

# The Logical Data Model

GABRIEL M. KUPER IBM T. J. Watson Research Center and MOSHE Y. VARDI IBM Almaden Research Center

We propose an object-oriented data model that generalizes the relational, hierarchical, and network models. A database scheme in this model is a directed graph, whose leaves represent data and whose internal nodes represent connections among the data. Instances are constructed from objects, which have separate names and values. We define a logic for the model, and describe a nonprocedural query language that is based on the logic. We also describe an algebraic query language and show that it is equivalent to the logical language.

Categories and Subject Descriptors: H.2.1 [Database Management]: Logical Design—data models; H.2.3 [Database Management]: Languages

General Terms: Languages, Theory

Additional Key Words and Phrases: Algebra, database schema, logic, relational database, tuple calculus

# 1. INTRODUCTION

Research in database theory during the 1970's and the early 1980's has focused mainly on the relational model [11], due to its elegance and mathematical simplicity. This very simplicity, however, has gradually been recognized as one of the major disadvantages of the relational model: it forces the stored data to have a flat structure that real data does not always have [12, 46]. This has motivated a great deal of research during the past decade on *structured* data models: the so-called *semantic data models* [21, 44], *nested relations* [15, 27], and *complex objects* [9]. The reader is referred to [19, 20, 41] for excellent surveys.

Two works that we found particularly inspiring are by Jacobs [25, 26] and by Hull and Yap [24]. Jacobs describes "database logic," a mathematical

© 1993 ACM 0362-5915/93/0900-0379 \$01.50

ACM Transactions on Database Systems, Vol. 18, No. 3, September 1993, Pages 379-413.

A preliminary version of this paper, under the title "A New Approach to Database Logic," appeared in *Proceedings of the 3rd ACM Symposium on Principles of Database Systems*, Waterloo, April 1984, pp. 86–96. For a more extensive coverage of the material presented here the reader is referred to the first author's Ph.D. dissertation. The Logical Data Model: A New Approach to Database Logic, Dept. of Computer Science, Stanford University, 1985.

Authors' addresses: G. M. Kuper, ECRC, Arabellastr. 17, 81925 Munich, Germany; M. Y. Vardi, IBM Almaden Research Center, San Jose, California.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

model for databases that claims to generalize all three principal data models. Hull and Yap [24] describe the "format model," which generalizes the relational and hierarchical models. In the format model, database schemes are viewed as trees, in which each leaf represents data and each internal node represents some connection among the data.

Both these models are unsatisfactory in their ability to restructure data, i.e., the ability to query the database. While Hull and Yap ignore the issue of a data manipulation language, Jacobs' treatment is an overkill—his query language enables one to write noncomputable queries [51].

Furthermore, both approaches fail to model a significant aspect of hierarchical and network database management systems, which is the ability to use *virtual records*. Virtual records are essentially pointers to physical records, and they are used to avoid redundancy in the database [49]. Note that virtual records introduce cyclicity not only in the schema level but also at the instance level.

In the model we propose here a database scheme is an arbitrary directed graph. As in the format model, leaves (i.e., nodes with no outgoing edges) represent data, and internal nodes represent connections among the data. While it is not hard to model cyclicity at the schema level, it is not quite apparent how to do it at the instance level without running into cyclic definitions. Our solution is to distinguish between object names and object values, or, equivalently, between the address space and the data space. This distinction goes back to Codd's notions of *surrogates* [12]. Thus, instances in our model consist of objects which have separate names and values. This enables us to give semantics to instances in a well-defined way.

A data model consists of several components [47]. The first is the database structure mentioned above, which describes the static portion of the database. The second component is a way to specify integrity constraints on the database that restrict the allowed instances of the schema. We describe a logic in which integrity constraints can be specified. Our logic is inspired by Jacob's *database logic* [26], but unlike database logic, our logic is effective. That is, given a database and a sentence in the logic, one can test effectively whether the sentence is true in the database or not.

The third component is a way to restructure data in order to describe user views, queries, and so on. We describe two such mechanisms, a logical, i.e., nonprocedural, query language and an algebraic, i.e., procedural, query language that are analogous to Codd's relational calculus and relational algebra; we prove that these two languages are equivalent. These languages have a novel feature: not only can they access a nonflat data structure, e.g., a hierarchy, but the answers they produce do not have to be flat either. Thus, the language really does have the ability to restructure data and not only to retrieve it, as opposed to the other approaches cited above.

# 2. INTRODUCTION TO THE LOGICAL DATA MODEL

The *logical data model* (LDM) is a generalization of Hull and Yap's format model [24]. The format model fails to model an important part of network and

hierarchical database systems, namely the ability to use virtual records. We model this by introducing cyclicity into the database schemas. An LDM schema is a labeled directed multigraph.<sup>1</sup> Each node has a particular *type*. The leaves of the schema (i.e., nodes with no outgoing edges) are all of the *basic* type, denoted graphically by  $\Box$ . The instance of each node contains the data values stored in the database. Each interior node has one of the following types:

- (1) *Product*, denoted graphically by  $\otimes$ . The domain of such a node is the Cartesian product of the domains of its children.
- (2) Power, denoted graphically by (2). Such a node has exactly one child. The domain of such a node is the set of all finite subsets of the domain of its child.
- (3) Union, denoted graphically by  $\odot$ . The domain of such a node is the disjoint union of the domains of its children.

*Example* 1. Figure 1 shows a genealogy database as a relation. We can represent the structure of this relation by the LDM schema in Figure 2. It consists of two nodes u and v of type  $\Box$  that correspond to the *Person* and *Parent* attributes, respectively, and one node w of type  $\otimes$  that contains pairs of related attributes.

For the moment, an instance I of an LDM schema will be an assignment to each node u of a set I(u) of values from the corresponding domain. An instance corresponding to the data in Figure 1 consists of the following assignments:

 $I(u) = \{\text{Rehoboam, Solomon, David}\}\$ 

 $I(v) = \{$ Solomon, David, Batsheba, Jesse $\}$ 

 $I(w) = \{(\text{Rehoboam}, \text{Solomon}), (\text{Solomon}, \text{David}), (\text{Solomon}, \text{Batsheba}), \}$ 

(David, Jesse)}

*Example* 2. The genealogy could be represented by the network in Figure 3. In this network there are two record types, *Person*, containing the names of the people in the database, and a dummy record, *PP*. There are two links (sets) that connect each dummy record to a person and his parents.

The idea behind the mapping from the network to the LDM schema in Figure 4 is as follows. Each record type  $R_i$  is mapped into a product node  $v_{R_i}$ . For each field of  $R_i$ ,  $v_{R_i}$  has a child of type  $\Box$ . For each link (set) in the network with  $R_i$  as a member, let  $R_j$  be the owner of the link. Then  $v_{R_j}$  is a child of  $v_R$ .

In Figure 4, w is  $v_{PP}$  and v is  $v_{Person}$ . u corresponds to the field of the *Person* record, i.e., the person's name, and the two arcs from w to v correspond to the two links.

<sup>&</sup>lt;sup>1</sup>In a *multigraph* one can have more than one edge between two nodes.

ACM Transactions on Database Systems, Vol. 18, No. 3, September 1993.

Fig

		Person	Parent
		Rehoboam	Solomon
1.	The Person-Parent relation.	Solomon	David
		Solomon	Batsheba
		David	Jesse

Fig. 2 The Person-Parent relation as an LDM schema.





Fig. 3. The genealogy as a network.

If the network had the same contents as the relation in Figure 1, the corresponding instance of the LDM schema in Figure 4 would be

$$\begin{split} I(u) &= \{ \text{Rehoboam, Solomon, David, Batsheba, Jesse} \} \\ I(v) &= \{ (\text{Rehoboam}), (\text{Solomon}), (\text{David}), (\text{Batsheba}), (\text{Jesse}) \} \\ I(w) &= \{ ((\text{Rehoboam}), (\text{Solomon})), ((\text{Solomon}), (\text{David})), \\ &\quad ((\text{Solomon}), (\text{Batsheba})), (((\text{David}), (\text{Jesse}))) \} \end{split}$$

*Example* 3. Figure 5 shows a hierarchical representation of the genealogy. In this hierarchy, each *Person* record is related to the linked list of his parents. Even though the hierarchical model uses linked lists, this is really just a matter of the implementation. Intuitively, the user should see only the connection between a person and the set of his parents. We therefore map each record type  $R_i$  into a product node  $v_{R_i}$  as we did for the network model, with a child of type  $\Box$  corresponding to each of  $R_i$ 's fields. However, if  $R_i$  is a member of the link  $(R_i, R_j)$ , then instead of connecting  $v_{R_i}$  to  $v_{R_j}$  directly, we connect them through a node of type  $\circledast$ . The corresponding LDM schema is then the cyclic schema in Figure 6.



When we have cyclic schemas, such as this one, we cannot define an instance in the same way. Writing down the instance, in this example, is very complicated, and in general, when the data as well as the schema is cyclic, we cannot write it down at all. This is similar to one of the problems with Jacobs' database logic. The mathematical theory we develop to deal with this problem is closely related to the non-well-founded sets of [5]. Our approach to defining an instance of a schema is to separate the concepts of object *name* and object *value*. Intuitively, an object name is an address and the object value is the content of that address. An instance  $\mathbf{I}$  then consists of

- (1) an assignment of a set I(u) of object names to each node u of the schema, and
- (2) an assignment of an object value val(l) to each name l in I(u).

I(u)		I(v)		I(w)		
l	val(l)	l	val(l)		l	val(l)
1	Rehoboam	4	Solomon		8	(1, 4)
2	Solomon	5	David		9	(2, 5)
3	David	6	Batsheba		10	(2, 6)
		7	Jesse		11	(3,7)

Fig. 7. Instance of the LDM schema that corresponds to a relation

I(u)		I(v)		I(w)	
l	val(l)	l	val(l)	l	val(l)
1	Rehoboam	6	(1, 11)	11	{7}
2	Solomon	7	(2, 12)	12	{8,9}
3	David	8	(3, 13)	13	{10}
4	Batsheba	9	(4, 14)	14	Ø
5	Jesse	10	(5, 14)		

Fig. 8. Instance of the LDM schema that corresponds to a hierarchy.

The names are taken from a fixed infinite set L, which will usually be the set of natural numbers. Values are taken from a fixed infinite set D of data values, or are built up from object names.

We now show what some of the instances in the previous examples look like when we use names and values.

*Example* 4. The relational instance in Example 1 consists of the following assignment of names to nodes.

$$I(u) = \{1, 2, 3\},\$$
  
 $I(v) = \{4, 5, 6, 7\},\$  and  $I(w) = \{8, 9, 10, 11\}.$ 

We then assign a value val(l) to each of these names. This assignment is shown in Figure 7.

*Example* 5. In Figure 8 we show the instance that corresponds to the hierarchy of Example 3.



Fig. 9. Nodes in LDM schemas.

## 3. FORMAL DESCRIPTION OF THE LOGICAL DATA MODEL

Definition 1. A schema is a tuple  $\mathbf{S} = \langle V, E, \langle \mu \rangle$ , where

- (1) (V, E) is a directed multigraph;
- (2) < is a total order on E;
- (3) μ is a function from V to the set of types {□, ⊗, ∞}, that satisfies the following conditions (see Figure 9):
  - (a)  $\mu(v) = \Box$  iff v is a leaf;
  - (b) if  $\mu(v) = \circledast$ , then v has exactly one child;
  - (c) if  $\mu(v) = \odot$ , then the children of v are distinct nodes (if  $\mu(v) = \otimes$ , however, there can be multiple edges from v to a node w).

The order on the edges is used to induce an order on the components of tuples.  $\mu(v)$  is called the *type* of v. For readability, we use the following abbreviations:

- (1)  $\mu(v) = (\circledast, w)$  is an abbreviation for " $\mu(v) = \circledast$  and its child is w."
- (2) (a)  $\mu(v) = (\otimes, n)$  is an abbreviation for " $\mu(v) = \otimes$  and v has n children."
  - (b)  $\mu(v) = (\otimes, n, v_1, \dots, v_n)$  is an abbreviation for " $\mu(v) = \otimes$ , there are exactly *n* edges  $e_1, \dots, e_n$  with tail *v*, these edges are in the order  $e_1 < \dots < e_n$  and their heads are  $v_1, \dots, v_n$ ."
- (3) (a)  $\mu(v) = (\odot, n)$  is an abbreviation for " $\mu(v) = \odot$  and v has n children."
  - (b)  $\mu(v) = (\odot, n, v_1, \dots, v_n)$  is an abbreviation for " $\mu(v) = \odot$ , there are exactly *n* edges  $e_1, \dots, e_n$  with tail *v*, these edges are in the order  $e_1 < \cdots < e_n$  and their heads are  $v_1, \dots, v_n$ ."

We are really overloading the symbol  $\mu$ , but in practice this will not cause any confusion. Some other abbreviations that we use include referring to elements of V and E as nodes and edges, respectively, of **S**, and referring to < as an order on the children of a node of **S**. We ignore the order < when it is clear from the context, and we often refer to a schema as  $\langle V, E, \mu \rangle$ .

An instance of S consists of two parts: An assignment of a set of object names to each node of S, and an assignment of an object value to each object name.

*Proviso.* We assume a fixed infinite set L of object names and a fixed infinite set D of data values.

