

# Fine Grained Information Integration with Description Logics

## Information Integration, Description Logics, Query Optimization

### Abstract

We outline an approach to query optimization in which a *description logic* (DL) reasoner serves a crucial strategic role, and present an example application of our approach in fine grained information integration. In particular, the application demonstrates how the internal structure of an unfolded B-tree can be captured as a terminology, and how an access plan that navigates this internal structure can be found with the aid of a DL reasoner.

## 1 Introduction

An *embedded control program* (ECP) is a new application area for database technology [Toman and Weddell, 2001b]. An ECP that is a *legacy* system presents an additional challenge. In particular, the ECP will already have code (sometimes a great deal of code) that specifies the internal encoding for a database of *control data*. To enable high-level access to this data using SQL-like languages, it becomes necessary to optimize queries over conceptual views. A query optimizer must therefore be capable of *fine grained information integration*; that is, it must be possible to supply descriptions of the internal encoding of control data as part of the input to the optimizer.

We have developed a novel resource bounded query optimizer in which integrity constraints that abstract the internal encoding of control data are used to extend the search space of possible query plans for a given source query. This paper links the optimizer with a powerful DL reasoner that enables complex query rewrites. In particular, we show how the data structures that constitute an internal control data encoding can be captured by a terminology in  $\mathcal{DLFDE}$ , a refinement of an earlier DL dialect called CFD [Khizder *et al.*, 2000]. We then show how this optimizer uses a DL reasoner for  $\mathcal{DLFDE}$  to help find access plans that navigate the data structures for queries over conceptual views. The process relies on dynamic construction of descriptions that characterize properties of subqueries generated during query optimization. A complicating factor, due to the nature of the application itself, is the presumption of a *bag semantics* for an underlying query algebra.

An unfolded B-tree data structure will be used as a running example. A schema for the B-tree is given in Figure 1;

squares and circles represent primitive classes in this schema, while unlabeled wide arcs denote sub-classing and labeled narrow arcs denote attributes. The  $P_i$  classes on the right correspond to data pages of the B-tree at level  $i$ , while the  $E_j$  classes in the middle correspond to employee records within the data pages. Classes  $E_jP$  denote *all* employee records in a data page at level  $i$  for  $i \geq j$ . The records are presumed to have a *name* field. A conceptual view of the data is represented by the EMP and STR classes. Also illustrated is a primitive concept on the left that abstracts a sample query over EMP and STR.

Squares have the added significance of denoting low-level indexes for which dotted outgoing arcs represent search parameters. For example, since no parameters for the root page  $P_0$  are required, a global variable provides access to  $P_0$ ; and because of their respective parameters  $a_0$  and  $p$ , access to employee record addresses in  $E_0$  and the level one pages  $P_1$  can be obtained by scanning an array within  $P_0$ .

Based on a terminological abstraction  $\mathcal{T}$  of the B-tree that is expressed in terms of  $\mathcal{DLFDE}$ , we show how the optimizer uses a DL reasoner to translate a request for all distinct employee names,

```
select distinct e.name as n from EMP as e,
```

to the following equivalent formulation:

```
select n
from [P0 as v0], (
  (select e, v0 from [E0 as e, e.a0 = v0] )
  union all
  (select e, v0 from [P1 as v1, v1.p = v0], (
    (select e, v1 from [E1 as e, e.a1 = v1] )
    union all
    (select e, v1 from [P2 as v2, v2.p = v1],
      [E2 as e, e.a2 = v2] )))
  )
[n = e.name].
```

To interpret the latter as an access plan, the subexpressions enclosed in all but the last box should be understood as index scans (note that some of the indices are supplied with parameters), and the last box as an assignment. The *union all*, *select*, and *from* operations should be understood as concatenations, as duplicate-preserving projections, and as (iterated) nested loop joins, respectively. With these assumptions, the second formulation obtains all employee names by per-

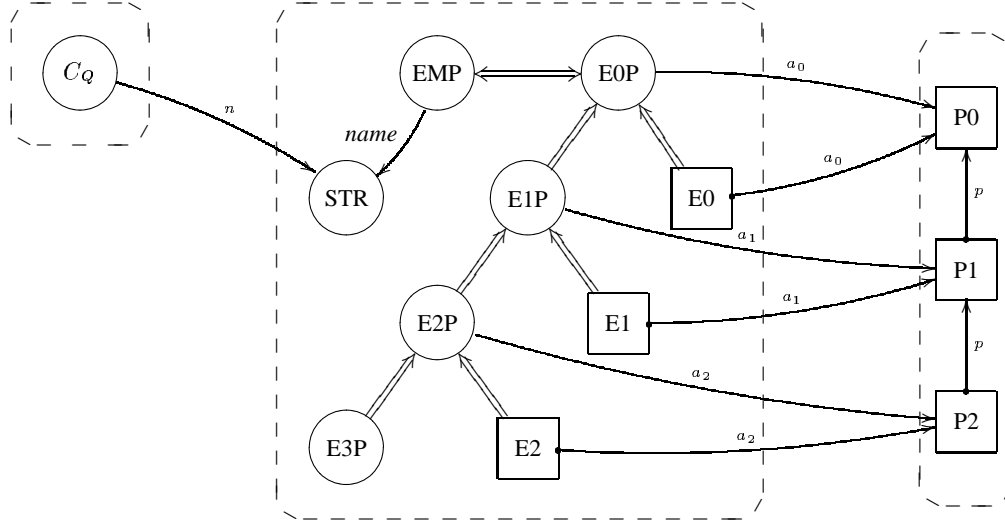


Figure 1: SCHEMA FOR THREE-LEVEL B-TREE ON CLASS EMP.

forming a preorder traversal of the index pages that comprise the B-tree.

