ON THE ALGEBRAIC SPECIFICATION OF DATABASES

Walter Dosch

Gianfranco Mascari

ari 🏋 Martin Wirsing

Institut für Informatik - Technische Universität München - Arcisstrasse 21 - D-8000 Munich 2 Fed. Rep. of Germany

ABSTRACT

For the conceptual level of database schemes a structured algebraic specification is presented. Within a uniform framework it comprises database abstractions, static and dynamic constraints, and a functional programming language for queries and updates. The specification is analysed w.r.t. syntactic and semantic aspects. Then the behaviour and the implementation of database specifications are discussed. Furthermore, algebraic conditions are presented that guarantee a sound specification according to criteria evolved from database theory.

O. INTRODUCTION

In the field of database specifications there is a growing interest in rigorous formal specifications (/Bjørner 80/, /Brodie, Zilles 81/, /Neuhold, Olnhoff 81/) which support a structured database design (/Ehrig, Fey 81/). Such specifications may essentially be classified according to the techniques used:

In constructive approaches a database scheme is specified by defining an abstract model, that is a specific mathematical structure onto which the universe of discourse is mapped.

In axiomatic approaches the behaviour is directly specified by means of a logic theory without resorting to a particular model. The abstract models of a constructive approach then provide for a concrete implementation.

The algebraic specification technique which combines tools from logic and universal algebra has proved to be a powerful and flexible tool for the formal definition of data structures (/Bauer, Wössner 82/) and programming language semantics. Therefore the development of algebraic database specifications (/Ehrig et al. 78/, /Paolini 81/) seems to be quite promising.

The algebraic specification technique may be employed for both approaches: It can serve to define specific abstract models (/Lockemann et al. 79/, /Hupbach 81/). In this way the meaning of (se"Istituto di Automatica - Università degli Studi di Roma - Via Eudossiana 18 - I-00185 Rome Italy

mantic) data models which often is only partially described can be completely specified. But above all algebraic types can be used to specify database schemes directly.

For every specification technique it is necessary to develop criteria for a sound database definition (/Casanova et al. 81/). One guideline is a syntactic and semantic analysis of database specification in order to check whether the class of admitted models indeed meets the informal ideas in mind.

The algebraic specification technique offers a uniform framework for various subjects for which in database theory different mathematical tools, e.g. set theory, relations (/Codd 70/), logic (/Gallaire 81/), are used side-by-side. Many of the concepts developed as part of this technique seem to adequately meet the requirements of database specifications:

- The signature of an algebraic type specifies a syntactic interface by naming the available sorts, for example attributes or entities, and the operations with their arities.
- The exioms of an algebraic type, being arbitrary first order formulas, provide flexibility in expressing static and dynamic constraints.
- The use of partial functions captures finite errors, which may be induced by constraints, and infinite errors, that is the non-termination of queries or updates.
- In a hierarchy of algebraic specifications the visible behaviour of the nonprimitive parts is described by mapping them into primitive types.
- By an encapsulation mechanism the representation of the objects and operations can be hidden from the user.
- Algebraic specifications allow full parameterization. In this way the specification of database schemes gets possible with full generality.
- Algebraic specification languages, like CLEAR (/Burstall, Goguen 80/), and ASL (/Wirsing 82/). give precise semantics to flexible mani-

Proceedings of the Eighth International Conference on Very Large Data Bases pulations with theories. For example, the specification of a large database scheme (/Albano, et al. 81/, /Hammer, Berkowitz 80/, /Ehrig et al. 78/) can be decomposed into parts of manageable complexity and, vice-versa, theories for small parts, like views of different users, may be put together to get an overall specification.

The paper is organized as follows: Section 1 contains a short overview on some basic notions of algebraic specifications. In section 2 external and conceptual database schemes are defined and related to each other by an abstract domain equation. In section 3 a conceptual database scheme is specified by a hierarchy of algebraic types using algebraically defined database abstractions. In section 4 a syntactic analysis of the specification leads to a classification of query and update operations as well as static and dynamic constraints. Section 5 to 7 provide tools for a semantic analysis: In section 5 criteria for the existence of Armstrong models are presented. In section 6 the behaviour of databases is formalized and the existence of behaviour models is discussed. In section 7 the implementation of database schemes is defined and some characteristic implementation steps are outlined. In section 8 a functional programming language for updates and queries is introduced. Finally in section 9 a notion of recursion complete language for updates and queries is algebraically defined, its equivalence to the "extended completeness " of /Chandra, Harel 80/ is outlined and the recursion completeness of the functional programming language (over set-like data structures) is shown.