#### 386 • G. M. Kuper and M. Y. Vardi

Definition 2. An instance of **S** is a tuple  $\mathbf{I} = \langle I, val \rangle$  that satisfies the following:

- (1)  $I: V \to 2^L$  is an assignment of sets of object names to nodes. We require that I(v) and I(w) be disjoint whenever v and w are distinct nodes of **S**.
- (2) val is a mapping with domain  $\bigcup_{v \in V} I(v)$ , i.e., from the set of all the names that are in the instance. The mapping val must satisfy the following:
  - (a) If  $\mu(v) = \Box$  and  $l \in I(v)$ , then val(l) is a member of the set D of data values.
  - (b) If  $\mu(v) = (\otimes, n, v_1, \dots, v_n)$  and  $l \in I(v)$ , then val(l) is a tuple  $(l_1, \dots, l_n)$  such that for each  $i, 1 \le i \le n, l_i$  is a element of  $I(v_i)$ .
  - (c) If  $\mu(v) = (\circledast, w)$  and  $l \in I(v)$ , then val(l) is a finite subset of I(w).
  - (d) If  $\mu(v) = (\bigoplus, n, v_1, \dots, v_n)$  and  $l \in I(v)$ , then  $val(l) \in I(v_1) \cup \dots \cup I(v_n)$ .

Definition 3. A finite instance of **S** is an instance  $\mathbf{I} = \langle I, val \rangle$  of **S** such that for each node v of **S**, I(v) is finite.

Here we are mostly interested in finite instances, since these correspond to real databases. If l is in  $\bigcup_{v \in V} I(v)$ , we say that it is a name *in* **I**, and val(l) is called *its* value. The set  $\bigcup_{v \in V} val[I(v)]$  is called the set of values *in* **I**.

Definition 4. Let **I** be an instance of the schema **S**, and let v be a node of **S** of type  $(\otimes, n, v_1, \ldots, v_n)$ . Let l be any name in I(v). If  $1 \le i \le n$ , then  $\prod_i(l)$  will be the *i*th component of val(l). We also use the notation  $\prod_{v_i}(l)$  for this component, whenever this does not result in ambiguity.

The following definition explains when it is meaningful to compare two names, that is, if v and w are nodes of **S**,  $l_1 \in I(v)$  and  $l_2 \in I(w)$ , is it possible for  $l_1$  and  $l_2$  to have the same value?

Definition 5. We say that two nodes v and w in a schema **S** are *similar* iff they are of the same type and have the same children, that is, if one of the following holds:

- (1)  $\mu(v) = \mu(w) = \Box;$
- (2) for some node u,  $\mu(v) = \mu(w) = \circledast, u$ ;
- (3) for some n and nodes  $u_1, \ldots, u_n$ ,  $\mu(v) = \mu(w) = (\otimes, n, u_1, \ldots, u_n)$ ;
- (4) for some n and nodes  $u_1, \ldots, u_n, \mu(v) = \mu(w) = (\bigoplus, n, u_1, \ldots, u_n)$ .

A query on an LDM schema S will involve the addition of some nodes to S. For this we need the following definitions.

Definition 6. Let  $\mathbf{S} = \langle V, E, \langle \mu \rangle$  be a schema.  $\mathbf{S}' = \langle V', E', \langle \mu' \rangle$  is an extension of  $\mathbf{S}$  iff

(1)  $V \subseteq V';$ (2) (a)  $E \subseteq E',$ 

- (b) if (v, v<sub>2</sub>) ∈ E' − E, then v<sub>1</sub> is in V', i.e. all new edges are either between new nodes, or from a new node to a node in V;
- (3)  $<'|_{E \times E} = <$ , i.e., <' is a conservative extension of <;
- (4)  $\mu'|_V = \mu$ , i.e.,  $\mu'$  is a conservative extension of  $\mu$ .

The intuition behind (1) and (2) is that the schema S' adds some extra nodes (V' - V) to V, and adds some edges (E' - E) to E, but it does not change the schema S, i.e., it does not add edges between nodes in V or from nodes in V to new nodes.

Let S' be an extension of S. We define an extension of I to an instance of S' as follows.

*Definition* 7. Let **S'** be an extension of **S**, and let  $\mathbf{I} = \langle I, val \rangle$  be an instance of **S**. We say that an instance  $\mathbf{I'} = \langle I', val' \rangle$  of **S'** is an extension of **I** to **S'** iff

(1) for all v in V, I'(v) = I(v);

(2) if v is a node of **S** and  $l \in I(v)$ , then val'(l) = val(l).

The proof of the following lemma is straightforward.

LEMMA 1. Let  $\mathbf{S}'$  be an extension of  $\mathbf{S}$ , and let  $\mathbf{I}'$  be an instance of  $\mathbf{S}'$ . Then there is a unique instance  $\mathbf{I}$  of  $\mathbf{S}$  such that  $\mathbf{I}'$  is the extension of  $\mathbf{I}$  to  $\mathbf{S}'$ . This instance is called the restriction of  $\mathbf{I}'$  to  $\mathbf{S}$ .

We conclude this section with a definition of isomorphism. Two instances are isomorphic if they are essentially the same, i.e., if they differ only by renaming. As we want to show that the result of a query is well defined up to isomorphism, we give a stronger definition of isomorphism. Let I be an instance of S, let S' be an extension of S, and let  $I_1$  and  $I_2$  be extensions of I to S'. We say that  $I_1$  and  $I_2$  are isomorphic relative to S if there is an isomorphism between  $I_1$  and  $I_2$  that leaves the objects of I fixed. In the case of a query, this means that an isomorphism relative to the database leaves the contents of the database fixed.

Definition 8. Let **S**' be an extension of **S**, and let  $\mathbf{I} = \langle I, val \rangle$  be an instance of **S**. Let  $\mathbf{I}_1 = \langle I_1, val_1 \rangle$ , and  $\mathbf{I}_2 = \langle I_2, val_2 \rangle$  be two extensions of **I** to **S**'. We say that  $\mathbf{I}_1$  and  $\mathbf{I}_2$  are *isomorphic relative to* **S** iff there is a mapping

$$g: \bigcup_{v \in \mathbf{S}'} I_1(v) \xrightarrow[]{\text{onto}} \bigcup_{v \in \mathbf{S}'} I_2(v)$$

such that

- (1) for each node v of **S**, g is the identity on I(v);
- (2) for each node v of S', g maps  $I_1(v)$  onto  $I_2(v)$ ;
- (3) if v is a node of **S**' and  $l \in I_1(v)$ , then
  - (a) if v is of type  $\Box$ , then  $val_2(g(l)) = val_1(l)$ ;

(b) if v is of type  $(\otimes, n)$ , then

$$val_2(g(l)) = (g(\Pi_1(val_1(l))), \dots, g(\Pi_n(val_1(l))))$$

- (c) if v is of type  $\odot$ , then  $val_2(g(l)) = g(val_1(l))$ ;
- (d) if v is of type  $\circledast$ , then  $g[val_2(l)] = val_1[g(l)]$ .

As a special case of this definition we get the definition of ordinary isomorphism.

Definition 9. Let  $\mathbf{I}_1 = \langle I_1, val_1 \rangle$  and  $\mathbf{I}_2 = \langle I_2, val_2 \rangle$  be instances of **S**. We say that  $\mathbf{I}_1$  and  $\mathbf{I}_2$  are *isomorphic* iff they are isomorphic relative to the empty schema, i.e., the schema with  $V = E = \mu = \emptyset$ .

### 4. LDM LOGIC

In this section we define the LDM logic. The logic is similar to relational tuple calculus, and will be used as part of the logical query language. The logic can also be used to specify integrity constraints on LDM schemas and to define views. Throughout this section  $\mathbf{S} = \langle V, E, \mu \rangle$  will be a fixed schema, and  $\mathbf{I} = \langle I, val \rangle$  a fixed instance of  $\mathbf{S}$ .

Each variable in the logic has a fixed *sort*, where the sorts are the elements of V. The sorts define the domains over which the variables range. For example, if x is a variable of sort v, then x ranges over I(v). The analog to this in relational calculus is a tuple variable that ranges over a specific relation. We usually write a variable with its sort as a subscript, e.g.,  $x_v$ . Variables with different subscripts denote distinct variables, so that  $x_u$  is a different variable from  $x_v$ . Even though variables range over object names, we think of them as ranging over objects. Thus, we refer to "the name of  $x_v$ " and to "the value of  $x_v$ ."

Definition 10. The atomic formulas over  $\mathbf{S}$  are the following:

- (1)  $x_v \pi_t y_w$ , where w is a node of type  $\otimes$  and v is its tth child;
- (2)  $x_v \rho y_w$ , where w is a node of type  $\odot$  and v is one of its children;
- (3)  $x_v \in y_w$ , where w is of type  $(\circledast, v)$ ;
- (4)  $x_v =_n y_v;$
- (5)  $x_v =_v y_w$ , where v and w are similar nodes;
- (6)  $x_v =_v d$ , where d is a data element in D, and v is of type  $\Box$ .

The atomic formula  $x_v \pi_t y_w$  means that the name of  $x_v$  is the *t*th component of the value of  $y_w$ . Note that we have to mention which component of w we are referring to, since there may be multiple edges from w to v. However, we also write  $x_v \pi y_w$  when this is unambiguous. The atomic formula  $x_v \rho y_w$  means that the value of  $y_w$  is  $x_v$ . Since there is only one edge from w to v, we use  $\rho$  rather than  $\rho_t$ . The atomic formula  $x_v \in y_w$  means that  $x_v$  is a member of the value of  $y_w$ .

There are several kinds of equality. The atomic formula  $x_v =_n y_v$  means that the names of  $x_v$  and  $y_v$  are equal. Since I(v) and I(w) are disjoint whenever  $v \neq w$ , we do not need atomic formulas of the form  $x_v =_n y_w$  for  $v \neq w$ . The atomic formula  $x_v =_v y_w$  means that the values of  $x_v$  and  $y_w$ are equal. This is meaningful only when v and w are similar nodes. Finally, the atomic formula  $x_v =_v d$  means that the data value of  $x_v$  is equal to the data element d.

Definition 11. A well-formed LDM formula over a schema S is

- (1) an atomic formula;
- (2)  $\phi_1 \lor \phi_2$ , where  $\phi_1$  and  $\phi_2$  are well-formed formulas;
- (3)  $\neg \phi_1$ , where  $\phi_1$  is a well-formed formula;
- (4)  $(\forall x_v)\phi_1$ , where  $\phi_1$  is a well-formed formula.

The free variables of  $\phi$  are defined in the same way as in first-order logic.

We use  $\phi_1 \wedge \phi_2, (\exists x_v)\phi, \phi_1 \Rightarrow \phi_2$  and  $\phi_1 \Leftrightarrow \phi_2$  with the standard meanings. Another useful abbreviation is the following.

Definition 12. The formula " $x_v =_v (x_{v_1}^1, \ldots, x_{v_n}^n)$ " where v is a node of type  $(\otimes, n, v_1, \ldots, v_n)$  means " $x_{v_1}^1 \pi_1 x_v \wedge \cdots \wedge x_{v_n}^n \pi_n x_v$ ."

We now define satisfaction of LDM formulas. Let  $\phi(x_{v_1}^1, \ldots, x_{v_n}^n)$  be an LDM formula whose free variables are  $x_{v_1}^1, \ldots, x_{v_n}^n$ . Let  $l_1, \ldots, l_n$  be an assignment of object names to the free variables in the formula, that is, each  $l_i$  is a member of the corresponding  $I(v_i)$ .  $\models_{\mathbf{I}} \phi(l_1, \ldots, l_n)$  means that  $\phi$  is satisfied by  $l_1, \ldots, l_n$  in the instance **I**. When **I** is clear from the context, we write  $\models$  instead of  $\models_{\mathbf{F}}$ 

Definition 13. Let  $\phi(x_{v_1}^1, \ldots, x_{v_n}^n)$  be a formula with free variables  $x_{v_1}^1, \ldots, x_{v_n}^n$ , and let  $l_i \in I(v_i)$  for all  $i, 1 \le i \le n$ . Then  $\vDash_{\mathbf{I}} \phi(l_1, \ldots, l_n)$  iff the following hold:

(1) if  $\phi$  is  $x_v^i \pi_t y_w^j$ , then  $\models_{\mathbf{I}} (x_v^i \pi_t x_w^j)(l_1, \ldots, l_n)$  iff  $l_i = \prod_t (l_i)$ ;

(2) if 
$$\phi$$
 is  $x_n^i \rho y_n^j$ , then  $\models_{\mathbf{r}} (x_n^i \rho x_n^j)(l_1, \ldots, l_n)$  iff  $l_i = val(l_i)$ ;

(2) If  $\phi$  is  $x_v p y_w$ , then  $\vdash_{\mathbf{I}} (x_v p x_w (t_1, \dots, t_n) \text{ if } t_i = bat(t_j)$ , (3) if  $\phi$  is  $x_v^i \in x_w^j$ , then  $\models_{\mathbf{I}} (x_v^i \in x_w^j)(t_1, \dots, t_n)$  iff  $t_i \in val(t_j)$ ;

(4) if  $\phi$  is  $x_v^i =_n x_v^j$ , then  $\vDash_{\mathbf{I}} (x_v^i =_n x_v^j)(l_1, \ldots, l_n)$  iff  $l_i = l_i$ ;

(5) if  $\phi$  is  $x_v^i =_v x_w^j$ , then  $\vDash_{\mathbf{I}} (x_v^i =_v x_w^j)(l_1, \ldots, l_n)$  iff  $val(l_i) = val(l_i)$ ;

(6) if  $\phi$  is  $x_v^i = d$ , then  $\vDash_I (x_v^i = d)(l_1, \dots, l_n)$  iff  $val(l_i) = d$ ;

- (7)  $\vDash_{\mathbf{I}} (\phi_1 \lor \phi_2)$  iff  $\vDash_{\mathbf{I}} \phi_1$  or  $\vDash_{\mathbf{I}} \phi_2$ ;
- (8)  $\models_{\mathbf{I}} \neg \phi$  iff  $\models_{\mathbf{I}} \phi$  does not hold;
- (9) if  $\phi$  is a formula with free variables  $x_{v_1}^1, \ldots, x_{v_w}^n, y_w$ , then

$$\models_{\mathbf{I}} ((\forall y_w)\phi)(l_1,\ldots,l_n) \text{ iff for all } l \in I(w), \quad \models_{\mathbf{I}} \phi(l_1,\ldots,l_n,l).$$

Definition 14. Let  $\phi$  be an LDM sentence. We say that I satisfies  $\phi$  iff  $\models_{\mathbf{I}} \phi$  holds.

