, 1998.
- [CLC95] J-Y. Chung, Y-J. Lin, and D.T Chang. Object and relational databases. In *ACM Conf on Object-Oriented Programming, Systems, Languages and Applications (OOPSLA Addendum)*, pages 164–169, 1995.
- [CLK97] D.W. Cheung, S.D. Lee, and B. Kao. A general incremental technique for maintaining discovered association rules. In *Intl Conf on Database Systems for Advanced Applications (DASFAA)*, pages 185–193, 1997.
- [CMN95] M.J. Carey, N.M. Mattos, and A.K. Nori. Object-relational database systems (tutorial); principles, products and challenges. In *ACM SIGMOD Intl Conf on Management of Data*, pages 455–459, 1995.
- [CS96] S. Chaudhuri and K. Shim. Optimizing queries with aggregate views. In *Conf on Extending Database Technology (EDBT)*, pages 167–182, 1996.
- [Dam85] L. Damas. *Type Assignment in Programming Languages*. PhD thesis, Edinburgh University, 1985.
- [DD88] J. Duhl and C. Damon. A performance comparison of object and relational databases using the Sun Benchmark. In *ACM Conf on Object-Oriented Programming, Systems, Languages and Applications (OOPSLA)*, pages 153–163, 1988.
- [EGLM01] J. Edmonds, J. Gryz, D. Liang, and R.J. Miller. Mining for empty rectangles in large data sets. In *Intl Conf on Database Theory (ICDT)*, pages 174–188, 2001.
- [Feg94] L. Fegaras. A uniform calculus for collection types. Technical Report 94-030, Oregon Graduate Institute, 1994.
- [Feg97] L. Fegaras. An experimental optimizer for OQL. Technical Report TR-CSE-97-007, University of Texas at Arlington, 1997.
- [Feg98] L. Fegaras. Query unnesting in object-oriented databases. In *ACM SIGMOD Intl Conf on Management of Data*, pages 49–60, 1998.

- [FM95] L. Fegaras and D. Maier. Towards an effective calculus for object query languages. In *ACM SIGMOD Intl Conf on Management of Data*, pages 47–58, 1995.
- [FM00] L. Fegaras and D. Maier. Optimizing object queries using an effective calculus. *ACM Transactions on Database Systems*, 2000.
- [FMMT96] T. Fukuda, Y. Morimoto, S. Morishita, and T. Tokuyama. Mining optimized association rules for numeric attributes. In *ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS)*, pages 182–191, 1996.
- [GGMR97] J. Grant, J. Gryz, J. Minker, and L. Raschid. Semantic query optimization for object databases. In *Intl Conf on Data Engineering (ICDE)*, pages 444–453, 1997.
- [GGZ01] P. Godfrey, J. Gryz, and C. Zuzarte. Exploiting constraint-like data characterizations in query optimization. In *ACM SIGMOD Intl Conf on Management of Data*, pages 582–592, 2001.
- [GKG⁺97] T. Grust, J. Kröger, D. Gluche, A. Heuer, and M.H. Scholl. Query evaluation in CROQUE - calculus and algebra coincide. In *15th British National Conference on Databases (BNCOD)*, pages 84–100, 1997.
- [GLQ99] J. Gryz, L. Liv, and X. Qian. Semantic query optimization techniques in IBM DB2: initial results. Technical Report CS-1999-01, York University, 1999.
- [GLR97] C. Galindo-Legaria and A. Rosenthal. Outerjoin simplification and reordering for query optimization. *ACM Transactions on Database Systems*, 22(1):43–74, March 1997.
- [GS96] T. Grust and M.H. Scholl. Translating OQL into monoid comprehensions — stuck with nested loops? Technical report, University of Konstanz, Dept. of Mathematics and Computer Science, 1996.
- [GS98a] T. Grust and M.H. Scholl. Hybrid strategies for query translation and optimization. Technical report, University of Konstanz, Dept. of Mathematics and Computer Science, 1998.
- [GS98b] T. Grust and M.H. Scholl. Query deforestation. Technical report, University of Konstanz, Dept. of Mathematics and Computer Science, 1998.
- [Gun92] C.A. Gunter. *Semantics of Programming Languages*. Foundations of Computing. MIT Press, 1992.