Part of what is needed for this to work is an ability to include certain kinds of equational constraints in a terminology. For example, there is a crucial logical relationship between employee records not occurring in the root page, E1P, and the data pages on the first level of the B-tree, P1. This dependency can be captured by asserting that *the  $a_0$  value of each employee record in class E1P is equal to its  $a_{1.p}$  value*, i.e., by adding the inclusion dependency  $E1P \sqsubseteq (a_0 = a_{1.p})$  to  $\mathcal{T}$ . Thus, another contribution is the incorporation of an equational concept constructor in  $\mathcal{DLFDE}$  in a way that retains decidability of the logical implication problem for  $\mathcal{DLFDE}$ , and that remains sufficiently expressive for fine grained data descriptions such as the unfolded B-tree case.

### 1.1 Related Work

A good survey of how DL reasoning can be useful in information systems, circa 1995, can be found in [Borgida, 1995]. Of particular relevance to our own work is the use of DLs as a lingua franca for various database schema languages [Call *et al.*, 2001; Calvanese *et al.*, 1998b], and then as a means for information integration [Calvanese *et al.*, 2001a; 2001c; 1998a]. Recently, there has been some work on using DLs to reason about set query containment [Calvanese *et al.*, 2000; Horrocks *et al.*, 2000] and about bag query equivalence [Khizder *et al.*, 2000] in the presence of database schema expressed in terms of DLs. Much of this more recent work also depends on an ability to express functional constraints [Calvanese *et al.*, 2001b; Khizder *et al.*, 2001].

The remainder of the paper is organized as follows. A DL dialect called  $\mathcal{DLFDE}$  and a query language called  $\mathcal{QL}$  is first introduced in Section 2. Section 3 then defines a number of bidirectional rules for rewriting queries. Our focus is on rules that involve reasoning with  $\mathcal{DLFDE}$  terminologies that in turn abstract database schema and the structural properties of the queries themselves. Section 3 also outlines a general

approach to query optimization based on these rules, and efficient ways to integrate reasoning in  $\mathcal{DLFDE}$  with query rewriting. Section 4 gives a brief summary and outline of further extensions of the approach.

## 2 Definitions

We begin with a description (or feature) logic that allows the use of a concept constructor for equational constraints in possibly cyclic terminologies. This is only permitted in a way that ensures decidability for the associated implication problem. The new DL generalizes an earlier description logic defined in [Khizder *et al.*, 2000].

**Definition 1 (Syntax and Semantics of  $\mathcal{DLFDE}$ )** *Let  $F$  be a set of attribute names. A path expression is defined by the grammar “ $\text{Pf} ::= f^{i,j} . \text{Pf} \mid \text{Id}$ ” for  $f^{i,j} \in F$ . Also, the right and left superscripts of consecutive attributes must match in every well formed path description.*

*Let  $C^i, i \in N$  be primitive concept description(s). We define derived concept descriptions by the grammar in Figure 2.*

*An inclusion dependency is of the form  $D^i \sqsubseteq E^i$ .*

*The semantics of expressions is defined with respect to a many-sorted structure  $(\Delta^i, \mathcal{I})$ , where  $\Delta^i$  are disjoint domains of “objects” and  $(\cdot)^{\mathcal{I}}$  an interpretation function that fixes the interpretations of primitive concepts  $C^i$  to be subsets of  $\Delta^i$  and primitive attributes  $f^{i,j}$  to be total functions  $(f^{i,j})^{\mathcal{I}} : \Delta^i \rightarrow \Delta^j$ . The interpretation is extended to path expressions,  $(\text{Id})^{\mathcal{I}} = \lambda x.x$ ,  $(f^{i,j} . \text{Pf})^{\mathcal{I}} = (\text{Pf})^{\mathcal{I}} \circ (f^{i,j})^{\mathcal{I}}$  and derived concept descriptions  $D^i$  and  $E^i$  as defined in Figure 2.*

*An interpretation satisfies an inclusion dependency  $D^i \sqsubseteq E^i$  if  $(D^i)^{\mathcal{I}} \subseteq (E^i)^{\mathcal{I}}$ .*

*A terminology  $\mathcal{T}$  consists of a finite set of inclusion dependencies, and is stratified if:*

1. *For each description  $\forall f^{i,j} . D$ , we have  $i \leq j$ ;*
2. *For each description  $D\{ \text{Pf}_1, \dots, \text{Pf}_k \} \rightarrow \text{Pf}$ , each  $\text{Pf}_l$*