*Example* 6. This example and the next one will be over the LDM schema of Figure 6 with the instance of Figure 8. The LDM formula  $\phi(x_u, y_v) =$ 

ACM Transactions on Database Systems, Vol. 18, No. 3, September 1993.

#### 390 . G. M. Kuper and M. Y. Vardı

 $(x_u \pi_1 y_v)$  says that the name of  $x_u$  is equal to the first component of the value of  $y_v$ .  $\models_{\mathbf{I}} \phi(l_1, l_2)$  holds for the  $(l_1, l_2)$  pairs (1, 6), (2, 7), (3, 8), (4, 9), and (5, 10).

*Example* 7. The following constraint says that each name in u is related to exactly one set in w. For example. "8" and "9" as parents of "2" must be in one set rather that in two different sets. The constraint is

$$\phi = (\forall x_u)(\forall y_v^1)(\forall y_v^2)(\forall z_w^1)(\forall z_w^2)$$
$$\left(y_v^1 =_v (x_u, z_w^1) \land y_v^2 =_v (x_u, z_w^2) \Rightarrow z_w^1 =_n z_w^2\right).$$

In other words, each name in  $u(x_u)$  has at most one name in  $w(z_w^1 \text{ and } z_w^2)$  associated with it. This association is through  $y_v^1$  and  $y_v^2$ .

Note that this constraint says that each *name* in u is associated with at most one set in w, rather than saying that each *person* in the database is associated with at most one such set. There could still be duplication in u, e.g., two names with the data value "Solomon." One way to prevent this is through the constraint

$$\psi = (\forall x_u^1)(\forall x_u^2)(x_u^1 = x_u^2) \Rightarrow x_u^1 = x_u^2).$$

The following lemma shows that a slightly restricted logic has the same power. This restricted logic does not have atomic formulas that compare data values of internal nodes. This lemma makes subsequent proofs and definitions simpler.

LEMMA 2. Let  $\phi(x_{v_1}^1, \ldots, x_{v_n}^n)$  be an LDM formula whose free variables are the variables  $x_{v_1}^1, \ldots, x_{v_n}^n$ . There is an LDM formula  $\phi(x_{v_1}^1, \ldots, x_{v_n}^n)$  with the same free variables that does not contain any atomic subformula of the form  $x_u =_v y_w$  with  $\mu(v)$  and  $\mu(w)$  different from  $\Box$ , such that  $\psi$  is equivalent to  $\phi$ . That is, for all instances **I** of **S** and all  $l_1, \ldots, l_n, l_i \in$  $I(v_i), \models_{\mathbf{I}} \phi(l_1, \ldots, l_n)$  iff  $\models_{\mathbf{I}} \psi(l_1, \ldots, l_n)$ .

**PROOF.** The proof is by induction on the size of  $\phi$ . We show how to construct  $\psi$  for formulas of the form  $x_u =_v y_v$ , where u and v are similar and not of type  $\Box$ . The result then follows immediately.

We distinguish between the possible types of v and w.

(1) if u and v are of type ((a), w), then  $\psi(x_u, y_v)$  is  $(\forall z_w)(z_w \in x_u \Leftrightarrow z_w \in y_v)$ , where  $z_w$  is a new variable. Let **I** be an instance of **S**. Then

$$\models_{\mathbf{I}} (x_u =_v y_v)(l_1, l_2) \Leftrightarrow val(l_1) = val(l_2) \Leftrightarrow \text{ for all } l \text{ in } I(w), l \in val(l_1) \Leftrightarrow l \in val(l_2) \Leftrightarrow \models_{\mathbf{I}} ((\forall z_w)(z_w \in x_u \Leftrightarrow z_w \in y_v))(l_1, l_2),$$

and therefore  $\phi \Leftrightarrow \psi$  is valid.

(2) If u and v are of type  $(\circledast, n, w_1, \dots, w_n)$ , then  $\psi(x_u, y_n)$  is

$$(\forall z_{w_1}^1) \cdots (\forall z_{w_n}^n) ( (z_{w_1}^1 \pi_1 x_u \Leftrightarrow z_{w_1}^1 \pi_1 y_v) \land \cdots \land (z_{w_n}^n \pi_n x_u \Leftrightarrow z_{w_n}^n \pi_n y_v) )$$

ACM Transactions on Database Systems, Vol. 18, No 3, September 1993

where  $z_{w_1}^1, \ldots, z_{w_n}^n$  are *n* distinct new variables. Let **I** be an instance of **S**. Then  $\models_{\mathbf{I}} (x_u =_v y_v)(l_1, l_2)$  is equivalent to  $val(l_1) = val(l_2)$ . If  $val(l_1) = (l_1^1, \ldots, l_1^n)$  and  $val(l_2) = (l_2^1, \ldots, l_2^n)$ ,  $val(l_1) = val(l_2)$  is equivalent to the conjunction of  $l_1^i = l_2^i$  for  $i = 1, \ldots, n$ . In other words,  $val(l_1) = val(l_2)$  iff, for  $i = 1, \ldots, n$ 

$$\vDash_{\mathbf{I}} \Big( (\forall z_w^i) \Big( z_{w_i}^i \pi_i x_u \Leftrightarrow z_{w_i}^i \pi_i y_v \Big) \Big) (l_1, l_2).$$

Therefore,  $\vDash_{\mathbf{I}} (x_{y} =_{v} y_{y})(l_{1}, l_{2})$  is equivalent to

$$\vdash_{\mathbf{I}} (\forall z_{w_1}^1) \cdots (\forall z_{w_n}^n) ((z_{w_1}^1 \pi_1 x_u \Leftrightarrow z_{w_1}^1 \pi_1 y_v) \\ \wedge \cdots \wedge (z_{w_n}^n \pi_n x_u \Leftrightarrow z_{w_n}^n \pi_n y_v)) (l_1, l_2),$$

i.e.,  $\phi \Leftrightarrow \psi$  is valid.

(3) If u and v are of type  $(\odot, n, w_1, \dots, w_n)$ , then  $\psi(x_u, y_v)$  is

$$(\exists z_{w_1}^1)(z_{w_1}^1\rho x_u \wedge z_{w_1}^1\rho y_v) \vee \cdots \vee (\exists z_{w_n}^n)(z_{w_n}^n\rho x_u \wedge z_{w_n}^n\rho y_v),$$

where  $z_{w_1}^1, \ldots, z_{w_n}^n$  are *n* distinct new variables. Let **I** be an instance of **S**. Then  $\vDash_{\mathbf{I}} (x_u =_v y_v)(l_1, l_2)$  is equivalent to  $val(l_1) = val(l_2) = l$ . This holds iff for some  $i, 1 \le i \le n, l \in I(w_i)$ , in which case

$$\models_{\mathbf{I}} \left( \left( \exists z_{w_i}^i \right) \left( z_{u_i}^i \rho x_u \wedge z_{w_i}^i \rho y_v \right) \right) (l_1, l_2)$$

and once more  $\phi \Leftrightarrow \psi$  is valid.  $\Box$ 

From now on we assume that  $x_u =_v y_v$  can appear as a subformula only when  $\mu(v) = \mu(w) = \Box$ , as far as proofs are concerned, but will use the more general form when convenient.

The proof of the following lemma, which says that satisfaction is preserved under isomorphism, is straightforward.

LEMMA 3. Let **S**' be an extension of **S**, and let  $\mathbf{I}_1$  and  $\mathbf{I}_2$  be extensions of  $\mathbb{I}$  to **S**'. Let g be an isomorphism from  $\mathbf{I}_1$  and  $\mathbf{I}_2$  relative to **S**, and let  $\phi(x_{v_1}^1, \ldots, x_{v_n}^n)$  be an LDM formula. Then

$$\vDash_{\mathbf{I}_{1}} \phi(l_{1},\ldots,l_{n}) \Leftrightarrow \vDash_{\mathbf{I}_{2}} \phi(g(l_{1}),\ldots,g(l_{n})).$$

Finally, we show that our logic is effective over finite instances.

LEMMA 4. Let  $\phi(x_{v_1}^1, \ldots, x_{v_n}^n)$  be an LDM formula over **S** whose free variables are the variables  $x_{v_1}^1, \ldots, x_{v_n}^n$ . Let  $\mathbb{I}$  be a finite instance and let  $l_i \in I(v_i)$  for all  $i, 1 \le i \le n$ . Then  $\models_{\mathbf{I}} \phi(l_1, \ldots, l_n)$  can be determined effectively.

PROOF. We show this by induction on the size of the formula. For atomic formulas, testing for satisfaction is straightforward. Testing for disjunction and negation is also clearly effective. For quantification we make use of the

ACM Transactions on Database Systems, Vol. 18, No 3, September 1993

finiteness of **I**. In order to test whether  $\vDash_{\mathbf{I}} ((\forall y_w)\phi)(l_1,\ldots,l_n)$ , we test whether  $\vDash_{\mathbf{I}} \phi(l_1,\ldots,l_n,l)$  for each l in the finite set I(w).  $\Box$ 

# 5. THE LOGICAL QUERY LANGUAGE

In the relational model the result of a query is a relation. A natural generalization would be for the result of a query in our model to be another LDM schema, henceforth called the *query schema*, together with an instance of that schema. We modify this slightly by not requiring that the query's schema be an independent schema, but instead allowing the successors of nodes in the query to be nodes in the database schema.

Continuing the analogy with the relational calculus, the natural thing to do would be to let the query be some LDM formula  $\phi$ . The resulting instance would then consist of those objects that satisfy  $\phi$ .

There are two problems with this approach. First, while it is clear how an LDM formula can select objects that satisfy certain conditions, it is not clear how an LDM formula can construct new objects. One solution would be to prevent the query from referring directly to object names, but rather have the query refer only to object values. We could then find all possible values that might appear in the result, assign them arbitrary names, and show that the result is unique up to isomorphism.

This still does not solve the second problem, which is how to deal with cyclicity. Not only do we need the ability to refer directly to object names in order to deal with cycles, but even then the result of the query is not always defined uniquely. For example, if the query schema is that of Figure 10, and the query specifies that I(u) and I(v) each contain at least two different names, then there is no way to choose between the two incomparable instances in Figures 11 and 12. Our solution to this problem is to restrict the queries to ones that do not contain cycles, while allowing cyclicity in the database. Furthermore, we allow the query to refer explicitly to names only in nodes of the database. For a more detailed discussion of the motivation underlying our approach, see [31]. More recent work, e.g., IQL ([3]), has shown how query languages could be defined to allow cyclic queries.

# 5.1 The Query Language

*Definition* 15. Let  $\mathbf{S} = \langle V, E, \mu \rangle$  be an LDM schema. A query on  $\mathbf{S}$  consists of a tuple  $\mathbf{Q} = \langle \mathbf{S}_{\mathbf{Q}}, \Phi_{\mathbf{Q}}, \prec_{\mathbf{Q}} \rangle$  where

- (1)  $\mathbf{S}_{\mathbf{Q}}$  is an extension of  $\mathbf{S}$ ;
- (2)  $\prec_{\mathbf{Q}}$  is a topological order on the nodes in  $V_{\mathbf{Q}} V$ , i.e.,  $\prec_{\mathbf{Q}}$  is a linear order such that if v is a child of w then  $v \prec_{\mathbf{Q}} w$ ;
- (3)  $\Phi_{\mathbf{Q}}$  is a set of pairs  $\langle v, \phi_v \rangle$  that assigns a formula  $\phi_v$  to the node v, for each node v in  $V_{\mathbf{Q}} V$ . The formula  $\phi_v$  that corresponds to the node v satisfies:

(a)  $\phi_v$  has only one free variable, and it is of sort v.

(b) all other variables in  $\phi_v$  are bound. Each of their sorts is either a node of the database schema **S** or is a query node that precedes v under  $\prec_{\mathbf{Q}}$ .



Fig. 10. LDM schema.

$$I_1(v)$$
  $I_1(v)$ 

l	val(l)	l	val(l)
1	(3)	3	(1)
2	(4)	4	(2)

Fig. 11. A possible result of the query.