This completes the aim of the paper to provide an abstract programming language model for database schemes. Several examples accompany the presentation. In the notation we largely follow the wide spectrum language CIP-L (/Bauer et al. 81/).

1. BASIC DEFINITIONS

Below we summarize some basic notions of algebraic specifications which extend the wellknown theory to partial algebras. For a more detailed treatment see for example /Wirsing et al. 80/ and /Wirsing 82/.

Readers interested in an informal survey of this paper may skip this section.

1.1 SIGNATURES

A signature $\Sigma = \langle S, F \rangle$ comprises a set S of sorts and a set F of operation symbols together with their arities of S* × S. The union of signatures $\langle S_1, F_1 \rangle$ means $\langle S_1 \cup S_2, F_1 \cup F_2 \rangle$. A pair $\langle \alpha, \beta \rangle$ forms a signature morphism σ : $\Sigma_1 \rightarrow \Sigma_2$ if α : $S_1 \rightarrow S_2$ is a total mapping for sorts and β : $F_1 \rightarrow F_2$ a family of mappings for operation symbols compatible with α . For a

S-indexed family X of variables $W(\Sigma, X)$ denotes the set of all (finite) Σ -terms including the ground terms $W(\Sigma) = W(\Sigma, \emptyset)$.

1.2 PARTIAL ALGEBRAS

In a (partial) Σ -algebra $A = \langle (S^A)_{G \in S}, (f^A)_{f \in F} \rangle$

with each sort s of S a carrier set A and with each operation symbol f: $sl_{\star}..._{\star}sk \rightarrow s \in F$ a (possibly partial) function $f^{A}:A_{s1}^{\star}..._{\star}A_{sk} \rightarrow A_{s}$ is associated. A is finitely generated if every element of its carrier sets can be obtained by the interpretation t^{A} of a ground term $t \in W(\Sigma)$. The term algebra $W(\Sigma, X)$ and the ground term algebra $W(\Sigma)$ are total finitely generated Σ -algebras. For a signature morphism $\sigma: \Sigma_{1} \rightarrow \Sigma_{2}$ and a Σ_{2} -algebra A^{2} the σ -restriction $A^{2}/_{\sigma}$ is the Σ_{1} -algebra A^{1} with $A_{S}^{1} = A_{\sigma(S)}^{2}$ and $f^{A^{1}} = \sigma(f)^{A^{2}}$ (s, $f \in \Sigma_{1}$).

1.3 HOMOMORPHISMS, EXTREMAL ALGEBRAS

A family $(\varphi_{S} : A_{S}^{1} \rightarrow A_{S}^{2})_{S \in S}$ of total mappings is called (strong) Σ -homomorphism $\varphi : A^{1} \rightarrow A^{2}$ between two Σ -algebras A^{1} , A^{2} if φ preserves the definedness (and undefinedness) of functions and is compatible with them.

In a class C of Σ -algebras the isomorphism classes I resp. Z of (strongly) initial resp. (strongly) terminal algebras are characterized by the existence of (strong) homomorphisms $\varphi : I \rightarrow A$ resp. $\varphi : A \rightarrow Z$ for all $A \in C$.

1.4 FORMULAS

 Σ -formulas are all first-order formulas built from the definedness predicate $D_S(t)$ and the strong equality $t_1 = t_2$ as atomic formulas using the quantifiers \forall , \exists and the connectives \land , \lor , \neg , \Rightarrow . The satisfaction $A \models \Phi$ of a Σ -formula \blacklozenge in a Σ -algebra A is defined as usual where for atomic formulas and allquantification

1.5 ALGEBRAIC TYPES

An algebraic type $T = \langle \Sigma, E \rangle$ consists of a signature Σ and a (countable) set E of closed Σ -formulas, called axioms. Its semantics Mod(T) is the class of all Σ -algebras where all axioms $\phi \in E$ are satisfied and $A \models \underline{true} + \underline{false}$. A consistent type , that is a type T with

Mod(T) $\neq \emptyset$, is called monomorphic, if Mod(T) comprises only isomorphic algebras, and polymorphic elsewise. For a type T a Σ -formula ϕ is provable $(T - \phi)$, if ϕ is deducible from E, valid $(T = \phi)$, if all models of T satisfy ϕ . A formula $\forall sx : \phi(x)$ where ϕ is quantifier-free is called maximal if for all ground terms t T = D(t) implies T $\models \phi(t)$ or T $\models \neg \phi(t)$.