SYNTAX	SEMANTICS: DEFN OF “ $(\cdot)^{\mathcal{I}}$ ”
$D^i ::= C^i$ $\quad   D_1^i \sqcap D_2^i$ $\quad   D_1^i \sqcup D_2^i$ $\quad   \neg D^i$ $\quad   \forall f^{i,j}. D^j$ $E^i ::= D^i$ $\quad   E_1^i \sqcap E_2^i$ $\quad   \forall f^{i,j}. E^j$ $\quad   (\text{Pf}_1 = \text{Pf}_2)$ $\quad   D^i \{ \text{Pf}_1, \dots, \text{Pf}_k \} \rightarrow \text{Pf}$	$(C^i)^{\mathcal{I}} \subseteq \Delta^i$ $(D_1^i)^{\mathcal{I}} \sqcap (D_2^i)^{\mathcal{I}}$ $(D_1^i)^{\mathcal{I}} \sqcup (D_2^i)^{\mathcal{I}}$ $\Delta^i \setminus (D^i)^{\mathcal{I}}$ $\{x : (f^{i,j})^{\mathcal{I}}(x) \in (D^j)^{\mathcal{I}}\}$ $(E_1^i)^{\mathcal{I}} \sqcap (E_2^i)^{\mathcal{I}}$ $\{x : (f^{i,j})^{\mathcal{I}}(x) \in (E^j)^{\mathcal{I}}\}$ $\{x : (\text{Pf}_1)^{\mathcal{I}}(x) = (\text{Pf}_2)^{\mathcal{I}}(x)\}$ $\{x : \forall y \in (D^i)^{\mathcal{I}}.$ $\quad \bigwedge_{i=1}^k (\text{Pf}_i)^{\mathcal{I}}(x) = (\text{Pf}_i)^{\mathcal{I}}(y) \Rightarrow (\text{Pf})^{\mathcal{I}}(x) = (\text{Pf})^{\mathcal{I}}(y)\}$

Figure 2: SYNTAX AND SEMANTICS OF  $\mathcal{DLFDE}$ .

consists only of primitive attributes  $f^{i,j}$  for which  $i = j$ ; and

3. For each description  $(\text{Pf}_1 = \text{Pf}_2)$ , each  $\text{Pf}_l$  consists only of primitive attributes  $f^{i,j}$  for which  $i \leq j$  and has at least one primitive attribute for which  $i < j$ .

The logical implication problem asks if  $\mathcal{T} \models D^i \sqsubseteq E^i$  holds; that is, if all interpretations that satisfy all constraints in  $\mathcal{T}$  also satisfy  $(D^i)^{\mathcal{I}} \subseteq (E^i)^{\mathcal{I}}$  (the posed question).

Herein, we simplify the notation for path expressions by omitting the superscripts for descriptions and primitive attributes whenever clear from the context, and allow a syntactic composition  $\text{Pf}_1 . \text{Pf}_2$  of path expressions that stands for their concatenation. The following theorem is an extension of an earlier result [Toman and Weddell, 2001a].

**Proposition 2** *Let  $\mathcal{T}$  be a stratified  $\mathcal{DLFDE}$  terminology. Then the  $\mathcal{DLFDE}$  implication problem  $\mathcal{T} \models D^i \sqsubseteq E^i$  is decidable and complete for DEXPTIME.*

Relaxing the restrictions on terminologies leads to undecidable implication problems. In particular, an absence of either condition (1) or (3) will allow equations of the form  $(f_1.f_2 = f_3)$  on a cyclic schema (which is well known to be undecidable [Machtey and Young, 1978]). Perhaps surprisingly, allowing path functional dependencies (and in turn keys) to constrain the behavior of  $f^{i,j}$ -attributes for which  $i < j$  also leads to undecidability, as illustrated by the following.

**Example 3** Consider a general (undecidable) decision problem of the form

$$\{C^0 \sqsubseteq \forall f_{i,j}^{0,0}. C^0\} \cup \boxed{\{C^0 \sqsubseteq (f_{i,1}^{0,0}.f_{i,2}^{0,0} = f_{i,3}^{0,0})\}}$$

$$\models C^0 \sqsubseteq (f_1^{0,0}.f_2^{0,0} = f_3^{0,0})$$

where  $0 < i \leq n$  and  $0 < j \leq 3$ . The inclusion dependencies enclosed in the box have the following equivalent formulation (note the use of  $f^{0,1}$ -attributes).

$$\{C^0 \sqsubseteq (f_{i,1}^{0,0}.f_{i,1}^{0,1} = f_{i,2}^{0,1}) \sqcap (f_{i,2}^{0,0}.f_{i,3}^{0,1} = f_{i,1}^{0,1})$$

$$\quad \sqcap (f_{i,3}^{0,0}.f_{i,3}^{0,1} = f_{i,2}^{0,1}) \sqcap C^0 \{f_{i,3}^{0,1}\} \rightarrow Id \}$$

The reformulated terminology now satisfies conditions (1) and (3), but fails to satisfy condition (2).