$I_2(u)$		$I_2(v)$	
l	val(l)	l	val(l)
1	(3)	3	(2)
2	(4)	4	(1)

Fig. 12. Another possible result of the query.

The order  $\prec_{\mathbf{Q}}$  is used to specify the order in which we define the result of the query. Each formula  $\phi_v$  specifies the contents of v in terms of database nodes and of query nodes that precede v.

Before continuing with the formal details, we give several examples of queries. The database schema in all these examples is the genealogy schema of Figure 6 with the instance of Figure 8.

*Example* 8. The schema of  $\mathbf{Q}_1$  is shown in Figure 13. The formula  $\phi_{u'}(x_{u'})$  is

$$(\exists y_u)(x_{u'} =_v y_u).$$

In other words, we want I(u') to be a copy of I(u). (We eliminate, however, any duplication that may be in I(u).) The result of the query is shown in Figure 14.<sup>2</sup>

 $<sup>^{2}</sup>$  In all these examples, the result is defined only up to isomorphism relative to S, i.e., the choice of names is arbitrary.

ACM Transactions on Database Systems, Vol. 18, No. 3, September 1993.





I(u')



*Example* 9. The schema of  $\mathbf{Q}_2$  is shown Figure 15. We want v' to contain the set of parents of Solomon, so we have the formulas

$$\phi_{u'}(x_{u'}) = (\exists y_u^1)(\exists y_u^2)(\exists z_v^1)(\exists z_v^2)(\exists z_w^3)((y_u^1 =_v x_{u'}) \land (y_u^2 =_v \text{"Solomon"}) \land (z_v^2 =_v (y_u^2, z_w^3)) \land (z_v^1 \in z_v^w) \land (y_u^1 \pi_1 z_v^1))$$

and  $\phi_{v'}(x_{v'}) = (\forall y_{u'})(y_{u'} \in x_{v'}).$ 

The intuition is that  $\phi_{u'}(x_{u'})$  says that there is some name  $(y_u^2)$  in I(u) with the value "Solomon" and another name  $(y_u^1)$  with value equal to  $x_{u'}$ . The rest of the formula says that  $y_u^1$  is a parent of  $y_u^2$ .  $\phi_{v'}(x_{v'})$  says that I(v') contains all the names in I(u') in one set. The result of the query is shown in Figure 16.



Fig. 16. Result of  $\mathbf{Q}_2$ .

We now formally define the result of a logical query. We start by looking at queries that add just one node to the schema. We call such queries *simple queries*.

Definition 16. A query **Q** is called a simple query if  $|V_{\mathbf{Q}} - V| = 1$ . We use the notation  $\mathbf{Q}_v$  for a simple query with  $V_{\mathbf{Q}} - V = \{v\}$ .

Let  $\mathbf{Q}_v$  be a simple query on a schema  $\mathbf{S}$  and let  $\mathbf{I}$  be an instance of  $\mathbf{S}$ . The result of  $\mathbf{Q}_v$  on  $\mathbf{I}$  is an extension  $\mathbf{I}_v$  of  $\mathbf{I}$  to  $\mathbf{S}_{\mathbf{Q}_v}$ . In order to define  $\mathbf{I}_v$  we have to define what names  $I_v(v)$  contains and what the values of these names are. We would like  $I_v(v)$  to contain all the objects that satisfy  $\phi_v(x_v)$ . The problem with using this as a definition of  $\mathbf{I}_v$  is that  $\phi_v(x_v)$  is satisfied by names, and since  $I_v$  has not yet been defined, it is meaningless to talk about the objects that satisfy  $\phi_v$ . It might seem that this is really a trivial problem, but suppose that  $\phi_v(x_v)$  included the conjunct  $(\forall y_v)(\forall z_v)(\forall z_v)$ , i.e., I(v) can contain at most one name. If the rest of  $\phi_v$  allowed several possibilities for the value of this name, there would be no way to choose which one should be in the result.

This is not a problem for us because such a formula is not allowed in our query language—all bound variables in our language must refer to database nodes or to nodes that precede v, but cannot refer to v itself. As a result of this restriction, it turns out that although  $\phi_v$  refers to names,  $\phi_v$  is actually a statement about values. We can therefore find the values that satisfy  $\phi_v$  and pick the names arbitrarily.

Definition 17. Let val be a value (i.e., anything that could be the value of  $x_v$ ). We say that val is a candidate value for  $v^3$  if the following holds. Let l be some new name, i.e., a name that does not appear in **I**. Let  $\mathbf{I}_v$  be the extension of **I** to  $\mathbf{S}_{\mathbf{Q}_i}$  with  $I_v(v) = \{l\}$  and  $val_v(l) = val$ . Then  $\models_{\mathbf{I}} \phi_v(l)$ .

By using this arbitrary name, we are able to express the fact that *val* is one of the objects that should be in the result of the query. Note that this definition is where we make use of the fact that the query is acyclic.

We first show that the particular choice of name is unimportant.

LEMMA 5. Let val be a data value and let  $\mathbf{I}_1$  and  $\mathbf{I}_2$  be two extensions of  $\mathbf{I}$  to  $\mathbf{S}_{\mathbf{Q}_n}$  defined by, respectively,  $I_1(v) = \{l_1\}$ ,  $val_1(l_1) = val$ , and  $I_2(v) = \{l_2\}$ ,  $val_2(l_2) = val$ . Then  $\models_{\mathbf{I}_1} \phi_v(l_1) \Leftrightarrow \models_{\mathbf{I}_2} \phi_v(l_2)$ .

 $<sup>{}^{3}</sup>$ This really depends on  ${f Q}$  and  ${f I}$  as well, but they should be clear from the context.

ACM Transactions on Database Systems, Vol 18, No. 3, September 1993

**PROOF.** By definition,  $\phi_v$  has only one free variable of sort v, i.e., the variable  $x_v$ . By inspection we can see that the only atomic formulas that can contain  $x_v$  are  $x_w \pi_i x_v$ ,  $x_w \rho x_v$ ,  $x_w \in x_v$ ,  $x_v =_v d$ ,  $x_v =_v x_w$  and  $x_v =_n x_v$ . The last of these is always true, and it is easy to see that the truth of the others depends only on the value of  $x_v$  and not on its name. The proof is then a straightforward induction.  $\Box$ 

We now define the result of  $\mathbf{S}_{\mathbf{Q}_v}$ . Take all the candidate values for v, pick a new name for each one of them, and put all of these names into  $I_v(v)$ . Notice that the set of candidate values can be infinite in principle, even when I is a finite instance. Queries with the property that the set of candidate values is finite correspond to the safe queries in the relational model. In the next section we look at this issue in more detail.

*Definition* 18. The result of  $\mathbf{Q}_v$  is the extension  $\mathbf{I}_v$  of  $\mathbf{I}$  to  $\mathbf{S}_{\mathbf{Q}_v}$  defined as follows. Let R be the set of all the candidate values for v and let  $\{l_{val} \mid val \in R\}$  be a set of new names, i.e., names that do not appear in  $\mathbf{I}$ .  $I_v(v)$  is defined as the set  $\{l_{val} \mid val \in R\}$  with  $val_v(l_{val}) = val$  for each  $val \in R$ .

We now show that this definition has the desired properties: the result is well defined up to isomorphism relative to **S**, everything in the result satisfies  $\phi_v$ , and we cannot add anything else that satisfies  $\phi_v$  to the result without introducing duplication.

The proof of the following lemma is similar to the proof of Lemma 5.

LEMMA 6. Let  $\mathbf{I}_1$  be an extension of  $\mathbf{I}$  to  $\mathbf{S}_{\mathbf{Q}_r}$ . Let l be an element of  $I_1(v)$ and let  $\mathbf{I}_2$  be the extension of  $\mathbf{I}$  to  $\mathbf{S}_{\mathbf{Q}_r}$  defined by  $\mathbf{I}_2(v) = \{l\}$  and  $val_2(l) = val_1(l)$ . Then  $\models_{\mathbf{I}_r} \phi_v(l) \Leftrightarrow \models_{\mathbf{I}_r} \phi_v(l)$ .

LEMMA 7. (1) Let  $\mathbf{I}_1$  and  $\mathbf{I}_2$  be two results of  $\mathbf{Q}_v$ . Then  $\mathbf{I}_1$  and  $\mathbf{I}_2$  are isomorphic relative to  $\mathbf{S}$ . (2) Let  $\mathbf{I}_v$  be the result of  $\mathbf{S}_{\mathbf{Q}_v}$ . Then for each l in  $I_v(v), \models_{\mathbf{I}_v} \phi_v(l)$ .

PROOF. (1) Let  $l_1$  be an element of  $I_1(v)$ . Since  $val_1(l_1)$  is a candidate value for v, there must be some  $l_2$  in  $I_2(v)$  with  $val_2(l_2) = val_1(l_1)$ . Since, by the definition of the result of the query, both  $I_1(v)$  and  $I_2(v)$  have no duplication, we immediately get a 1-to-1 correspondence between the names of  $I_1(v)$  and  $I_2(v)$ . It is straightforward to show that this correspondence is an isomorphism.

(2) Let l be an arbitrary element of  $I_v(v)$ , and let  $\mathbf{I}^* = \langle I^*, val^* \rangle$  be the extension of  $\mathbf{I}$  to  $\mathbf{S}_{\mathbf{Q}_v}$  defined by  $I^*(v) = \{l\}$  and  $val^*(l) = val(l)$ . By Lemma 6:

$$\vDash_{\mathbf{I}_{v}} \phi_{v}(l) \Leftrightarrow \vDash_{\mathbf{I}^{v}} \phi_{v}(l).$$

Since val(l) is a candidate value for v, we can extend  $\mathbf{I}$  to an instance  $\mathbf{I}^{**}$  of  $\mathbf{S}_{\mathbf{Q}_v}$  by defining  $I^{**}(v) = \{l^{**}\}, val(l^{**}) = val(l)$ . We then have  $\vDash_{\mathbf{I}} \phi_v(l^{**})$ . By Lemma 5,  $\vDash_{\mathbf{I}} \phi_v(l)$ , and therefore  $\vDash_{\mathbf{I}} \phi_v(l)$ .  $\Box$ 

We now define the result of an arbitrary query  $\mathbf{Q}$ . To do this, we first define composition of queries.

ACM Transactions on Database Systems, Vol 18, No 3, September 1993

Definition 19. Let  $\mathbf{Q}_1$  be a query and let  $\mathbf{Q}_2$  be a query on  $\mathbf{S}_{\mathbf{Q}_1}$ .  $\mathbf{Q}_2 \circ \mathbf{Q}_1$  is the query on  $\mathbf{S}$  that we get by composing them in the following way.  $\mathbf{Q}_2 \circ \mathbf{Q}_1$  consists of  $\mathbf{S}_{\mathbf{Q}_2 \circ \mathbf{Q}_1} = \mathbf{S}_{\mathbf{Q}_2}$ ,  $\Phi_{\mathbf{Q}_2 \circ \mathbf{Q}_1} = \Phi_{\mathbf{Q}_1} \cup \Phi_{\mathbf{Q}_2}$ , and

$$\prec_{\mathbf{Q}_2 \circ \mathbf{Q}_1} = \prec_{\mathbf{Q}_1} \cup \prec_{\mathbf{Q}_2} \cup \{(v, w) \mid v \in V_{\mathbf{Q}_1}, w \in V_{\mathbf{Q}_2}\}.$$

LEMMA 8.  $\mathbf{Q}_2 \circ \mathbf{Q}_1$  is a query on **S**.

Let the nodes added by the query  $\mathbf{Q}$  be  $V_{\mathbf{Q}} - V = \{v_1, \ldots, v_n\}$  where  $v_1 \prec \cdots \prec v_n$ . We define a sequence of simple queries  $\mathbf{Q}_{v_1}, \ldots, \mathbf{Q}_{v_n}$  as follows. Each  $\mathbf{Q}_{v_i}$  is a query on the schema of  $\mathbf{Q}_{v_{i-1}}$ .  $\mathbf{Q}_{v_i}$  adds the node  $v_i$  to that schema, and the formula for  $v_i$  is  $\phi_{v_i}$ . It is easy to see that  $\mathbf{Q} = \mathbf{Q}_{v_n} \circ \cdots \circ \mathbf{Q}_{v_1}$ , and we use this to define the result of  $\mathbf{Q}$ .

Definition 20. The result of the query **Q** on **I** is the result of applying the queries  $\mathbf{Q}_{v_1}, \ldots, \mathbf{Q}_{v_n}$  successively to **I**.

LEMMA 9. The result of  $\mathbf{Q}$  is unique up to isomorphism. In other words, different choices of names at each step yield isomorphic results.

Proof. This is a straightforward application of the first part of Lemma 7.  $\square$ 

THEOREM 10. Let  $I_{Q}$  be the result of the query Q on the instance I.

- (1) Let v be a node added by **Q** and let l be an element of  $I_{\mathbf{Q}}(v)$ . Then  $\models_{\mathbf{I}_{\mathbf{Q}}} \phi_{v}(l)$ .
- (2) Let v be a node added by  $\mathbf{Q}$  and let  $l_1$  and  $l_2$  be two different names in I(v). Then  $val(l_1) \neq val(l_2)$ , i.e., there is no duplication in the result.
- (3)  $\mathbf{I}_{\mathbf{Q}}$  is a maximal extension of  $\mathbf{I}$  to  $\mathbf{S}_{\mathbf{Q}}$  that satisfies (1)–(2). This means that there is no extension  $\mathbf{I}_{\mathbf{Q}}^*$  with  $I_{\mathbf{Q}}^*(v) \supseteq I_{\mathbf{Q}}(v)$  for all  $v \in V_{\mathbf{Q}} V$  that satisfies (1)–(2) and such that for at least one v the inclusion is proper.