1.6. HIERARCHICAL ALGEBRAIC TYPES

To structure the specification in a hierarchical algebraic type $T = \langle \Sigma, E, T^1 \rangle$ a primitive (sub-) type $T^1 = \langle \Sigma^1, E^1, T^2 \rangle$ is designated consisting of a primitive signature $\Sigma^1 \subseteq \Sigma$ of primitive sorts $S^1 \subseteq S$ and primitive operations $F^1 \subseteq F$, a subset $E^1 \subseteq E$ of primitive axioms and a (possibly empty) primitive subtype T^2 .

The primitive terms $W(\Sigma^1, X^1)$ and the terms of primitive sort form subsets of $W(\Sigma, X)$. In a hierarchical type the non-primitive objects are specified by their visible behaviour under output operations, that is by operations leading from a nonprimitive sort into a primitive sort.

The visible behaviour should be specified sufficiently precise. T is called weakly sufficiently complete, if for all ground terms $t \in W(\Sigma)$ of primitive sort there exists a primitive term $p \in W(\Sigma^1)$ with $\exists \vdash D(t) = T \vdash t = p$, and sufficiently complete, if furthermore $T \vdash D(t)$ or $T \vdash \neg D(t)$. A related model-oriented notion is : A type $T = \langle \Sigma, E, T^1 \rangle$ is hierarchy persistent if for every model A¹ of the primitive type T¹ there is a model A of T such that $A/_{\Sigma^1}$ is A¹, and for every model B of T the reduct $B/_{\Sigma^1}$ is a model of T¹.

1.7 STRUCTURED ALGEBRAIC SPECIFICATIONS

In an algebraic specification language for the manipulation of algebraic types composed type expressions can be built from signatures and sets of axioms as atomic expressions using

- the quotient T : E of a type expression T and a set E of axioms,
- the sum T¹+ T² of two type expressions T¹, T²,
- the restriction T/ σ of a type expression and a signature morphism $\sigma: \Sigma' \rightarrow \Sigma$. If $\Sigma' \subseteq \Sigma$ and σ is the inclusion operation, we write T/ $_{\Sigma'}$ instead of T/ $_{\sigma}$.
- the abstraction (PAR): T of a type expression and a type variable PAR yielding a parameterized type or type scheme with parameter type PAR ,
- the application (imstantiation) T(ARG) of a type scheme T to an argument type ARG,
- the constraint <u>data(T)</u> to finitely generated and minimally defined models of I.

This kernel language can be extended by notational variants like (data-)enrich by. The result of enriching a type $T = \langle \Sigma , E \rangle$ by $\langle \Sigma', E' \rangle$ is the type $T' = \langle \Sigma \cup \Sigma', E \cup E' \rangle$. A type T is data-enriched by $\langle \Sigma', E' \rangle$ if for every model A' of T' the reduct A' $/_{\Sigma}$ is a finitely generated minimally defined model of T.

The precondition notation

<u>funct</u> $(s x : pre(x)) \le f$

in the signature of a type is an abbreviation for the declaration

and the axiom

$$pre(x) = false \Rightarrow D(f(x)) = false$$
.

2. EXTERNAL AND CONCEPTUAL LEVEL OF DATABASES

At the external level a database (scheme) may be seen as a black box. The only way to talk about the information contained in it is to put queries to it. There are different (groups of) users each interacting with the database using a specific set of queries and updates. In this situation both the complete database DB and the view V_i of a user may be described by hierarchical algebraic types. For each user there is a signature Σ_i naming the available sorts and operations. From the conceptual level DB an external view V_i is obtained by restricting DB to Σ_i ; conversely a database is the sum of its views where certain interferences among the views V_i are respected.

Definition A conceptual database scheme DB and an external database scheme or view V_i are parameterized hierarchical algebraic types such that

(1) $DB = (V_1 + ... + V_n)$: Interference $(V_1, ..., V_n)$ (2) $V_i = DB/_{\Sigma_i}$ for i = 1, 2, ..., n holds,

where $\Sigma_i = sig(V_i)$ are signatures and

Interference (V_1, \ldots, V_n) are $\Sigma_1 \cup \ldots \cup \Sigma_n$ - axioms.

This definition reflects two directions in the conceptual design:

- a) First the external view V_i and their interferences are specified separately and then integrated into a conceptual database scheme DB.
 Algebraically this means to solve the abstract domain equation
- (3) $DB = (DB/\Sigma_1 + ... + DB/\Sigma_n)$: Interference $(DB/\Sigma_1, ..., DB/\Sigma_n)$ $(n \ge 1)$ for the unknown type variable DB.