- 1:  $\text{EMP} \sqsubseteq \forall \text{name}. \text{STR}$
- 2:  $\text{EMP} \sqsubseteq \text{EMP}\{\text{name}\} \rightarrow Id$
- 3:  $\text{P0} \sqsubseteq \overline{\text{P0}}\{\} \rightarrow Id$
- 4:  $\text{P1} \sqsubseteq \forall p. \text{P0}$
- 5:  $\text{P2} \sqsubseteq \forall p. \text{P1}$
- 6:  $\text{EMP} \sqsubseteq \text{E0P}$
- 7:  $\text{E0P} \sqsubseteq \text{EMP} \sqcap (\forall a_0. \text{P0}) \sqcap (\text{E0} \sqcup \text{E1P})$
- 8:  $\text{E0} \sqsubseteq \text{E0P}$
- 9:  $\text{E1P} \sqsubseteq \text{E0P} \sqcap (\forall a_1. \text{P1}) \sqcap (\text{E1} \sqcup \text{E2P}) \sqcap (a_1.p = a_0)$
- 10:  $\text{E0} \sqcap \text{E1P} \sqsubseteq \perp$
- 11:  $\text{E1} \sqsubseteq \text{E1P}$
- 12:  $\text{E2P} \sqsubseteq \text{E1P} \sqcap (\forall a_2. \text{P2}) \sqcap (\text{E2} \sqcup \text{E3P}) \sqcap (a_2.p = a_1)$
- 13:  $\text{E1} \sqcap \text{E2P} \sqsubseteq \perp$
- 14:  $\text{E2} \sqsubseteq \text{E2P}$
- 15:  $\text{E3P} \sqsubseteq \perp$

Figure 3: AN UNFOLDED B-TREE SCHEMA IN  $\mathcal{DLFDE}$ .

For the unfolded B-tree, however, the restriction to stratified terminologies does not cause any problems. In particular, a terminology that corresponds to the graphical representation of this data structure given earlier in Figure 1 is listed in Figure 3. An interpretation  $(\cdot)^{\mathcal{I}}$  induces a (class of) *database instances* over which we can formulate queries. The query language we use for this purpose is a positive fragment of an object-relational first-order language with duplicate semantics:

**Definition 4 (Object Relational Queries)** *Figure 4 defines the syntax (left) and semantics (right) of the query language  $\mathcal{QL}$ . In the syntax, the symbols  $a$ ,  $a_i$ , and  $b$  stand for query variables (identifiers). The semantics is given with respect to an interpretation  $(\cdot)^{\mathcal{I}}$ . In addition we assume standard syntactic safety conditions to be satisfied by the queries.*

For the remainder of the paper, we use parentheses and common abbreviations to make queries more readable; e.g., from  $Q_1, \dots, Q_k$  stands for iterated binary joins.

### 3 Query Optimization

We first present a collection of bidirectional rules for rewriting queries expressed in  $\mathcal{QL}$ , with a particular focus on the rules that depend critically on an ability to reason about

$Q ::= D \text{ as } a$ $\left  \begin{array}{l} a. \text{Pf}_1 = b. \text{Pf}_2 \\ \text{select } a_1, \dots, a_n \ Q \\ \text{elim } a_1, \dots, a_n \ Q \\ \text{true} \\ \text{from } Q_1, Q_2 \\ \text{empty } a_1, \dots, a_n \\ Q_1 \text{ union all } Q_2 \end{array} \right.$	$\llbracket Q \rrbracket \mathcal{I} := \{ \langle a : v \rangle : v \in (D)^{\mathcal{I}} \}$ $\{ \langle a : v, b : w \rangle : (\text{Pf}_1)^{\mathcal{I}}(v) = (\text{Pf}_2)^{\mathcal{I}}(w) \}$ $\{ \langle a_1 : v @ a_1, \dots, a_n : v @ a_n \rangle : v \in \llbracket Q \rrbracket \mathcal{I} \}$ $\{ \langle a_1 : v @ a_1, \dots, a_n : v @ a_n \rangle : v \in \llbracket Q \rrbracket \mathcal{I} \}$ $\{ \langle \rangle \}$ $\llbracket Q_1 \rrbracket \mathcal{I} \bowtie \llbracket Q_2 \rrbracket \mathcal{I}$ $\emptyset$ $\llbracket Q_1 \rrbracket \mathcal{I} \uplus \llbracket Q_2 \rrbracket \mathcal{I}$
--	--

Figure 4: SYNTAX AND SEMANTICS OF  $\mathcal{QL}$ .

$\mathcal{DLFDE}$  terminologies. We then outline an optimization process that uses these rules, primarily by illustrating their behavior on the B-tree case.

### 3.1 Transformation Rules

Typically, rule-based query optimizers use rules that are designed to apply universally to all subqueries, possibly with respect to integrity constraints that hold in a database schema. Often, however, a subquery nested within another fixed query, what we shall call a *query context*, can be rewritten to another subquery that preserves equivalence *only with respect to the enclosing query context*. Our formulation of query rules apply in this more general circumstance, and therefore depend on the following definition.

**Definition 5 (Context)** *An expression  $Q[\ ]$  in the language  $\mathcal{QL}$  enriched by an additional terminal symbol  $\llbracket \ \rrbracket$  is called a query context. For a query  $Q \in \mathcal{QL}$ , the expression  $Q[Q']$ , denotes the syntactical substitution of  $Q'$  for  $\llbracket \ \rrbracket$ . We also say that  $Q'$  is compatible with  $Q[\ ]$  if  $Q[Q'] \in \mathcal{QL}$ .*

To benefit from terminological reasoning, the rules refer to descriptions that abstract both subqueries  $Q$  and enclosing query contexts  $Q[\ ]$ , denoted  $E_Q$  and  $E_{Q[\ ]}$ , respectively. These descriptions capture important structural information about subquery results and considerably enhance the power of the rules. In particular, by adding inclusion dependencies of the form  $C_Q \sqsubseteq E_Q$  and  $C_Q \sqsubseteq E_{Q[\ ]}$  to a given terminology ( $C_Q$  denotes a primitive concept anchor for the descriptions), a  $\mathcal{DLFDE}$  reasoner can deduce additional information needed by the rules. Definitions of  $E_Q$  and  $E_{Q[\ ]}$  are as follows.