PROOF.

- (1) Let  $\mathbf{Q}^*$  be the query  $\mathbf{Q}_{v_h} \circ \cdots \circ \mathbf{Q}_{v_1}$  where  $v = v_h$  and let  $\mathbf{I}_{\mathbf{Q}}$  be the result of  $\mathbf{Q}^*$ . By Lemma 7,  $\models_{\mathbf{I}_{\mathbf{Q}}} \phi_v(l)$ . It is easy to see that  $\mathbf{I}_{\mathbf{Q}}$  is an extension of an isomorphic image of  $\mathbf{I}_{\mathbf{Q}^+}$  and that extending  $\mathbf{I}_{\mathbf{Q}^+}$  to  $\mathbf{I}_{\mathbf{Q}}$  does not affect the satisfaction of  $\phi_v$ .
- (2) Obvious.
- (3) Assume that such an  $\mathbf{I}^*$  exists. Let  $v = v_h$  be the first of the nodes  $v_1, \ldots, v_n$  for which  $I^*_{\mathbf{Q}}(v) \neq I_{\mathbf{Q}}(v)$  and let  $\mathbf{Q}^*$  be the query  $\mathbf{Q}_{v_h} \circ \cdots \circ \mathbf{Q}_{v_1}$ . From (1) and (2) it follows immediately that both  $\mathbf{I}_{\mathbf{Q}}$  and  $\mathbf{I}^*_{\mathbf{Q}}$  restricted to  $\mathbf{S}_{\mathbf{Q}^*}$  are results of  $\mathbf{Q}^*$ . Lemma 9 then implies that  $\mathbf{I}^*_{\mathbf{Q}}$  and  $\mathbf{I}_{\mathbf{Q}}$  are isomorphic; a contradiction.  $\Box$

5.2 Safe Queries

In the previous section we observed that the set of candidate values at a query node can be infinite in principle, even over finite instances. For example:

*Example* 10.  $\phi_u(x_{u'})$  is  $(x_{u'} \neq_v \text{``David''})$ . This query is unsafe since the set of candidate data values is  $R = D - \{\text{David}\}$ , an infinite set.

Definition 21. A query  $\mathbf{Q}$  on a schema  $\mathbf{S}$  is *safe on* a finite instance  $\mathbf{I}$  of  $\mathbf{S}$  if the set of candidate data values at each node is finite.  $\mathbf{Q}$  is *safe* if it is safe on every finite instance  $\mathbf{I}$  of  $\mathbf{S}$ .

Let v be a query node. Assume that we have defined the result of  $\mathbf{Q}$  on all the nodes that precede v, and that each of these preceding nodes contains a finite set of names. If  $\mu(v) = \otimes, \otimes$  or  $\odot$ , the set of candidate values for v is a subset of, respectively, the Cartesian product, powerset, or union of the instances of v's children, and therefore is finite. The only case when the set of candidate values may be infinite is when  $\mu(v) = \Box$ , since the set D of data values is infinite.

LEMMA 11. **Q** is safe on **I** iff for every query node of type  $\Box$  the set of candidate data values for v is finite.

The following two examples use the database and query schema of Figure 13 and the database instance of Figure 8.

*Example* 11.  $\phi_u(x_{u'})$  is  $(\exists y_u)(x_{u'} =_v y_u) \lor (x_{u'} =_v \text{"Absalom"})$ . This query is safe, since the set of candidate data values is  $R = \{\text{Jesse, David, Batsheba}, \text{Solomon, Rehoboam, Absalom}\}$ .

Testing safety of relational queries can be reduced to testing safety of LDM queries [31]. As a consequence, testing for safety is, in general, undecidable. It is decidable, however, to test whether a given query  $\mathbf{Q}$  on a schema  $\mathbf{S}$  is safe on a given instance  $\mathbf{I}$ . We now describe the decision procedure.

LEMMA 12. Let  $w_1, \ldots, w_n$  be all of the nodes of the database schema **S** that are of type  $\Box$ , and let  $\{d_1, \ldots, d_k\}$  be the constants that occur in the formula of **Q**. **Q** is safe on **I** iff for each query node v of type  $\Box$ , each candidate value for v is either (a) the value of an element of some  $I(w_i)$  or (b) one of the  $d_i$ 's.

**PROOF.** One direction is obvious—if this condition holds then  $\mathbf{Q}$  is safe on **I**. We prove the converse by induction on the query nodes  $V_{\mathbf{Q}} - V =$  $\{v_1, \ldots, v_n\}$  where  $v_1 \prec \cdots \prec v_n$ . Let  $v = v_i$  be a query node of type  $\Box$ , and assume that the lemma holds for the nodes that precede  $v_i$  and that  $\mathbf{Q}$  is safe on **I**. Let  $\mathbf{I}_{i-1}$  be the result of  $\mathbf{Q}_{v_i}$ .

on I. Let  $\mathbf{I}_{\iota-1}$  be the result of  $\mathbf{Q}_{v_{\iota-1}}$ . Since  $\mathbf{Q}_v$  is safe on I, the set of candidate data values for v is a finite set R. We have to show that

$$R \subseteq \{d_1, \ldots, d_k\} \cup \bigcup_{\substack{\mu(w) = \ \square \\ w \in V}} I(w).$$

Call the right hand of this equation S. If the lemma is false, then there is some element val in R - S. By the induction hypothesis,

$$S = \{d_1, \dots, d_k\} \cup \bigcup_{\substack{\mu(w) = \square \\ w \in V \text{ or} \\ w \in V_{\mathbf{Q}}, w \prec v}} I_{i-1}(w).$$

Since *val* is a candidate data value for v, if we extend  $I_{i-1}$  to an instance  $\mathbf{I}_{i}^{1}$  of  $\mathbf{S}_{\mathbf{Q}_{v}}$  by defining  $I_{i}^{1}(v) = \{l\}$  and  $val_{1}(l) = val$ , we have  $\models_{\mathbf{I}_{i}^{1}} \phi_{v}(l)$ . Let val' be an arbitrary element of D - S, and extend  $\mathbf{I}_{i-1}$  to an instance  $\mathbf{I}_{i}^{2}$  of  $\mathbf{S}_{\mathbf{Q}_{v}}$  by defining  $I_{i}^{2}(v) = \{l\}$  and  $val_{2}(l) = val'$ . Since *val* and *val'* do not appear either in the database, in preceding query nodes, or in the query formulas, an induction on the size  $\phi_{v}(x_{v})$  shows that

$$\vDash_{\mathbf{I}^1} \phi_v(l) \Leftrightarrow \vDash_{\mathbf{I}^2} \phi_v(l).$$

The key point in the induction is that  $x_v$  can occur in  $\phi_v(x_v)$  only in

- (1) atomic formulas of the form  $x_v =_v d_j$  and  $x_v =_v y_w$ , where w is a node of type  $\Box$  that is either in V or is one of the nodes  $v_1, \ldots, v_{i-1}$ —such a formula is false whenever the data value of  $x_v$  is not in S;
- (2)  $x_v$  are  $x_v =_n x_v$  and  $x_v =_v x_v$ —these formulas are always true.

We have shown that all the elements of the infinite set D - S are candidate values, contradicting the assumption that **Q** is safe on **I**.  $\Box$ 

The technique of the proof gives us an effective procedure for determining whether a simple query  $\mathbf{Q}_v$  is safe on a finite instance I. Take some data value  $d_0 \in D$  that does not occur anywhere in the database or in the query formulas. Test if  $d_0$  is a candidate value (it is not difficult to see that this can be done effectively). In a similar way to the proof of this lemma, we can show that  $\mathbf{Q}_v$  is safe on I iff  $d_0$  is not a candidate value for v. Intuitively, if some such  $d_0$  is in the result, the result is infinite since  $d_0$  cannot be distinguished from any other such data value. This proves the following theorem.

THEOREM 13. Let  $\mathbf{Q}$  be a query on  $\mathbf{S}$  and let  $\mathbf{I}$  be a finite instance of  $\mathbf{S}$ . There is a decision procedure to test whether  $\mathbf{Q}$  is safe on  $\mathbf{I}$ . If  $\mathbf{Q}$  is safe on  $\mathbf{I}$ , then the result can be computed effectively.

# 6. THE ALGEBRAIC QUERY LANGUAGE

## 6.1 The Algebraic Operators

In this section we define a set of algebraic operators. We then show that any safe logical query is equivalent to some sequence of algebraic operations, and, conversely, each sequence of algebraic operations is equivalent to a safe logical query.

Since a logical query adds some nodes to the database schema and leaves the instance of the database schema unchanged, each algebraic operator must do the same. Therefore, a selection operator, for example, should not delete tuples that do not satisfy the selection condition, but should rather create a copy of the database node. That copy should contain only those

Fig. 17. The algebraic operation  $w \leftarrow \Box(v)$ .



tuples that satisfy the condition. In fact, this copying of tuples is what is really done in the relational model—a query does not throw away those tuples in the database that do not meet a selection condition, but copies those tuples that do. This issue is not addressed explicitly in relational database theory, because the theory does not deal with what happens to temporary relations that are created while computing the result of a query. As we shall see in the next section, it is necessary to delete nodes in certain circumstances. We would still like to make the deletions explicit, rather than hide them in other algebraic operations.

In this section **S** is a database schema with instance **I**. The algebra consists of operations of the form  $w \leftarrow \alpha(v_1, \ldots, v_n)$ . Here  $\alpha$  is the name of the operator, and its arguments  $v_1, \ldots, v_n$  are nodes in the schema **S**.  $\alpha$  adds the node w to the schema, and extends **I** to the new schema. We define each operator as a simple logical query, by giving

- (1) the types of its arguments;
- (2) the type of w and the list of its children;
- (3) an LDM formula  $\phi_w(x_w)$  that specifies the contents of I(w).

6.1.1 Operators that Copy and Combine Existing Nodes.

- (1)  $w \leftarrow \Box(v)$  creates a copy of the node v, as is shown in Figure 17. (In all these figures the schema **S** is shown in the box on the right, and the node that is created by the operation is on the left.) For each distinct data value in I(v), I(w) contains exactly one name with this data value. Note that duplication in I(v) is eliminated in I(w).
  - (a) v is a node of **S** that has type  $\Box$ ;
  - (b) w is of type  $\Box$ ;
  - (c)  $\phi_w(x_w)$  is  $(\exists y_v)(x_w =_v y_v)$ .
- (2)  $w \leftarrow \Box$  (d) creates a node of type  $\Box$  that contains just the data value d.
  - (a) d is a data value in the data domain D;
  - (b) w is of type  $\Box$ ;
  - (c)  $\phi_w(x_w)$  is  $x_w =_v d$ .
- (3)  $w \leftarrow \circledast(v)$  creates a node that contains all finite subsets of I(v) (see Figure 18).
  - (a) v is any node in the schema **S**;
  - (b) w is of type ( $\circledast$ , v);
  - (c)  $\phi_w(x_w)$  is T (i.e., always true).



- (4) w ← ⊗(v<sub>1</sub>,...,v<sub>n</sub>) creates a node that contains the Cartesian product I(v<sub>1</sub>) × … × I(v<sub>n</sub>) (see Figure 19).
  - (a)  $v_1, \ldots, v_n$  are any *n* nodes in the schema **S**;
  - (b) w is of type  $(\otimes, n, v_1, \dots, v_n)$ ;
  - (c)  $\phi_w(x_w)$  is T.
- (5)  $w \leftarrow \bigcirc (v_1, \ldots, v_n)$  creates a node that contains the disjoint union  $I(v_1) \cup \cdots \cup I(v_n)$  (see Figure 20).
  - (a)  $v_1, \ldots, v_n$  are *n* distinct nodes of the schema **S**;
  - (b) w is of type  $(\odot, n, v_1, \ldots, v_n)$ ;
  - (c)  $\phi_w(x_w)$  is T.

6.1.2 Selection Operators. The LDM algebra has two selection operators.

- (1) The operation  $w \leftarrow \sigma_{\iota\theta_j}(v)$  is similar to the selection operation in the relational algebra. This operator selects those tuples in v whose *i*th and *j*th components are related by  $\theta$  (see Figure 21).
  - (a) v is a node of **S** of type  $(\otimes, n, v_1, \dots, v_n)$  and  $i\theta j$  is one of the relations  $i \in j, i\pi_t j, i\rho j, i =_n j$  and  $i =_v j$ ;
  - (b) w is of type  $(\otimes, n, v_1, \dots, v_n)$ ;
  - (c)  $\phi_w(x_w)$  is

 $(\exists x_v)(\exists y_{v_i})(\exists y_{v_i})(y_{v_i}, \pi_{v_i}x_v \wedge y_{v_j}\pi_{v_j}x_v \wedge y_{v_i}\theta y_{v_i} \wedge x_v =_v x_w).$ 





Alternatively, the selection condition may be of the form  $i =_v d$  where d is a data value in D. Then  $\phi_w(x_w)$  is  $(\exists x_v)(\exists y_{v_i})(y_{v_i}\pi_{v_i}x_v \land y_{v_i} =_v d \land x_v =_v x_w)$ .