Proceedings of the Eighth International Conference on Very Large Data Bases

Mexico City, September, 1982

b) First the conceptual level DB is specified and then the different external views V_i are derived from it according to equation (2).

As it is well known from fixpoint theory of recursive functions the domain equation (3) in general has a lot of solutions. When using the method b) the specified conceptual scheme often contains more detailed information than needed for the different external schemes, then DB is not a minimal solution of equation (3). Furthermore when using database abstractions to construct a conceptual scheme, in general not a minimal solution of (3) is obtained.

3. SPECIFICATION OF THE CONCEPTUAL LEVEL

Below we specify the conceptual level of a database scheme by defining a hierarchy of algebraic types for its components.

3.1 OBJECT TYPES

The conceptual database specification as well as the views are modular composed (/Schiel et al. 82/) of various subtypes such as several object types OBJ_1, \ldots, OBJ_n that define the data structures of the data items.

In general the basic object types OBJ_i are not monomorphic. Nevertheless the following requirements should be fulfilled:

<u>Claim</u> The hierarchical types OBJ_i for the basic objects should be specified such that

- the definedness predicate is model-independent,
- 2) they are hierarchy persistent,
- for their semantics only finitely generated models are regarded.

Condition 1) ensures that the definedness of terms is model independent, 2) allows to implement the overall type by using implementations of its primitive types, 3) ensures that every object has a finite denotation and can be finitely computed, i.e., there is "no junk" in the models.

3.2 DATABASE ABSTRACTIONS

In general the basic object types OBJ_i are them-

selves obtained in a hierarchical way based on other primitive types. To abbreviate type schemes frequently occurring in the construction of database specifications database abstractions have been introduced (/Hammer 76/, /Smith, Smith 77/, /dos Santos et al. 80/, /Brodie 81/). They are similar to generic mode constructors in programming languages (/Schmidt 80/). Below we specify aggregation, generalization, and correspondence.

3.2.1 AGGREGATION

<u>Definition</u> An aggregation is the type scheme <u>type</u> AGGREGATION = $(OBJ_1, ..., OBJ_n, \sigma, PRE, \Sigma, E)$: <u>data-enrich</u> PROD $(OBJ_1, ..., OBJ_n)/\sigma$: PRE <u>by</u> Σ , E endoftype,

where the type scheme PROD forms the cartesian product of its parameter types (/Bauer, Wössner 82/)

 $\underline{type} \ PROD = (OBJ_1, \dots, OBJ_n) :$ $\underline{data-enrich} \ OBJ_1, \dots, OBJ_n, BOOL \ \underline{by}$ $\underline{sort \ prod},$ $\underline{funct} \ (\underline{obj}, x_1, \dots, \underline{obj}, x_n :$ $\underline{pre(x_1, \dots, x_n)} \ \underline{prod} \ mk,$ $\underline{funct} \ (\underline{prod}) \ \underline{obj}_i \ sel_i,$ $\underline{funct} \ (\underline{obj}, \dots, \underline{obj}_n) \ \underline{bool} \ pre,$ $\underline{pre(x_1, \dots, x_n)} \rightarrow sel_i \ (mk(x_1, \dots, x_n)) = x_i$

endoftype

If the precondition pre is constant true,

 $PRE = \{ pre(x_1, ..., x_n) = true \}$

then mk is a total operation; in this case we omit pre and PRE and write TPROD for PROD .

Example

A secretary may be characterized by the components name, age (over 18) and typing speed. Furthermore an operation to increase the age is provided.

type SECRETARY =:

data-enrich (secr, mksecr, sname, sage, speed,

isadult) ≡

PROD(NAME, AGE, TYPING_SPEED) :

 $isadult(n, a, s) = (a \ge 18)$

by funct (secr) secr incrage,

incrage(mksecr(n, a, s)) = mksecr(n, a+1, s)

endoftype

Here the notation

(secr, mksecr, ..., isadult) PROD(NAME, AGE, TYPING_SPEED)

abbreviates the signature morphism σ of PROD/ σ by establishing the correspondence

secr \rightarrow prod, mksecr \rightarrow mk, ..., pre \rightarrow isadult.

Similarly one may specify a type DIRECTOR consisting of the component types NAME, AGE, and TELEPHONE.