$$E_Q = \begin{cases} \forall a.D, & \text{if } Q = D \text{ as } a; \\ (a_1. \text{Pf}_1 = a_2. \text{Pf}_2), & \text{if } Q = (a_1. \text{Pf}_1 = a_2. \text{Pf}_2); \\ E_{Q'}, & \text{if } Q = \text{select } V \ Q' \text{ or } \text{elim } V \ Q'; \\ \top, & \text{if } Q = \text{true}; \\ E_{Q_1} \sqcap E_{Q_2}, & \text{if } Q = \text{from } Q_1, Q_2; \\ \perp, & \text{if } Q = \text{empty } V; \text{ or} \\ E_{Q_1} \sqcup E_{Q_2}, & \text{if } Q = Q_1 \text{ union all } Q_2. \end{cases}$$

$$E_{Q[\ ]} = \begin{cases} \top, & \text{if } Q[\ ] = \llbracket \ \rrbracket; \\ E_{Q[\ ]} \sqcap E_{Q'}, & \text{if } Q[\ ] = Q[\text{from } Q', \llbracket \ \rrbracket] \text{ or } Q[\text{from } \llbracket \ \rrbracket, Q']; \\ E_{Q[\ ]}, & \text{if } Q[\ ] = Q[\text{elim } V[\ ]], Q[\text{select } V[\ ]], \\ & Q[Q' \text{ union all } \llbracket \ \rrbracket], \text{ or} \\ & Q[\llbracket \ \rrbracket \text{ union all } Q']. \end{cases}$$

For every query  $Q$  and query context  $Q[\ ]$ , we also define the set  $\alpha_Q$  of *free variables* of  $Q$  and the set  $\alpha_{Q[\ ]}$  of *variables captured* by the context  $Q[\ ]$ . These two sets are defined inductively on the structure of  $Q$  and  $Q[\ ]$ , respectively.

The rules are listed in Figure 5, and operate as follows.

$\bowtie$ -intro: Given an existing query variable  $a$ , the DL reasoner infers that the hypothetical query object  $C_Q$  satisfies the constraint  $C_Q \sqsubseteq \forall a.D$ . This means that all valuations of  $a$  in  $Q[\ ]$  must also belong to  $(D)^{\mathcal{I}}$  and thus we can add or remove a conjunct “ $D$  as  $a$ ” to or from the query. Also note that the terminal `true` can be introduced or removed anywhere in  $Q$  using a standard rule for join.

=-intro: As above, an equational constraint on query variables that is implied in the context  $Q[\ ]$  by the terminology and the abstraction  $C_Q \sqsubseteq E_{Q[\ ]}$  can be freely added or removed.

var-intro: Existing equational constraints in queries can be split by introducing a new query variable with a unique name. Although this rule does not depend on DL reasoning, it enables the introduction of further information from the database schema into the query by subsequent use of rules  $\bowtie$ -intro and =-intro.

$\uplus$ -intro: Query variables labeled with disjunctive descriptions can be translated to the `union all` operation. Note the use of the duplicate elimination operator `elim` to account for the discrepancy between the set semantics of the description and the bag semantics of the `union all` operation.

$\varepsilon$ -elim: The (usually expensive) `elim` operation can be removed from queries when its input is duplicate free. This is guaranteed by the semantics of the atomic queries “ $D$  as  $a$ ” and “ $a. \text{Pf}_1 = b. \text{Pf}_2$ ”.

$\bowtie$ - $\varepsilon$ -dist: An `elim` operator can be distributed over a join to subqueries. However, this is only possible if the DL reasoner can deduce that the query variables of one of the subqueries are functionally determined by the “remaining” variables. This rule generalizes an earlier version [Khizder *et al.*, 2000].

$\uplus$ - $\varepsilon$ -dist: An `elim` operator can also be distributed over a `union all` if the DL reasoner can deduce that the two subqueries must be disjoint by virtue of their abstractions.

The above observations, together with an induction on query contexts, yield the following result.

**Proposition 6** *The rules in Figure 5 are sound.*

In the running example (and to achieve completeness for the conjunctive case) we also use several *administrative rules*, e.g., commutativity and associativity of joins, distributivity of unions over joins, absorption of nested duplicate elimination operators, etc.