- (2)  $w \leftarrow \sigma_{in}(u, v)$ . Here u is a child of v, and w contains those objects of I(u) that actually appear in I(v) (see Figure 22). For example, if v is of type  $\circledast$ , each element of I(v) becomes a set of elements from I(u). The result of  $\sigma_{in}(u, v)$  selects from I(v) those elements that are in at least one of these sets.
  - (a) u and v are nodes of **S** and u is a child of v;
  - (b) w is of the same type as u and has the same children;
  - (c)  $\phi_w(x_w)$  depends on the type of v. Note that v cannot be of type  $\Box$  since it has a child u.
    - (i) If v is of type  $\otimes$  with u as its *i*th child then  $\phi_w(x_w)$  is

$$(\exists x_u)(\exists x_v)(x_u = x_w \land x_u \pi_i x_v).$$

(Strictly speaking,  $\sigma_{\rm in}(u, v)$  is under-specified here, in case there are multiple edges from v to u, since we have to specify the edge to which we refer. In this case we use the notation  $\sigma_{\rm in}(u, v, i)$  to mean: use the *i*th edge with tail v.)

(ii) If v is of type  $\bigcirc$ , then  $\phi_w(x_w)$  is

$$(\exists x_u)(\exists x_v)(x_u =_v x_w \land x_u \rho x_v).$$

(iii) If v is of type  $\circledast$ , then  $\phi_w(x_w)$  is

$$(\exists x_u)(\exists x_v)(x_u = x_w \land x_u \in x_v).$$

6.1.3 Union, Difference, and Projection.

- (1) The union operator is similar to the relational union. The syntax we use is  $w \leftarrow \cup (v_1, \ldots, v_n)$  (see Figure 23).
  - (a)  $v_1, \ldots, v_n$  are *n* nodes of **S** that are of the same type and have the same children;
  - (b) w has the same type and the same children as the  $v_i$ 's;

(c)  $\phi_w(x_w)$  is  $(\exists x_{v_1})(x_{v_1} = x_w) \lor \cdots \lor (\exists x_{v_n})(x_{v_n} = x_w)$ .



- (2) For difference we use infix notation, i.e., we write  $w \leftarrow v_1 v_2$  rather than  $-(v_1, v_2)$ .
  - (a)  $v_1$  and  $v_2$  are nodes of **S** that are of the same type and have the same children;
  - (b) w has the same type and the same children as  $v_1$  and  $v_2$ ;
  - (c)  $\phi_w(x_w)$  is  $(\exists x_{v_1})(x_{v_1} = x_w) \land (\forall x_{v_2})(x_{v_2} \neq x_w)$ .
- (3) The projection operation is similar to projection in the relational algebra. The syntax we use is  $w \leftarrow \prod_A(v)$ , where A is an ordered multiset of edges with tail v.
  - (a) v is a node of **S** of type  $(\otimes, n, v_1, \dots, v_n)$  and A is an ordered multiset of edges with tail v;
  - (b) let  $A = \{e_1, \ldots, e_k\}$  where  $e_j$  is the edge  $(v, v_j)$ . Then w is of type  $(\otimes, k, v_{\iota_1}, \ldots, v_{\iota_k});$
  - (c)  $\phi_w(x_w)$  is

$$(\exists x_v)(\exists x_{v_1})\cdots(\exists x_{v_n})(x_{v_{i_1}}\pi_1x_w\wedge\cdots\wedge(x_{v_{i_k}}\pi_kx_w)$$
$$\wedge x_v =_v (x_{v_1},\ldots,x_{v_k})).$$

When it does not cause any ambiguity, we use a set A of nodes rather than of edges, as in Figure 24.

## 6.2 Equivalence of the Logical and Algebraic Query Languages

We now use these algebraic operators to define an algebraic query language. An algebraic query is a sequence  $\{\alpha_1, \ldots, \alpha_n\}$ , where each  $\alpha_i$  is an algebraic operator on the result of  $\alpha_{i-1}$ . We would like to be able to show that this query language is equivalent to the logical query language. In other words,

ACM Transactions on Database Systems, Vol 18, No. 3, September 1993.



Fig. 24. The algebraic operation  $w \leftarrow \prod_{\{v_1, v_2\}} (v)$ .

for each logical query on a schema S, there should exist a sequence of algebraic operations, and vice versa, with the property that the schemas created by these two queries are identical, and, for every database instance I, the results are isomorphic relative to S. Unfortunately, as the next example shows, this is not quite true.

*Example* 12. Let **S** consist of a node u of type  $\Box$  and let **Q** be the logical query that adds a node v of type  $\Box$  to **S**. Let  $d_1$  and  $d_2$  be two data values, and let  $\phi_v(x_v)$  be  $(x_v =_v d_1 \lor x_v =_v d_2)$ . The candidate values for v are then  $\{d_1, d_2\}$ . There is no algebraic query equivalent to **Q**. If there was such a query, it would consist of one algebraic operation alone, since each operator adds a new node to the schema. By inspection, we can see that no single algebraic operator is equivalent to **Q**.

How can we modify the definition to get an equivalent query? If  $\mathbf{Q}_A$  is the algebraic query that consists of the operations  $w_1 \leftarrow \Box(d_1), w_2 \leftarrow \Box(d_2)$ , and  $v \leftarrow \cup(w_1, w_2)$ , it is clear that the instance of v is what we are after. If we were then to restrict the result of the query to the schema that consists of the nodes u and v, we get the instance we want. We have essentially used the two nodes  $w_1$  and  $w_2$  for temporary storage while computing the result of the query. In fact, the same thing occurs in the relational model, since temporary relations are used there for subexpressions and then deleted at the end. It is therefore reasonable to expect the same thing to happen in the logical data model.

To be able to use temporary nodes, we extend the algebraic query language by adding a "delete" operator. This operator deletes a node from the schema and restricts the instance of the original schema to the new schema. We have to make sure that we never delete a node that is the child of some other node, since in that case the result would not be a legal schema. The operator that deletes the node v will be written D(v).

Definition 22. Let **S** be an LDM schema with instance **I**. The algebraic operator D(v) is legal when v is a node with no parent. The result of D(v) is the schema **S**' that consists of deleting v from **S**, together with the instance that we get by restricting **I** to **S**'.

In the algebraic query language we must take care not to delete database nodes, i.e., we must only allow user to delete nodes that have been constructed by the user query. We call the language with the deletion operator the *extended algebraic query language*.

Definition 23. Let **S** be an LDM schema. An extended algebraic query on **S** is a sequence  $\mathbf{Q}_A = (\alpha_1, \dots, \alpha_n)$  where each  $\alpha_i$  is either

- an operation of the form w<sub>i</sub> ← β<sub>i</sub>(v<sup>1</sup><sub>i</sub>,...,v<sup>i<sub>i</sub></sup>), where β<sub>i</sub> is an algebraic operator other than the deletion operator and v<sup>1</sup><sub>i</sub>,...,v<sup>i<sub>i</sub></sup> are either node of S or are nodes that were created by some previous β<sub>j</sub> and have not been deleted;
- (2) the operator  $D(v_i)$ , where  $v_i$  is a node that was created by a previous algebraic operator in the sequence  $\beta_1, \ldots, \beta_{i-1}$  and has not yet been deleted.

Definition 24. Let  $\mathbf{Q}_A$  be an extended algebraic query on  $\mathbf{S}$ , and let  $\mathbf{Q}_B$  be an extended algebraic query on the result of  $\mathbf{Q}_A$ . The query  $\mathbf{Q}_B \circ \mathbf{Q}_A$  is the composition of  $\mathbf{Q}_A$  and  $\mathbf{Q}_B$ , formed simply by concatenating the lists of algebraic operators.

Obviously, the delete operator itself is not equivalent to any logical query, since every logical query adds nodes to the schema. This by itself does not necessarily mean that we cannot find a logical query equivalent to any extended algebraic query. After all, an extended algebraic query does not delete database nodes, since the only nodes that are deleted are those that were constructed by previous algebraic operations. It might still be the case (as in Example 12) that there is an equivalent logical query that does not use temporary nodes. Nevertheless, in [31] it is shown that the extended algebraic query language is strictly more powerful than the logical query language. In order to get an equivalent query language, we modify the logical language to include temporary nodes as well.

Definition 25. Let **S** be an LDM schema. An extended logical query on **S** is a tuple  $\mathbf{Q} = \langle S_{\mathbf{Q}}, \Phi_{\mathbf{Q}}, \prec_{\mathbf{Q}}, D_{\mathbf{Q}} \rangle$  where

- (1)  $\langle S_{\mathbf{Q}}, \Phi_{\mathbf{Q}}, \prec_{\mathbf{Q}} \rangle$  is a logical query on **S**;
- (2)  $D_{\mathbf{Q}}$  is the set of temporary nodes used in the query.  $D_{\mathbf{Q}}$  is a subset of the query nodes  $V_{\mathbf{Q}} V$  that we can delete and still get an LDM schema. In other words, there is no edge with tail outside  $D_{\mathbf{Q}}$  and head in  $D_{\mathbf{Q}}$ , i.e., if  $(e_1, e_2) \in E_{\mathbf{Q}}$  and  $e_2 \in D_{\mathbf{Q}}$ , then  $e_1 \in D_{\mathbf{Q}}$ .

*Definition* 26. Let **Q** be the extended logical query  $\langle S_{\mathbf{Q}}, \Phi_{\mathbf{Q}}, \prec_{\mathbf{Q}}, D_{\mathbf{Q}} \rangle$ , and let **I** be an instance of **S**. The result of this query consists of

- (1) the schema  $\mathbf{S}^*_{\mathbf{Q}}$  consisting of
  - (a) the nodes in  $V_{\mathbf{Q}} D_{\mathbf{Q}}$ ;
  - (b) the relevant edges, i.e., all those edges of  $S_Q$  whose head and tail are both in  $V_Q D_Q$ ;
  - (c) the restriction of the type assignment  $\mu$  to  $V_{\mathbf{Q}} D_{\mathbf{Q}}$ .
- (2) The result of **Q** on **I** is defined as follows. Let  $\mathbf{I}_{\mathbf{Q}}$  be the result of  $\langle \mathbf{S}_{\mathbf{Q}}, \Phi_{\mathbf{Q}}, \prec_{\mathbf{Q}} \rangle$  on **I**. The result of **Q** on **I** is then the restriction of  $\mathbf{I}_{\mathbf{Q}}$  to  $\mathbf{S}_{\mathbf{Q}}^*$ .

We now show that every extended algebraic query is equivalent to some extended logical query.

ACM Transactions on Database Systems, Vol. 18, No. 3, September 1993.

LEMMA 14. Let  $\mathbf{Q}_A = (\alpha_1, \dots, \alpha_n)$  be an extended algebraic query on  $\mathbf{S}$ . There exists a safe extended logical query  $\mathbf{Q}_L$  on  $\mathbf{S}$  such that for every instance  $\mathbf{I}$  of  $\mathbf{S}$ , the results of  $\mathbf{Q}_A$  and  $\mathbf{Q}_L$  on  $\mathbf{I}$  are isomorphic relative to  $\mathbf{S}$ .

PROOF. The schema of  $\mathbf{Q}_L$  consists of all those nodes that are created by the operations in the query  $\mathbf{Q}_A$ . The set of temporary nodes  $D_{\mathbf{Q}_L}$  is the set of nodes deleted in  $\mathbf{Q}_A$ , i.e.,  $\{v_i \mid \text{the operator } \alpha_i \text{ is } D(v_i)\}$ . Since we are only allowed to delete nodes that are not in  $\mathbf{S}$  and that have no parent, it is easy to see that there is no edge whose head is in  $D_{\mathbf{Q}_L}$  and whose tail is not. Each  $\alpha_i$ that is not a delete operator must be of the form  $w_j \leftarrow \beta_j(w_i^1, \ldots, w_i^{i_j})$ . We define an order on the nodes of  $V_{\mathbf{Q}} - V$  as follows:  $w_i \prec w_j$  whenever i < j.  $\phi_{w_i}(x_{w_i})$  is the formula that was used to define the operator  $\beta_j$ . It is easy to verify that the results of  $\mathbf{Q}_A$  and  $\mathbf{Q}_L$  on any instance  $\mathbf{I}$  are indeed isomorphic.  $\Box$ 

To show the converse, let  $\mathbf{Q}_L$  be a logical query on **S**. Let **I** be a fixed instance of **S** such that  $\mathbf{Q}_L$  is safe on **I**. The definition of  $\mathbf{Q}_A$  will not depend on **I**, but the results of  $\mathbf{Q}_A$  and  $\mathbf{Q}_L$  will only be isomorphic on those instances of **S** on which  $\mathbf{Q}_L$  is safe. We keep **I** fixed so that we will be able to prove various lemmas about the results of the queries as we go along.

We first look at the case when  $\mathbf{Q}_L$  is a simple query  $\mathbf{Q}_w$ . We start by creating a node  $w_{\text{dom}}$ , that contains the "domain" of w, i.e., all those objects that might be candidate values for w, if we were to ignore everything except the type of w and the fact that  $\mathbf{Q}_L$  is safe on **I**. We define  $w_{\text{dom}}$  as follows.