- [Hel98] J.M. Hellerstein. Optimization techniques for queries with expensive methods. *ACM Transactions on Database Systems*, 23(2):113–157, 1998.
- [HHCF96] J. Han, Y. Huang, N. Cercone, and Y. Fu. Intelligent query answering by knowledge discovery techniques. *IEEE Transactions on Knowledge and Data Engineering*, 8(3):373–390, 1996.
- [HIT97] Z. Hu, H. Iwasaki, and M. Takechi. Formal derivation of efficient parallel programs by construction of list homomorphisms. *ACM Transactions on Programming Languages and Systems*, 19(3):444–461, 1997.
- [HK96] C.N. Hsu and C. Knoblock. Using inductive learning to generate rules for semantic query optimization. In *Advances in Knowledge Discovery and Data Mining*, pages 425–445, 1996.
- [HK98] C.N. Hsu and C. A. Knoblock. Discovering robust knowledge from databases that change. *Data Mining and Knowledge Discovery*, 2(1):69–95, 1998.
- [HK00] C.N. Hsu and C. A. Knoblock. Semantic query optimization for query plans of heterogeneous multidatabase systems. *Knowledge and Data Engineering*, 12(6):959–978, 2000.
- [Hsu96] C.N. Hsu. *Learning effective and robust knowledge for semantic query optimization*. PhD thesis, Dept of Computer Science, University of Southern California, 1996.
- [Hsu97] C.N. Hsu. Tradeoff in rule induction for semantic query optimization. In *AAAI Workshop on Building Resource-Bounded Reasoning Systems*, pages 61–68, 1997.
- [Hul97] R. Hull. Managing semantic heterogeneity in databases: A theoretical perspective. In *ACM Symposium on Principles of Database Systems (PODS)*, 1997.
- [Ioa96] Y.E. Ioannidis. Query optimization. *ACM Computing Surveys*, 28(1):121–123, 1996.
- [JS97] V. Josifovski and S.Y.W. Su. Incorporating association pattern and operation specification in ODMG’s OQL. In *ACM Intl Conf on Information and Knowledge Management (CIKM)*, pages 332–339, 1997.
- [Kin81a] J.J. King. *Query optimization by semantic reasoning*. PhD thesis, Stanford University, 1981.
- [Kin81b] J.J. King. Quist: A system for semantic query optimization in relational databases. In *Intl Conf on Very Large Databases (VLDB)*, pages 510–517, 1981.

- [KMR⁺94] M. Klemettinen, H. Mannila, P. Ronkainen, H. Toivonen, and A.I. Verkamo. Finding interesting rules from large sets of discovered association rules. In *ACM Intl Conf on Information and Knowledge Management (CIKM)*, pages 401–407, 1994.
- [KPH98] J. Kröger, S. Paul, and A. Heuer. Query optimization: On the ordering of rules. Technical report, CS Department, University of Rostock, 1998.
- [LS95] A.Y. Levy and Y. Sagiv. Semantic query optimization in datalog programs (extended abstract). In *ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS)*, pages 163–173, 1995.
- [MC96] I. Min and A. Chen. Query answering using discovered rules. In *Intl Conf on Data Engineering (ICDE)*, pages 402–411, 1996.
- [Mil78] R. Milner. A theory of type polymorphism in programming. *Journal of Computer Systems Sciences*, 17:348–375, 1978.
- [MY97] R.J. Miller and Y. Yang. Association rules over interval data. In *ACM SIGMOD Intl Conf on Management of Data*, 1997.
- [NLH98] R.T. Ng, L.V.S. Lakshmanan, and J. Han. Exploratory mining and pruning optimizations of constrained association rules. In *ACM SIGMOD Intl Conf on Management of Data*, pages 13–24, 1998.
- [OHMS92] J. Orenstein, S. Haradhvala, B. Margulies, and D. Sakahara. Query processing in ObjectStore database system. In *ACM SIGMOD Intl Conf on Management of Data*, pages 403–412, 1992.
- [Par95] J.S. Park. An effective hash-based algorithm for mining association rules. In *ACM SIGMOD Intl Conf on Management of Data*, pages 175–186, 1995.
- [PC94] J. Plevyak and A.A. Chien. Precise concrete type inference for object-oriented languages. In *ACM Conf on Object-Oriented Programming, Systems, Languages and Applications (OOPSLA)*, pages 324–340, 1994.
- [PHH92] H. Pirahesh, J.M. Hellerstein, and W. Hasan. Extensible/rule based query rewrite optimization in Starburst. In *ACM SIGMOD Intl Conf on Management of Data*, pages 39–48, 1992.
- [PLH97] H. Pirahesh, T.Y.C. Leung, and W. Hasan. A rule engine for query transformation in Starburst and IBM DB2 C/S DBMS. In *Intl Conf on Data Engineering (ICDE)*, pages 391–400, 1997.