3.2.2 GENERALIZATION

Definition A generalization is the type scheme

<u>type</u> GENERALIZATION = $(OBJ_1, \ldots, OBJ_n, \sigma, PRE, \Sigma, E)$: <u>data-enrich</u> SUM(OBJ ,...,OBJ)/ $_{\sigma}$: PRE by Σ , E endoftype . where the type scheme SUM defines the direct sum (disjoint union) of its parameters: <u>type</u> SUM = (OBJ_1, \dots, OBJ_n) : data-enrich OBJ , ..., OBJ , BOOL by sort sum, funct (obj; x;: pre;(x;)) sum mk; , funct (sum) bool is, , funct (sum s : is_i(s)) obj_i pr_i , funct (obj;) bool pre; , $pre_i(x_i) \Rightarrow [is_i(mk_i(x_i)) = true],$ $is_i(mk_i(x_i)) = \underline{false}$, $(i \neq j)$ $pr_{i}(mk_{i}(x_{i})) = x_{i}$] endoftype

Example

The staff of a company consists of either secretaries with excellent typing speed or directors older than 40 years. It additionally may have an operation to access the name of the staff directly.

```
type STAFF = :
```

3.2.3 CORRESPONDENCE

Many (semantic) data models widely use sets but nevertheless do not support the data structure set explicitly. Sets as objects are required in the conceptual modelling whenever sets have specific properties exceeding the properties of their elements.

```
Definition A correspondence is the type scheme
type CORRESPONDENCE = (OBJ, \sigma, PRE, \Sigma, E) :
     data-enrich FINSET(OBJ)/\sigma : PRE by \Sigma, E
endoftype
where the type scheme FINSET defines finite sets
of a member type OBJ (cf. /Wirsing et al. 80/) :
type FINSET = (OBJ) :
   data-enrich OBJ, BOOL by
   sort finset .
   funct finset emptyset ,
   funct (finset s, obj x :
          preins(s, x)) finset insert,
   funct (finset s, obj x :
          predel(s, x)) finset delete ,
   funct (finset) bool isempty ,
   funct (finset, obj) bool iselem ,
   funct (finset s : -, isempty(s)) obj choose ,
   funct (finset, obj) bool preins, predel ,
   iselem(emptyset, x) = false ,
   preins(s, y) =
   iselem(insert(s,y), x)= (eq(x,y) v iselem(s,x)),
   predel(s, y) ⇒
   iselem(delete(s,y), x) = (\neg eq(x,y) \land iselem(s,x)),
   isempty(s) \leftrightarrow \neg \exists obj x : iselem(s, x),
   ¬ isempty(s) ⇒ iselem(s, choose(s)) = true
endoftype
```

Here eq denotes an equality operation on OBJ. Note that the result of the partial operation choose ("Give me some element of the set") is undefined for the empty set, uniquely determined for sets with one element, and ambiguously for sets with two or more elements. Furthermore, delete is specified as a constructor (cf. also section 5).

Example

A secretary belongs to a set of candidates when she submits to work and is no longer a candidate when she withdraws. This has to respect the following constraints: A secretary withdrawing (submitting) must (not) be a candidate.

type CANDIDATES = (SECRETARY) :

data (cands, emptycand, csubmit, cwithdraw, cisempty,

iscand, cchoose, precsubmit, precwithdraw)

= FINSET(SECRETARY) :

precsubmit(c, s) = - iscand(c, s),

precwithdraw(c, s) = iscand(c, s)

endoftype

3.2.4 AN EXAMPLE: EMPLOYMENT AGENCY

In this section for the well-known example of an employment agency (/Veloso et al. 81/) a conceptual scheme is specified by a hierarchy of algebraic types using database abstractions. It is based on the primitive types SECRETARY and COMPANY for secretaries and companies.

A secretary working for a company is an employee. The following type assumes no constraints, that is every secretary may work for every company.

```
type EMPLOYEE = (SECRETARY, COMPANY) :
    (empl, < .,.>, secr, comp) =
        TPROD(SECRETARY, COMPANY)
```

endoftype

In the next step the set of employees is specified with the following constraints:

- A secretary cannot be hired by a company when she is already working for any company.
- She only can be fired by a company when she is working for this company.

type EMPLOYEES = (EMPLOYEE) :

data (empls, emptyempls, ehire, efire, eisempty, worksfor, echoose, preehire, preefire) = FINSET(EMPLOYEE) :

preehire(e, <s, r >) =

 $(\forall \underline{comp} r_1: \neg worksfor(e, <s,r_1>))$, preefire(e, <s,r>) = worksfor(e, <s,r>))

endoftype

Then an employment agency consists of a set of candidates and a set of employees respecting the following integrity constraint :

(PRE) No secretary can simultaneously be a candidate and an employee.