(1) If w is of type  $\Box$ , let  $v_1, \ldots, v_t$  be all the nodes in **S** that are of type  $\Box$  and let  $d_1, \ldots, d_k$  be the constants that occur in  $\phi_w(x_w)$ . Define  $w_{\text{dom}}$  by the algebraic query:

$$s_{1} \leftarrow \Box(v_{1})$$

$$\vdots$$

$$s_{t} \leftarrow \Box(v_{t})$$

$$s_{t+1} \leftarrow \Box(d_{1})$$

$$\vdots$$

$$s_{t+k} \leftarrow \Box(d_{k})$$

$$w_{\text{dom}} \leftarrow \cup(s_{1}, \dots, s_{t+k})$$

$$D(s_{1})$$

$$\vdots$$

$$D(s_{t+k})$$

- (2) If  $\mu(w) = (\otimes, k, v_1, \dots, v_k)$  define  $w_{\text{dom}}$  by  $w_{\text{dom}} \leftarrow \otimes (v_1, \dots, v_k)$ .
- (3) If  $\mu(w) = (\circledast, v)$  define  $w_{\text{dom}}$  by  $w_{\text{dom}} \leftarrow \circledast(v)$ .
- (4) If  $\mu(w) = (\odot, k, v_1, \dots, v_k)$  define  $w_{\text{dom}}$  by  $w_{\text{dom}} \leftarrow \odot(v_1, \dots, v_k)$ .

Call this algebraic query  $\mathbf{Q}_{dom}$ .

Lemma 15.

- (1) The schema created by  $\mathbf{Q}_{dom}$  is equal to the schema of  $\mathbf{S}$  together with a node  $w_{dom}$  of the same type and with the same children as the node w in the original logical query  $\mathbf{Q}_{w}$ .
- (2) Let  $\mathbf{I}_{dom}$  be the result of  $\mathbf{Q}_{dom}$  on  $\mathbf{I}$  and let  $\mathbf{I}_w$  be the result of  $\mathbf{Q}_w$  on  $\mathbf{I}$ . If val is a value of an object in  $I_w(w)$ , then val is also a value of an object in  $\mathbf{I}_{dom}(w_{dom})$ .

**PROOF.** If w is of type  $\otimes$ ,  $\odot$ , or  $\circledast$ , the lemma is obvious. If w is of type  $\Box$ , the first part follows from the fact that all the nodes except  $w_{\text{dom}}$  that are created by  $\mathbf{Q}_{\text{dom}}$  are also deleted by  $\mathbf{Q}_{\text{dom}}$ . The second part is an immediate consequence of Lemma 12 and the definition of  $\mathbf{Q}_{\text{dom}}$ .  $\Box$ 

We may assume (if necessary by renaming some bound variables) that all the bound variables in the formula  $\phi_w(x_w)$  are distinct. Let these variables be  $x_{w_1}^1, \ldots, x_{w_k}^k$ . The algebraic query  $\mathbf{Q}_{\text{prod}}$  on the result of  $\mathbf{Q}_{\text{dom}}$  consists of the operation

$$w_{\text{prod}} \leftarrow \otimes (w_1, \dots, w_k, w_{\text{dom}}).$$

For the purpose of defining  $\mathbf{Q}_A$ , we label the edges with tail  $w_{\text{prod}}$ , as follows. The *i*th edge with tail  $w_{\text{prod}}$  will be labeled  $x_{w_i}^i$ . These labels are used only to define the algebraic query, and are not themselves part of the query. Their purpose is to tell us which bound variable the edge corresponds to.

In certain cases, when we create a new node using some algebraic operation, the outgoing edges from the new node will inherit the labels of the corresponding edges whose head is one of the arguments of the operator. We only use this inheritance in cases when it is unambiguous, i.e., in cases when all the arguments have the same labeling. The operations for which labels will be inherited are  $\sigma_{\iota\theta j}$ , difference and union. When we use the projection operation the new edges will also inherit the labeling of the corresponding edges whose head is the argument of the projection.

Arrange all the well-formed subformulas of  $\phi_w(x_w)$  in a list  $\psi_1, \ldots, \psi_m$ , where  $\psi_m = \phi_w(x_w)$ , and  $\psi_i$  precedes  $\psi_j$  whenever it is a subformula of  $\psi_j$ . For each such subformula, we define an extended algebraic query  $\mathbf{Q}_{\psi_i}$  on the result of  $\mathbf{Q}_{\psi_{i-1}}$ .  $\mathbf{Q}_{\psi_1}$  is a query on the result of  $\mathbf{Q}_{\text{prod}}$ . The node  $w_{\psi_i}$  is of type  $(\otimes, k + 1, w_{j_1}, \ldots, w_{j_k}, w_{\text{dom}})$ , and contains, intuitively, those tuples  $(l_1, \ldots, l_k, l_d)$  for which  $\vDash_{\mathbf{I}_{\psi_i}} \psi_i(l_1, \ldots, l_k, l_d)$ . The edges with tail  $w_{\psi}$  will be labelled with the variables that might be free in  $\psi$ —i.e., those that have not yet been bound by  $\psi$ . The definition of  $\mathbf{Q}_{\psi_i}$  is as follows.

- (1)  $\psi_i$  is  $x_{w_a}^a \theta x_{w_b}^b$ .  $\mathbf{Q}_{\psi_i}$  is  $w_{\psi_i} \leftarrow \sigma_{a\theta b}(w_{\text{prod}})$ .
- (2)  $\psi_i$  is  $x_{w_a}^a \theta x_w$ .  $\mathbf{Q}_{\psi_i}$  is  $w_{\psi_i} \leftarrow \sigma_{a \theta k+1}(w_{\text{prod}})$ .
- (3)  $\psi_i$  is  $x_w \theta x_w$ .  $\mathbf{Q}_{\psi_i}$  is  $w_{\psi_i} \leftarrow \sigma_{k+1\theta k+1}(w_{\text{prod}})$ .
- (4)  $\psi_i$  is  $x_{w_a}^a =_v d$ .  $\mathbf{Q}_{\psi_i}$  is  $w_{\psi_i} \leftarrow \sigma_{a=,d} (w_{\text{prod}})$ .
- (5)  $\psi_i$  is  $x_{\omega}^{a} =_{v} d$ .  $\mathbf{Q}_{\psi_i}$  is  $w_{\psi_i} \leftarrow \sigma_{(k+1)=vd} (w_{\text{prod}})$ .
- (6)  $\psi_i$  is  $\psi_{j_1} \vee \psi_{j_2}$ . Let  $A_1$  be the (ordered multi)set of edges with tail  $w_{\psi_{j_1}}$  that have the same label as some edge with tail  $\psi_{\psi_{j_2}}$ . Let  $A_2$  be the

corresponding set of edges with tail  $w_{\psi_{j_2}}$ .  $\mathbf{Q}_{\psi_j}$  is the following extended algebraic query

$$s_1 \leftarrow \Pi_{A_1}(w_{\psi_{j_1}})$$

$$s_2 \leftarrow \Pi_{A_2}(w_{\psi_{j_2}})$$

$$w_{\psi_i} \leftarrow \cup (s_1, s_2)$$

$$D(s_1)$$

$$D(s_2)$$

 $(s_1 \text{ and } s_2 \text{ are different temporary nodes from those used above, and from similarly named nodes used below.) Note that the way we defined <math>A_1$  and  $A_2$  guarantees that there is no ambiguity in labeling the edges of the result, at least as long as the labels of the edges in  $A_1$  and  $A_2$  are in the same order. We show later that this is indeed the case.

(7)  $\psi_i$  is  $\neg \psi_j$ . Let A be the (ordered multi)set of edges with tail  $w_{\text{dom}}$  that have the same label as some edge with tail  $w_{\psi_j}$ .  $\mathbf{Q}_{\psi_i}$  is the following extended algebraic query

$$s_1 \leftarrow \Pi_A(w_{\psi_{\text{dom}}})$$
$$w_{\psi_i} \leftarrow s_1 - w_{\psi_j}$$
$$D(s_1).$$

As in the previous case, we show that we can label the edges with tail  $w_{\psi_i}$  without any ambiguity.

(8)  $\psi_i$  is  $(\exists x_{w_a}^a)(\psi_j)$ . Let A be the (ordered multi)set of all edges with tail  $w_{\psi_j}$  except for the edge labeled  $x_{w_a}^a$ . We show below that there must be exactly one edge with this label.  $\mathbf{Q}_{\psi_i}$  then consists of the algebraic operation  $w_{\psi_i} \leftarrow \Pi_A(w_{\psi_i})$ .

LEMMA 16. Let  $\psi = \psi_i$  be one of these well-formed subformulas of  $\phi_w(x_w)$ . Let  $x_{w_{a_1}}^{a_1}, \ldots, x_{w_{a_i}}^{a_i}$  be those variables in the above list that are not bound in  $\psi_i$ . Note that some of the  $x_{w_{a_i}}^{a_i}$ 's may not actually occur in  $\psi_i$ . Then  $w_{\psi_i}$  is of type  $(\otimes, j + 1, w_{a_1}, \ldots, w_{a_i}, w_{dom})$ . Furthermore, the th edge with tail  $w_{\psi_i}$  has head  $w_{a_i}$  and is labeled with the variable  $x_{w_{a_i}}^{a_i}$ . As a consequence of this, all the labelings of edges are in the same order, and the assumptions that we made on the labelings when we defined the  $w_{\psi_i}$ 's hold.

PROOF. The proof is a fairly straightforward induction using the definition of  $w_{\psi_i}$ . The tricky case is when  $\psi_i$  is  $\psi_{j_1} \vee \psi_{j_2}$ . Then the children of  $w_{\psi_{j_1}}$  correspond to the bound variables of  $\phi_w$  that are not bound in  $\psi_{j_1}$  and the children of  $w_{\psi_{j_2}}$  to the bound variables of  $\phi_w$  not bound in  $\psi_{j_2}$ . Since a variable is not bound in  $\psi_{j_1} \vee \psi_{j_2}$  iff it is not bound in  $\psi_{j_1}$  and it is not bound in  $\psi_{j_2}$ , we see that the result does hold in this case.<sup>4</sup>

<sup>&</sup>lt;sup>4</sup>Note that "not bound" is not the same as "free." A variable that is not bound in one of these formulas is either free in this formula, or does not appear in it at all.

ACM Transactions on Database Systems, Vol. 18, No. 3, September 1993

LEMMA 17. Let  $w_{\psi_i}$  be of type  $(\otimes, j, w_{j_1}, \ldots, w_{j_k}, w_{dom})$ . Let  $\mathbf{I}_{\psi_i}$  be the result of  $\mathbf{Q}_{\psi_i}$  on  $\mathbf{I}$ , let  $l_d$  be a member of  $I_{\psi}(w_{dom})$  and let  $l_t$  be a member of  $I_{\psi_i}(w_{j_i})$ for  $t = 1, \ldots, k$ . Then there exists an l in  $I_{\psi_i}(w_{\psi_i})$  with  $val(l) = (l_1, \ldots, l_k, l_d)$  if and only if  $\models_{\mathbf{I}_{\psi_i}} \psi_i(l_1, \ldots, l_k, l_d)$ . Intuitively,  $(l_1, \ldots, l_k, l_d)$  is a candidate value iff it satisfies  $\psi_i$ .

**PROOF.** A straightforward induction on the structure of  $\psi_i$ .

The extended algebraic query  $\mathbf{Q}_{\text{final}}$  on the result of  $\mathbf{Q}_{\phi_w}$  consists of

$$w_A \leftarrow \sigma_{in}(w_{dom}, w_{\phi})$$
$$D(w_{\phi})$$
$$\vdots$$
$$D(w_{\psi_1})$$
$$D(w_{prod})$$
$$D(w_{dom}).$$

We finally define the algebraic query  $\mathbf{Q}_A$  as

$$\mathbf{Q}_{\text{final}} \circ \mathbf{Q}_{\phi} \circ \mathbf{Q}_{\psi_{m-1}} \circ \mathbf{Q}_{\psi_1} \circ \mathbf{Q}_{\text{prod}} \circ \mathbf{Q}_{\text{dom}}$$

LEMMA 18. Let  $\mathbf{I}_1$  be the result of  $\mathbf{Q}_w$  on  $\mathbf{I}$  and let  $\mathbf{I}_2$  be the result of the algebraic query  $\mathbf{Q}_A$  on  $\mathbf{I}$ . Then  $\mathbf{I}_1$  and  $\mathbf{I}_2$  are isomorphic relative to  $\mathbf{S}$ .

**PROOF.** First note that the schemas are equal. The only node created but not deleted by  $\mathbf{Q}_A$  is the node  $w_A$ . This node is similar to the node  $w_{\text{dom}}$ , and hence to w.

We have to show that the instances of  $w_A$  and w are isomorphic, i.e., that at the point in evaluating the queries where we compute the instances of these nodes, they have the same candidate values. We assume that we are at the point in the evaluation of  $\mathbf{Q}_A$  just before the final round of deletions.

Let val be a candidate data value for w. Extend I to an instance  $\mathbf{I}_w$  of  $\mathbf{S}_{\mathbf{Q}_w}$  by defining  $I_w(w) = \{l\}$  and val(l) = val. Then  $\models_{\mathbf{I}_w} \phi_w(l)$ . Let  $\mathbf{I}_{\phi_w}$  be the result of  $\mathbf{Q}_{\phi_w}$  on I. By the second clause of Lemma 15, val is a candidate data value for  $w_{\text{dom}}$ , and so for some  $l_d$  in  $I_{\phi_w}(w_{\text{dom}})$ ,  $val(l_d) = val$ . By Lemma 6,  $\models_{\mathbf{I}_w} \phi_w(l)$  implies  $\models_{\mathbf{I}_{\phi_w}} \phi_w(l_d)$ , and therefore, by Lemma 17, for some  $l_\phi$  in  $I_{\phi_w}(w_{\text{dom}})$ ,  $val(l_\phi) = val$  is a candidate data value for  $w_A$ .