- [PS91] J. Palsberg and M. I. Schwartzbach. Object-oriented type inference. In *ACM Conf on Object-Oriented Programming, Systems, Languages and Applications (OOPSLA)*, pages 146–161, 1991.
- [PS96] A. Poulouvasilis and C. Small. Algebraic query optimisation for database programming languages. *VLDB Journal*, 5(2):119–132, 1996.
- [PS97] A. Poulouvasilis and C. Small. Formal foundations for optimising aggregation functions in database programming languages. In *Workshop on Database Programming Languages (DBPL)*, pages 299–318, 1997.
- [Pul96] E. Pulido. Recursive query processing using graph traversal techniques. In *ACM Intl Conf on Information and Knowledge Management (CIKM)*, pages 37–44, 1996.
- [Rob65] J. A. Robinson. A machine-oriented logic based on the resolution principle. *Journal of the ACM*, 12(1):23–41, January 1965.
- [RS93] C. Rich and M.H. Scholl. Query optimization in oodbms. In *BTW*, pages 266–284, 1993.
- [RS97] H. Riedel and M.H. Scholl. A formalization of ODMG queries. In *Seventh IFIP 2.6 Working Conference on Database Semantics (DS-7)*, 1997.
- [SA96] R. Srikant and R. Agrawal. Mining quantitative association rules in large relational tables. In *ACM SIGMOD Intl Conf on Management of Data*, pages 1–12, 1996.
- [Sam00] P.R.F. Sampaio. *Design and Implementation of a Deductive Query Language for ODMG Compliant Object Databases*. PhD thesis, University of Manchester, 2000.
- [SdBB96] H.J. Steenhagen, R.A. de By, and H.M. Blanken. Translating OSQL-queries into efficient set expressions. In *Conf on Extending Database Technology (EDBT)*, pages 183–200, 1996.
- [SHKC93] S. Shekar, B. Hamidzadeh, A. Kohli, and M. Coyle. Learning transformation rules for semantic query optimization. *IEEE Transactions on Knowledge and Data Engineering*, 5(6):950–964, 1993.
- [SLR94] P. Seshadri, M. Livny, and R. Ramakrishnan. Sequence query processing. In *ACM SIGMOD Intl Conf on Management of Data*, pages 430–441, 1994.
- [SO95a] D.D. Straube and M.T. Özsu. Queries and query processing in object-oriented database systems. *ACM Transactions on Information Systems*, 8(4):387–430, 1995.

- [SO95b] D.D. Straube and M.T. Özsu. Query optimization and execution plan generation in object-oriented data management systems. *IEEE Transactions on Knowledge and Data Engineering*, 7(2):210–227, 1995.
- [SP00] P.R.F. Sampaio and N.W. Paton. Query processing in DOQL: A deductive database language for the ODMG model. *Data and Knowledge Engineering*, 2000.
- [SSS92] M. Siegel, E. Sciore, and S. Salveter. A method for automatic rule derivation to support semantic query optimization. *ACM Transactions on Database Systems*, 17(4):563–600, December 1992.
- [ST95] A. Shrufi and T. Topaloglou. Query processing for knowledge bases using join indices. In *ACM Intl Conf on Information and Knowledge Management (CIKM)*, pages 158–166, 1995.
- [Sul00] M. Sulzmann. A general type inference framework for Hindley/Milner style systems. Technical Report 2000/15, Dept. of Computer Science and Software Engineering, The University of Melbourne, July 2000.
- [TB01] A. Trigoni and G.M. Bierman. Inferring the principal type and the schema requirements of an OQL query. In *18th British National Conference on Databases (BNCOD)*, pages 185–201, 2001.
- [TM01] A. Trigoni and K. Moody. Using association rules to add or eliminate query constraints automatically. In *Intl Conf on Scientific and Statistical Database Management (SS-DBM)*, pages 124–133, 2001.
- [Tof87] M. Tofte. *Operational semantics and polymorphic type inference*. PhD thesis, University of Edinburgh, 1987.
- [TUA⁺98] D. Tsur, J. Ullman, S. Abiteboul, C. Clifton, R. Motwani, S. Nestorov, and A. Rosenthal. Query flocks: a generalization of association-rule mining. In *ACM SIGMOD Intl Conf on Management of Data*, pages 1–12, 1998.
- [Van92] B. Vance. Towards an object-oriented query algebra. Technical report, Oregon Graduate Institute, 1992.
- [WM99] L.B. Warshaw and D.P. Miranker. Rule-based query optimization, revisited. In *ACM Intl Conf on Information and Knowledge Management (CIKM)*, pages 267–275, 1999.

- [WYY00] W. Wang, J. Yang, and P.S. Yu. Efficient mining of weighted association rules (war). In *ACM SIGKDD Intl Conf on Knowledge Discovery and Data Mining (KDD)*, pages 270–274, 2000.
- [Yoo94] S.C. Yoon. Intelligent query answering in deductive and object-oriented databases. In *ACM Intl Conf on Information and Knowledge Management (CIKM)*, pages 244–251, 1994.
- [YS89] C.T. Yu and W. Sun. Automatic knowledge acquisition and maintenance for semantic query optimization. *IEEE Transactions on Knowledge and Data Engineering*, 1(3), 1989.
- [YSP95] S.C. Yoon, I.Y. Song, and E.K. Park. Semantic query processing in object-oriented databases using deductive approach. In *ACM Intl Conf on Information and Knowledge Management (CIKM)*, pages 150–157, 1995.
- [YSP97] S.C. Yoon, I.Y. Song, and E.K. Park. Intensional query processing using data mining approaches. In *ACM Intl Conf on Information and Knowledge Management (CIKM)*, pages 201–208, 1997.
- [ZM90] S. Zdonik and D. Maier. *Readings in Object-Oriented Database Systems*. Morgan Kaufmann, 1990.
- [ZRL96] T. Zhang, R. Ramakrishnan, and M. Livny. Birch: an efficient data clustering method for very large databases. In *ACM SIGMOD Intl Conf on Management of Data*, pages 103–114, 1996.