$$\begin{array}{l}
\bowtie\text{-intro: } \frac{Q[\text{true}]}{Q[D \text{ as } a]} \quad \mathcal{T} \cup \{C_Q \sqsubseteq E_{Q[1]}\} \models C_Q \sqsubseteq \forall a.D \text{ and } a \in \alpha_{Q[1]} \\
=\text{-intro: } \frac{Q[\text{true}]}{Q[a.Pf_1 = b.Pf_2]} \quad \mathcal{T} \cup \{C_Q \sqsubseteq E_{Q[1]}\} \models C_Q \sqsubseteq (a.Pf_1 = b.Pf_2) \text{ and } a, b \in \alpha_{Q[1]} \\
\text{var-intro: } \frac{Q[a_1.Pf_1 = a_2.Pf_2.Pf_3]}{Q[\text{from } a_1.Pf_1 = b.Pf_3, b = a_2.Pf_2]} \quad b \notin \alpha_{Q[1]} \cup \{a_1, a_2\} \\
\uplus\text{-intro: } \frac{Q[(E_1 \sqcup E_2) \text{ as } a]}{Q[\text{elim } a ((E_1 \text{ as } a) \text{ union all } (E_2 \text{ as } a))]} \\
\varepsilon\text{-elim: } \frac{Q[\text{elim } a D \text{ as } a]}{Q[D \text{ as } a]} \qquad \varepsilon\text{-elim: } \frac{Q[\text{elim } a, b a.Pf_1 = b.Pf_2]}{Q[a.Pf_1 = b.Pf_2]} \\
\bowtie\text{-}\varepsilon\text{-dist: } \frac{Q[\text{elim } V \text{ from } Q_1, Q_2]}{Q[\text{select } V \text{ from } (\text{elim } \alpha_{Q_1} Q_1), (\text{elim } W Q_2)]} \quad \begin{array}{l} \mathcal{T} \cup \{C_Q \sqsubseteq E_{Q[1]} \sqcap E_{Q_1} \sqcap E_{Q_2}\} \\ \models C_Q \sqsubseteq C_Q \sqcup \{V \cup \alpha_{Q[1]}\} \rightarrow \alpha_{Q_1} \\ \text{where } W = (V \cup \alpha_{Q_1}) \sqcap \alpha_{Q_2} \end{array} \\
\uplus\text{-}\varepsilon\text{-dist: } \frac{Q[\text{elim } V (Q_1 \text{ union all } Q_2)]}{Q[(\text{elim } V Q_1) \text{ union all } (\text{elim } V Q_2)]} \quad \mathcal{T} \cup \{C_Q \sqsubseteq E_{Q[1]} \sqcap E_{Q_1} \sqcap E_{Q_2}\} \models C_Q \sqsubseteq \perp
\end{array}$$

Figure 5: QUERY TRANSFORMATION RULES.

### 3.2 Implementation

Figure 6 illustrates a sequence of applications of rules in Figure 5 that progressively transforms the “names of employees” query to the query plan given in the introduction. The order of application of these rewrites generally follows the approach outlined in [Toman and Weddell, 2001b] for a query language with set semantics. In particular, there are three phases that altogether find query plans for conjunctive fragments.

1. An initial *query expansion* phase is applied for conjunctive subexpressions (e.g., Steps 1-7 in Figure 6). The phase introduces additional existentially quantified variables and conditions whose existence is implied by the terminology. Observe the use of rules  $\bowtie$ -intro, =-intro and var-intro.
2. The second phase essentially relates to so-called *join order selection*. A “choice-point” during join order selection is to select from among the variables and conditions that relate to index use. For example, index P0 is selected (Step 8) and, if possible, removed from the scope of an `elim` operator (Step 9). An important detraction from normal practice, however, is an option to explore disjunctive descriptions using rule  $\uplus$ -intro in lieu of “selecting next index” (Step 12). In this latter case, administrative rules (Step 13) and (if applicable) rule  $\uplus$ - $\varepsilon$ -dist (Step 14) are then used to “prepare” each disjunct for a recursive application of the three phases (Steps 15-16 for the first disjunct; the remaining steps for the second disjunct).
3. The third and final *query contraction* phase then removes redundant variables and conditions not used in the actual query plan. Observe the use of rules  $\bowtie$ -intro, =-intro and var-intro in reverse (Steps 16 and 19).

Another feature of our implementation of this process is the use of a graph encoding of queries that is also used directly by the  $\mathcal{DLFDE}$  model building procedure. Also, query contraction is instead accomplished by applying query expansion on a separate *plan graph* constructed during join order selection.

### 4 Summary

We have outlined an approach to query optimization for a positive fragment of a query language with duplicate semantics for which a DL reasoner is used to search for possible index use, including cases that require exploring horizontally partitioned data, and to reason about duplicate elimination and its interaction with the join and union operations. We have also presented  $\mathcal{DLFDE}$ , a boolean-complete DL dialect with a DEXPTIME-complete implication problem that includes concept constructors for both functional and equational constraints, and allows a stratified use of these constructors to occur in possibly cyclic terminologies.

There are many opportunities for future work. For example, we are currently exploring alternative restrictions on  $\mathcal{DLFDE}$  terminologies that continue to guarantee decidability. Completeness results that relate to query equivalence defined in terms of our rules are another ongoing avenue of research. Another direction is the incorporation of *order dependencies* [Toman and Weddell, 2001a] for reasoning about order and aggregate optimizations and their interaction with the current rules. Finally, we are continuing our efforts on evaluating an experimental implementation of our query optimizer on a number of real-world test cases such as the Linux kernel data structures.