For the converse, suppose that val is a candidate data value for  $w_A$ . Let  $\mathbf{I}_{\phi_w}$  be the result of  $\mathbf{Q}_{\phi_w}$  on I. Since val is a candidate data value for  $w_A$ , for some  $l_{\phi}$  in  $I_{\phi_w}(w_{\phi_w})$  and some  $l_d$  in  $I_{\phi_w}(w_{\mathrm{dom}})$ ,  $val(l_{\phi}) = (l_d)$  and  $val(l_d) = val$ . Since  $l_{\phi}$  is in  $I_{\phi_w}(w_{\phi_w})$ , Lemma 15 implies that  $\models_{\mathbf{I}_{\phi_w}} \phi_w(l_d)$ . Restrict  $\mathbf{I}_{\phi_w}$  to an instance  $\mathbf{I}_{\mathrm{dom}}$  of the schema of  $\mathbf{Q}_{\mathrm{dom}}$ . Then  $\models_{\mathbf{I}_{\mathrm{dom}}} \phi_w(l_d)$ , and so by Lemma 6,  $val(l_d) = val$  is a candidate data value for w.  $\Box$ 

Since an arbitrary logical query can be viewed as a sequence of simple queries, we can easily extend the above construction to arbitrary queries by concatenating the algebraic queries for the individual simple queries. If we

ACM Transactions on Database Systems, Vol 18, No 3, September 1993.

have an extended logical query, we have to add deletion operations to delete those nodes in the delete set of the query. This yields the following theorem.

THEOREM 19. The extended algebraic query language and the extended logical query language are equivalent, i.e., for every extended algebraic query on  $\mathbf{S}$  there exists a safe extended logical query on  $\mathbf{S}$  and for every extended logical query on  $\mathbf{S}$  there exists an extended algebraic query on  $\mathbf{S}$  such that both queries define the same schema, and for every database instance  $\mathbf{I}$  on which the logical query is safe, the results of both queries are isomorphic relative to  $\mathbf{S}$ .

## 7. CONCLUDING REMARKS

We have described a new model of data, the Logical Data Model, that is designed to combine the advantages of the existing data models. On the one hand, it enables the database to describe more of the semantics of the data than is possible using the relational model of data. On the other hand, we do not lose the nice properties that relational databases have, in particular the ability to query the database using equivalent nonprocedural and procedural languages. The complexity of our query language is studied in [23, 33, 34].

Our work unifies and generalizes a long sequence of previous works on semantic data models [1, 13–18, 21, 27, 29, 38, 37, 40, 43–46, 48]. The preliminary publication of our results in [32], stimulated later work, [2, 30, 39]. All these models can be viewed as special cases of LDM, designed to make certain classes of queries easier to express or more efficient to implement.

Our model incorporates two important *object-oriented* features: *object identity*, captured by the distinction between object names and object values, and *strong typing*, our types are built from the product, power, and union operations. It does not, however, incorporate other features considered important to object-oriented databases, such as "inheritance," "methods," "encapsulation" (see [7]). An important feature of our model is the separation of schema and instance, unlike the models in [8, 36].

One of the features of LDM—the ability to describe cyclic data—is lacking in the query language. This brings up two questions: Is it possible to eliminate cycles from databases? Is it possible to add cycles to queries? The first question is addressed in [33]. More recently, [3] shows how to extend the query language to allow cycles in queries.

The query languages that we have described are based on the paradigm of first-order logic. It is by now recognized that in the framework of the relational model, first-order languages are too weak for the task of database querying [6]. This motivated the study of more powerful query languages, based on the paradigm of logic programming [10, 22, 42, 50]. More recently, similar query languages have been developed for models similar to LDM [2, 3, 4, 8, 28, 35].

## ACKNOWLEDGMENTS

We are indebted to Jeff Ullman for some of the basic ideas underlying this work. We would also like to thank Rick Hull, Paris Kanellakis, and Dave

Maier, for helpful discussions and suggestions, and the anonymous referees for suggestions for improvements.

#### REFERENCES

- 1. ABITEBOUL, S., AND BIDOIT, N. Non-first normal form relations: An algebra allowing data restructuring. JCSS 33 (1980), 361–393.
- 2. ABITEBOUL, S., AND BEERI, C. On the power of languages for the manipulation of complex objects. Tech. Rep. 846, INRIA, 1988.
- 3. ABITEBOUL, S., AND KANELLAKIS, P. Object identity as a query language primitive. To appear in *Proceedings of the ACM Conference on Management of Data* (1989).
- 4. ABITEBOUL, S., AND GRUMBACH, S. A rule based language with functions and sets. ACM Trans. Database Syst. 16 (1991), 1-30.
- ACZEL, P. Non-Well-Founded Sets. Stanford University, Center for Study of Language and Information Lecture Notes, no. 14, 1988.
- AHO, A. V., AND ULLMAN, J. D. Universality of data retrieval languages. In Conference Record of the Sixth Annual ACM Symposium on Principles of Programming Languages (1979), 110–120.
- 7. BANCILHON, F. Object-oriented database systems. In Proceedings of the 6th ACM Symposium on Principles of Database Systems (1987), 152-162.
- BANCILHON, F., AND KHOSHAFIAN, S. A calculus for complex objects. In Proceedings of the 5th ACM Symposium on Principles of Database Systems (1986), 53–60.
- 9. BANCILHON, F., NAQVI, S., RAMAKRISHNAN, R., SHMUELI, O., AND TSUR, S. Sets and negation in a logic database language (LDL1). In *Proceedings of the 6th ACM Symposium on Principles* of Database Systems (1987), 21–37.
- 10. BANCILHON, F., AND RAMAKRISHNAN, R. An amateur's introduction to recursive query processing strategies. In *Proceedings of the ACM Conference on Management of Data* (Washington, 1986), 16–52.
- 11. CODD, E. F. A relational model of data for large shared data banks. Commun. ACM 13, 6 (1970), 377-387.
- CODD, E. F. Extending the database relational model to capture more meaning. ACM Trans. Database Syst. 4 (1979), 397-434.
- 13. DAYAL, U., AND BERNSTEIN, P. A. On the correct translation of update operations on relational views. ACM Trans. Database Syst. 7, 3 (1982), 381-416.
- 14. FURTADO, A. L., AND KERSCHBERG, L. An algebra of quotient relations. In Proceedings of the ACM International Conference on Management of Data (Toronto, Ont., 1977), 1–8.
- FISCHER, P. C., AND THOMAS, S. J. Operators for non-first-normal-form relations. In Proceedings of the IEEE Computer Software Applications Conference (1983), 464–475.
- 16. GANGOPADHYAY, D., DAYAL, U., AND BROWNE, J. C. Semantics of network data manipulation languages: an object-oriented approach. In *Proceedings of the Eighth International Conference on Very Large Data Bases* (1982), IEEE.
- GRAHAM, M. H. NETS: operations and logic. In A Panache of DBMS Ideas II, F. H. Lochovsky, Ed., 152–179. Tech. Rep. CSRG-101, Computer Systems Research Group, Univ. of Toronto, 1979.
- 18. HARDGRAVE, W. T. Ambiguity in processing boolean queries on TDMS tree-structures: A study of four different philosophies. Tech. Rep. IFSM TR-35, Univ. of Maryland, 1978.
- HULL, R. A survey of theoretical research on typed complex database objects. In *Databases*, J. Paredaens, Ed., Academic Press, 1987, 193-256.
- HULL, R., AND KING, R. Semantic database modeling: survey, applications, and research Issues. ACM Comput. Surv. 19, 3 (Sept. 1987), 201-260.
- 21. HAMMER, M., AND MCLEOD, D. Database description with SDM: a semantic database model. ACM Trans. Database Syst. 6 (1981), 351-386.
- 22. HENSCHEN, L. J., AND NAQVI, S. A. On compiling queries in recursive first-order databases. J. ACM 31 (1984), 47-85.

#### 412 . G. M Kuper and M. Y. Vardı

- HULL, R, AND SU, J. Domain independence and the relational calculus. Tech. Rep. 88-64. Univ. of Southern California, 1989.
- HULL, R., AND YAP, C. K. The format model: A theory of database organization. J. ACM 31, 3 (1984), 518-537.
- JACOBS, B. E. Application of database logic to database design. Tech. Rep. TR-892, Univ. of Maryland at College Park, 1979.
- 26. JACOBS, B. E. On database logic. J. ACM 29, 2 (1982), 310-332.
- JAESCHKE, G., AND SCHEK, H.-J. Remarks on the algebra of non first normal form relations. In Proceedings of the First Annual ACM Symposium on Principles of Database Systems (Los Angeles, Calif., 1982), 124–138.
- KIFER, M., AND LAUSEN, G. F-Logic: A higher order language for reasoning about objects, inheritance, and scheme. In *Proceedings of the ACM Conference on Management of Data* (1989), 134-146.
- 29. KOBAYASHI, I An overview of the database management technology Tech. Rep. TRCS-4-1, Sanno College, Kanagawa, Japan, 1980, 259–11.
- KORTH, H. F., ROTH, M. A., AND SILBERSCHATZ, A. Extended algebra and calculus for non-first-normal-form relational databases. Dept. of Computer Science, TR-84-36, Univ. of Texas at Austin, 1984.
- 31. KUPER, G. M. The logical data model: A new approach to database logic. Ph.D. dissertation, Stanford Univ., Stanford, Calif., 1985.
- KUPER, G. M., AND VARDI, M. Y. A new approach to database logic. In Proceedings of the Third Annual ACM Symposium on Principles of Database Systems (Waterloo, Ont., 1984), 86–96.
- KUPER, G. M., AND VARDI, M. Y. On the expressive power of the logical data model. In Proceedings of the ACM International Conference on Management of Data (Austin, Tex., 1985), 180-189.
- 34. KUPER, G. M, AND VARDI, M. Y. On the complexity of queries in the logical data model. In Proceedings 2nd International Conference on Database Theory. Lecture Notes in Computer Science, 326, Springer-Verlag, 1988, 267–280.
- 35. KUPER, G. M. Logic programming with sets J. Comput. Syst. Sci. 41, 1 (1990), 44-64.
- MAIER, D. A logic for objects. In Proceedings of the Workshop on Foundations of Deductive Databases and Logic Programming (Washington, D.C., 1986).
- MAKINOUCHI, A. A consideration on normal form of not-necessarily normalized relations in the relational data model. In *Proceedings of the Third International Conference on Very Large Data Bases* (Tokyo, 1977), IEEE, 447–453.
- 38 MANOLA, F., AND PIROTTE, A. CQLF-A query language for CODASYL-type databases. In Proceedings of the ACM International Conference on Management of Data (Orlando, Fla, 1982), 94-103
- OZSOYOGLU, Z. M., OZSOYOGLU, M., AND MATOS, V. Extending relational algebra and relational calculus with set-valued attributes and aggregate functions. ACM Trans. Database Syst. 12 (1987), 566-592.
- OZSOYOGLU, Z. M., AND YUAN, L.-Y. A normal form for nested relations In Proceedings Fourth Annual ACM Symposium on Principles of Database Systems (Portland, Ore., 1985), 251-260.
- PECKHAM, J., AND MARYANSKI, F. Semantic data model. ACM Comput Surv. 20, 3 (Sept. 1988), 153-190.
- 42. REITER, R. Deductive question answering in relational databases. In *Logic and Databases*, H. Gallaire and J. Minker, Eds., Plenum Press, 1978, 147–177.
- RAFANELLI, M., AND RICCI, F. L. A data definition language for a statistical database. Tech. Rep. TR-62, IASI-CNR, July 1983.
- 44. SHIPMAN, D. The functional model and the data language DAPLEX. ACM Trans. Database Syst. 6 (1981), 140-173.
- 45. SCHECK, H.-J., AND PISTOR, P. Data structures for an integrated data base management and information retrieval system. In *Proceedings of the IEEE Fourth International Conference on Very Large Data Bases* (1982).

- SMITH, J. M., AND SMITH, D. C. P. Database abstractions: Aggregation and generalization. ACM Trans. Database Syst. 2, 2 (1977), 105-133.
- 47. TSICHRITZIS, D. C., AND LOCHOVSKY, F. H. Data Models. Prentice-Hall, Englewood Cliffs, N.J., 1982.
- 48. TSICHRITZIS, D. C. LSL: A link and selector language. In Proceedings of the ACM International Conference on Management of Data (Washington, D. C., 1976), 123-133.
- 49. ULLMAN, J. D. Principles of Database Systems. Computer Science Press, Rockville, Md., 1982.
- 50. ULLMAN, J. D. Implementation of logical query languages for databases. In Proceedings of the ACM International Conference on Management of Data (Austin, Tex., 1985).
- VARDI, M. Y. Review of "On database logic." J. ACM 29, 2 (1982), 310-332. Zentralblatt für Mathematik, 497.68061, 1983.

Received July 1989; revised July 1992; accepted August 1992.