Moreover there are the following constraints:

(PRES) A secretary submitting can neither be a candidate nor be working for a company.

(PREH) A secretary can only be hired by a company if she is a candidate and does not work for any company.

(PREF) A secretary getting fired by a company must work for this company.

type EMPL_AGENCY = (CANDIDATE, EMPLOYEE) :

<u>data-enrich</u> (<u>ag</u>, <.,.>, cands, empls, preag) = PROD(CANDIDATES, EMPLOYEES) :

funct (ag a, empl e : prefire(a, e)) ag fire, funct (ag, empl) bool prehire, prefire,

preag(c, e) 👄

(PRES) [presubmit(< c,
$$e$$
 >, s) = ($\neg iscand(c,s) \land \forall comp r : \neg worksfor(e, < s, r>))$,

PREH) prehire(,) = (iscand(c,s)
$$\land$$

 $\forall comp r_1 : \neg worksfor(e,$

(H) prehire(,)
$$\Rightarrow$$
 hire(,) =)>,

endoftype

For the consistency proof it is necessary to check that the update operations submit, hire, and fire preserve the integrity constraint (PRE):

presubmit(, s)
$$\Rightarrow$$
 preag(csubmit(c,s), e) ,

prehire(<c,e >,<s,r>) +

preag(cwithdraw(c,s), ehire(e, <s,r>)) ,

prefire(<c,e>, <s,r>) →

preag(csubmit(c,s), efire(e, <s,r>))

This can be proved using the preconditions (PRE), (PRES), (PREH), and (PREF). For example, the application of submit keeps the integrity whereas csubmit may violate it. For the final specification of the employment agency all updates have to be hidden which possibly destroy the integrity:

```
<u>type</u> EMPLOYMENT_AGENCY = (CANDIDATE, EMPLOYEE) :
EMPL_AGENCY(CANDIDATE, EMPLOYEE)/
```

endoftype where

 $\Sigma = sig(EMPL_AGENCY) >$

{csubmit, cwithdraw, ehire, efire }

For this type requirements of /Veloso et al. 81/ are provable (as postconditions), for example: (POSTS) a = submit(al, s) → iscand(cands(a), s) ∧ ∨ comp r : ¬ worksfor(empls(a),<s, r>).

More complex operations like

funct (ag, secr) bool works

('Does a secretary work?') may be specified at this level by an enrichment using descriptive formulations like

works(a,s) = $\exists comp r : worksfor(a, <s,r>)$

or structural induction

```
works(<c, emptyempls>,s) = false ,
```

preag(c, ehire(e, <s,,r>)) ⇒

works(<c, ehire(e, $\langle s_1, r \rangle \rangle$, s_2) =

 $(eqs(s_1,s_2) \lor works(\langle c,e \rangle,s_2)),$ preag(c, efire(e, $\langle s_1, r \rangle)) \Rightarrow$

works(<c, efire(e, $\langle s_1, r \rangle \rangle$, s_2) =

 $(\neg \operatorname{eqs}(s_1, s_2) \land \operatorname{works}(\langle c, e \rangle, s_2)),$ Alternatively these operations may be defined in the next level using a programming language, see section 8.

4. SYNTACTIC ANALYSIS

In this section algebraic specifications for database schemes are analysed under syntactic aspects.

The hierarchical structure of algebraic types induces a syntactic classification of the operations as well as of the axioms:

Definition

For a database specification $DB = \langle \Sigma, E, OBJ_1, \dots, OBJ_n \rangle$ with primitive subtypes OBJ_i we call an operation symbol f of arity funct $(\underline{s}_1, \dots, \underline{s}_k) \leq \underline{s}_k$ query operation if the range \underline{s}_i of f is a primitive sort, that is $\underline{s}_i \in \bigcup_{i=1}^{n} sig(OBJ_i)$, and update operation otherwise.

Example

In the type SECRETARY sname, sage, speed, and isadult are query operations, whereas mksecr and incrage are update operations.

In the type EMPLOYMENT_AGENCY for example iscand and worksfor are queries whereas hire, fire, and submit are updates.

Now let DB be a hierarchical algebraic type specifying a (view of a) conceptual database (scheme) which is based on the primitive object types OBJ₁, ..., OBJ_n. Then the signature

$$\tilde{\Sigma}$$
-DB = sig(DB) $\sim \bigcup_{i=1}^{n}$ sig(OBJ_i)

contains all sorts and operation symbols of DB without the sorts and operations of the primitive object types.