### References

- [Borgida, 1995] Alexander Borgida. Description logics in data management. *IEEE Transactions on Knowledge and Data Engineering*, 7(5):671–682, 1995.
- [Call et al., 2001] Andrea Call, Diego Calvanese, Giuseppe De Giacomo, and Maurizio Lenzerini. Reasoning on UML Class Diagrams in Description Logics. In *Proc. of IJCAR Workshop on Precise Modelling and Deduction for Object-oriented Software Development (PMD 2001)*, 2001.
- [Calvanese et al., 1998a] Diego Calvanese, Giuseppe De Giacomo, Maurizio Lenzerini, Daniele Nardi, and Riccardo Rosati. Description Logic Framework for Information Integration. In *Prin-*

1. $\text{elim } n \text{ from EMP as } e, n = e.name$	
2. $\text{elim } n \text{ from EMP as } e, n = e.name, \text{EOP as } e$	$\bowtie$ -intro and 6
3. $\text{elim } n \text{ from EMP as } e, n = e.name, \text{EOP as } e, (\text{E0} \sqcup \text{E1P}) \text{ as } e$	$\bowtie$ -intro and 7
4. $\text{elim } n \text{ from EMP as } e, n = e.name, \text{EOP as } e, (\text{E0} \sqcup \text{E1P}) \text{ as } e, e.a_0 = e.a_0$	=-intro
5. $\text{elim } n \text{ from EMP as } e, n = e.name, \text{EOP as } e, (\text{E0} \sqcup \text{E1P}) \text{ as } e, e.a_0 = v_0, v_0 = e.a_0$	var-intro
6. $\text{elim } n \text{ from EMP as } e, n = e.name, \text{EOP as } e, (\text{E0} \sqcup \text{E1P}) \text{ as } e, e.a_0 = v_0$	rev=-intro
7. $\text{elim } n \text{ from EMP as } e, n = e.name, \text{EOP as } e, (\text{E0} \sqcup \text{E1P}) \text{ as } e, e.a_0 = v_0, \text{P0 as } v_0$	$\bowtie$ -intro and 7
8. $\text{select } n \text{ from elim } v_0 \text{ P0 as } v_0, \text{elim } n, v_0 \text{ from EMP as } e, n = e.name, \text{EOP as } e, (\text{E0} \sqcup \text{E1P}) \text{ as } e, e.a_0 = v_0$	$\bowtie$ - $\epsilon$ -dist
9. $\text{select } n \text{ from P0 as } v_0, \text{elim } n, v_0 \text{ from EMP as } e, n = e.name, \text{EOP as } e, (\text{E0} \sqcup \text{E1P}) \text{ as } e, e.a_0 = v_0$	$\epsilon$ -elim
10. $\text{select } n \text{ from P0 as } v_0,$ $\quad \text{select } n, v_0 \text{ from elim } e, v_0 \text{ from EMP as } e, \text{EOP as } e, (\text{E0} \sqcup \text{E1P}) \text{ as } e, e.a_0 = v_0, \text{elim } n, e n = e.name$	$\bowtie$ - $\epsilon$ -dist
11. $\text{select } n \text{ from P0 as } v_0, n = e.name,$ $\quad \text{elim } e, v_0 \text{ from EMP as } e, \text{EOP as } e, (\text{E0} \sqcup \text{E1P}) \text{ as } e, e.a_0 = v_0$	$\epsilon$ -elim+administrative rules
12. $\text{select } n \text{ from P0 as } v_0,$ $\quad \text{elim } e, v_0 \text{ from EMP as } e, \text{EOP as } e, e.a_0 = v_0, \text{elim } e (\text{E0 as } e) \text{ union all } (\text{E1P as } e)$	$\cup$ -intro
13. $\text{select } n \text{ from P0 as } v_0, n = e.name,$ $\quad \text{elim } e, v_0 (\text{from EMP as } e, \text{EOP as } e, e.a_0 = v_0, \text{E0 as } e) \text{ union all}$ $\quad \quad (\text{from EMP as } e, \text{EOP as } e, e.a_0 = v_0, \text{E1P as } e)$	administrative rules
14. $\text{select } n \text{ from P0 as } v_0, n = e.name,$ $\quad (\text{elim } e, v_0 \text{ from EMP as } e, \text{EOP as } e, e.a_0 = v_0, \text{E0 as } e) \text{ union all}$ $\quad (\text{elim } e, v_0 \text{ from EMP as } e, \text{EOP as } e, e.a_0 = v_0, \text{E1P as } e)$	$\cup$ - $\epsilon$ -dist
15. $\text{select } n \text{ from P0 as } v_0, n = e.name,$ $\quad (\text{select } e, v_0 \text{ from EMP as } e, \text{EOP as } e, e.a_0 = v_0, \text{E0 as } e) \text{ union all}$ $\quad (\text{elim } e, v_0 \text{ from EMP as } e, \text{EOP as } e, e.a_0 = v_0, \text{E1P as } e)$	$3\times \bowtie$ - $\epsilon$ -dist+ $4\times \epsilon$ -elim
16. $\text{select } n \text{ from P0 as } v_0, n = e.name,$ $\quad (\text{select } e, v_0 \text{ from E0 as } e, e.a_0 = v_0) \text{ union all}$ $\quad (\text{elim } e, v_0 \text{ from EMP as } e, \text{EOP as } e, e.a_0 = v_0, \text{E1P as } e)$	$2\times \text{rev-}\bowtie$ -intro and 7, 8
17. $\text{select } n \text{ from P0 as } v_0, n = e.name,$ $\quad (\text{select } e, v_0 \text{ from E0 as } e, e.a_0 = v_0) \text{ union all}$ $\quad (\text{elim } e, v_0 \text{ from EMP as } e, \text{EOP as } e, e.a_0 = v_0,$ $\quad \quad \text{E1P as } e, (\text{E1} \sqcup \text{E2P}) \text{ as } e, e.a_1 = v_1, \text{P1 as } v_1)$	analogous to steps 3-7
18. $\text{select } n \text{ from P0 as } v_0, n = e.name,$ $\quad (\text{select } e, v_0 \text{ from E0 as } e, e.a_0 = v_0) \text{ union all}$ $\quad (\text{elim } e, v_0 \text{ from EMP as } e, \text{EOP as } e, e.a_0 = v_0,$ $\quad \quad \text{E1P as } e, (\text{E1} \sqcup \text{E2P}) \text{ as } e, e.a_1 = v_1, \text{P1 as } v_1, v_1.p = v_0)$	=-intro
19. $\text{select } n \text{ from P0 as } v_0, n = e.name,$ $\quad (\text{select } e, v_0 \text{ from E0 as } e, e.a_0 = v_0) \text{ union all}$ $\quad (\text{elim } e, v_0 \text{ from EMP as } e, \text{EOP as } e,$ $\quad \quad \text{E1P as } e, (\text{E1} \sqcup \text{E2P}) \text{ as } e, e.a_1 = v_1, \text{P1 as } v_1, v_1.p = v_0),$	rev=-intro
20. $\text{select } n \text{ from P0 as } v_0,$ $\quad (\text{select } e, v_0 \text{ from E0 as } e, e.a_0 = v_0) \text{ union all}$ $\quad (\text{select } e, v_0 \text{ from P1 as } v_1, v_1.p = v_0,$ $\quad \quad (\text{select } e, v_1 \text{ from E1 as } e, e.a_1 = v_1) \text{ union all}$ $\quad \quad (\text{select } e, v_1 \text{ from P2 as } v_2, v_2.p = v_1, \text{E}_2 \text{ as } e, e.a_2 = v_2)),$ $n = e.name$	analogous to steps 8-19 applied two additional times