Definition

An operation $f \in sig(OBJ_i)$ is called object operation, and an operation $f \in \Sigma$ -DB is called

holding operation.

These two definitions allow a two dimensional classification of the database operations into holding updates and holding queries as well as in object updates and object queries.

Example

In the EMPLOYMENT_AGENCY there are the object queries sname, sage, speed, the object update incrage, the holding queries iscand and worksfor, and holding updates submit, hire, and fire.

This classification of the operations may also be used to discriminate two kinds of constraints.

Definition

In a database specification DB a formula ϕ with DB $\models \phi$ is called static constraint iff there does not occur any operation f in ϕ , which is a hciding update. Otherwise ϕ is called <u>dynamic</u> constraint.

Example

In the EMPLOYMENT_AGENCY the (provable) formula

mksecr(sname(s), sage(s), speed(s)) = s

as well as the axiom (PRE) are static constraints, whereas the preconditions (PRES), (rREH), (PREF) and the postcondition (POSTS) are dynamic constraints.

In the following sections the specification for database schemes is semantically analysed with respect to notions evolved from database theory. These notions are also related to algebraic concepts.

5. ARMSTRONG MODELS

In order to get an overview of the properties of a specification special models are studied.

Definition

An algebra A is called Armstrong model (/Makowsky 81/, /Fagin 82/) of a database specification DB, if all first order formulas which are provable in DB and which only contain terms the definedness of which is provable, are valid in A and no others.

The existence of Armstrong models is often accepted as a criterium for a sound specification since they help the database designer to see what he has implied by the axioms of his specification. Thus a major concern of research in database theory is to characterize the class of formulas for the description of data dependencies which allow or guarantee Armstrong models. For algebraic types with partial functions there is the result:

Theorem

Let a database be specified by a consistent hierarchical type $DB = \langle \Sigma, E, OBJ \rangle$ where all axioms in E have positive conditional form

$$\forall \underline{s} \times_1 \dots \underline{s}_n \times_n : \bigwedge_{i=1}^{n} D(p_i) \wedge p_i = q_i \Rightarrow C$$
,

where C is an atomic formula p = q or D(p).

- DB has an Armstrong model, iff DB is weakly sufficiently complete.
- (2) Every Armstrong model is initial in DB.
- (3) An Armstrong model is strongly initial in DB if for all ground terms $t \in W(\Sigma)$ either DB $\models D(t)$ or DB $\models \neg D(t)$.

Proof

See /Wirsing et al. 80/.

Q

For the database abstractions presented in section 3 the existence of Armstrong models is ensured whenever the preconditions and the parameter types behave well.

Proposition

For the database abstraction AGGREGATION (GENER-ALIZATION; CORRESPONDENCE) there exist Armstrong models, iff the operations pre (pre₁; preins,

predel, choose) specified are weakly sufficiently complete and their parameter type(s) OBJ,..., OBJ_n (OBJ_1 ,..., OBJ_n ; OBJ) allow an Armstrong model.

Example

If cchoose is specified weakly sufficiently complete, then CANDIDATES has an Armstrong model I, where for example

(*) csubmit(emptycand, s) +

csubmit(cwithdrwaw(csubmit(emptycand,s),s),s)

holds. More generally, Armstrong models of type FINSET (if they exist) are based on sequences where multiplicity and ordering of the sequence elements is relevant for the equality. Also for implementing delete(s, x) the element xmust be appended to the sequence s with a marker for being no longer present. Thus in I the candidate's history is completely reflected.

Also the type EMPLOYMENT_AGENCY has a (strongly initial) Armstrong model if cchoose and echoose are specified sufficiently complete. Similar to CANDIDATES in an Armstrong model of EMPLOYMENT_ AGENCY the equality between objects of sort ag is determined by the sequence of updates submit, hire, fire leading to them. Thus Armstrong models contain redundant information which cannot be retrieved by queries. In the partial algebra approach Armstrong models are minimally defined and defined terms are interpreted as different as possible. Thus the existence of Armstrong models is not endangered when adding inequalities such as (*) to a database specification as long as it remains consistent.

6. BEHAVIOUR OF DATABASES

In a view point dual to Armstrong models a database is significantly characterized by its behaviour, that is by all effects which are visible in the primitive types. Thus a model of a database specification is called a behaviour model if two elements are distinguished in the model only if they can be discriminated by queries. Such models can be defined using the behaviour view of a database which abstracts from the equations holding on its nonprimitive sorts.