Figure 6: QUERY COMPILATION.

*ciples of Knowledge Representation and Reasoning (KR'98)*, pages 2–13, 1998.

- [Calvanese *et al.*, 1998b] Diego Calvanese, Maurizio Lenzerini, and Daniele Nardi. Description Logics for Conceptual Data Modelling. In Jan Chomicki and Gunter Saake, editors, *Logics for Databases and Information Systems*, chapter 8. Kluwer, 1998.
- [Calvanese *et al.*, 2000] Diego Calvanese, Giuseppe De Giacomo, and Maurizio Lenzerini. Answering Queries Using Views over Description Logics Knowledge Bases. In *Proc. of the 16th Nat. Conf. on Artificial Intelligence (AAAI 2000)*, pages 386–391, 2000.
- [Calvanese *et al.*, 2001a] Diego Calvanese, Giuseppe De Giacomo, and Maurizio Lenzerini. Description Logics for Information Integration. In A. Kakas and F. Sadri, editors, *Computational Logic: From Logic Programming into the Future (In honour of Bob Kowalski)*, Lecture Notes in Computer Science. Springer-Verlag, 2001. To appear.
- [Calvanese *et al.*, 2001b] Diego Calvanese, Giuseppe De Giacomo, and Maurizio Lenzerini. Identification Constraints and Functional Dependencies in Description Logics. In *Proc. of the 17th Int. Joint Conf. on Artificial Intelligence (IJCAI 2001)*, pages 155–160, 2001.
- [Calvanese *et al.*, 2001c] Diego Calvanese, Giuseppe De Giacomo, and Maurizio Lenzerini. Ontology of Integration and Integration of Ontologies. In *Proc. of the 2001 Description Logic Workshop (DL 2001)*, pages 10–19. CEUR-WS Vol.49, 2001.

- [Horrocks *et al.*, 2000] Ian Horrocks, Ulrike Sattler, Sergio Tessaris, and Stephan Tobies. How to decide Query Containment under Constraints using a Description Logic. In *Proceedings of the 7th International Workshop on Knowledge Representation meets Databases (KRDB 2000)*, CEUR-WS vol.29, pages 59–72, 2000.
- [Khizder *et al.*, 2000] Vitaliy L. Khizder, David Toman, and Grant E. Weddell. Reasoning about Duplicate Elimination with Description Logic. In *Rules and Objects in Databases, DOOD 2000 (part of Computational Logic 2000)*, pages 1017–1032, 2000.
- [Khizder *et al.*, 2001] Vitaliy L. Khizder, David Toman, and Grant E. Weddell. On Decidability and Complexity of Description Logics with Uniqueness Constraints. In *International Conference on Database Theory ICDT'01*, pages 54–67, 2001.
- [Machtey and Young, 1978] Michael Machtey and Paul Young. *An Introduction to the General Theory of Algorithms*. North-Holland Amsterdam, 1978.
- [Toman and Weddell, 2001a] David Toman and Grant E. Weddell. On Attributes, Roles, and Dependencies in Description Logics and the Ackermann Case of the Decision Problem. In *Description Logics 2001*, pages 76–85. CEUR-WS vol.49, 2001.
- [Toman and Weddell, 2001b] David Toman and Grant E. Weddell. Query Processing in Embedded Control Programs. In *2nd Int. Workshop on Databases in Telecommunications*, pages 68–87. Springer LNCS 2209, 2001.