For simplicity throughout section 6.1 and 6.2 let $DB = \langle \Sigma, E, OBJ \rangle$ be a database specification with exactly one nonprimitive sort $\underline{db} \in \Sigma$ where the primitive type OBJ contains one sort obj.

Definition

The behaviour view Beh(DB, V) of a database scheme DB = $<\Sigma$, E, OBJ > w.r.t. a subsignature $V \subseteq \Sigma$ is given by

Beh(DB, V) = Σ + DB/sig(OBJ) + OUT(V) : EQU(V) where

$$OUT(V) = \{ \underline{funct obj }_{o} \mid t \in W(V) \text{ of sort } \underline{obj } \}$$

$$EQU(V) = \{ t = t_{o} \mid t \in W(V) \text{ of sort } \underline{obj } \}$$

Thus in a behaviour view of a database scheme DB the axioms of the sort \underline{db} are neglected and all terms of W(V) of primitive sort \underline{obj} are added as constants.

6.1 BEHAVIOUR MODELS

The identity of the nonprimitive parts of a database specification may be completely determined by their visible behaviour in the primitive object type OBJ. On the term level, one gets a visible behaviour of a nonprimitive object by putting the corresponding term into a context of a primitive sort.

Definition

- a) A Σ -algebra A of a database specification DB = < Σ , E, OBJ> is called fully abstract, if for all t, t' $\in W(\Sigma)$ of sort db
 - (1) A ⊨ D(t) iff there exists a context c[x] ∈ W(DB,{x}) of sort <u>obj</u> with A ⊨ D(c[t]).
 - (2) A ⊨ t ≠ t' iff there exists a context c[x] ∈ W(DB,{x}) of sort obj with A ⊨ c[t] ≠ c[t'].

b) A fully abstract model of Beh(DB, Σ) is called behaviour model of DB.

Proposition

- a) If a database specification $DB = \langle \Sigma, E, OBJ \rangle$ is sufficiently complete, then all behaviour models of DB are isomorphic.
- b) If furthermore all axioms of DB are of positive conditional form with maximal premisses and the definedness predicate D is model independent in DB, then the following is equivalent for an algebra A with signature of Beh(DB, Σ) :
- (1) A/ $_{\tau}$ is a strongly terminal model of DB.
- (2) A is a strongly terminal model of Beh(DB, Σ)
- (3) A is a behaviour model of DB.

Corollary

For the database abstractions AGGREGATION and GENERALIZATION Armstrong models and behaviour models coincide if the precondition operations preresp. pre_i are sufficiently complete. In this

case, AGGREGATION and GENERALIZATION are parameter monomorphic.

For the database abstraction CORRESPONDENCE all behaviour models are isomorphic if the operations preins, predel, and choose are sufficiently complete.

Example

The types SECRETARY and STAFF are parameter monomorphic.

The type CANDIDATES has behaviour models satisfying (in contrast to Armstrong models) for example the formula

csubmit(emptycand, s) =

csubmit(cwithdraw(csubmit(emptycand,s),s),s)

since its left and right hand side are not distinguishable by contexts of primitive sort.

6.2 DISTINCTION BY QUERIES

In a database specification the effects of update operations are visible to a user only by the change of results of queries. Thus internal states with different behaviour should be distinguishable by query operations only.

Definition

For a set Q of query operations a database specification DB = $\langle \Sigma, E, OBJ \rangle$ is called query distinctive, if any two states db, db' $\in W(\Sigma)$ of sort <u>db</u> are distinguishable by queries, that is <u>Det(DD S)</u>

 $Beh(DB, \Sigma)/_{\Sigma^i} = DB^i$

where Σ' is the signature of DB' = Beh(DB, $\{\underline{db}, \underline{obj}\}, \mathbb{Q}\}$). For database specifications that are sufficiently complete and query distinctive, behaviour models are completely characterized by queries.

Proposition

Let DB = $\langle \Sigma, E, OBJ \rangle$ be a database specification which is sufficiently complete and query distinctive. A Σ -algebra A is a behaviour model of DB iff for all db, db' $\in W(\Sigma)$ of sort db A \models db \ddagger db' \Leftrightarrow there exists a holding query q $\in \Sigma$ with DB \models q(db) \ddagger q(db') A \models D(db) \Leftrightarrow there exists a holding query q $\in \Sigma$

For the database abstractions of section 3 we get in particular:

with $DB \models D(q(db))$.

Proposition

AGGREGATION, GENERALIZATION, and CORRESPONDENCE are query distinctive if the signature morphism σ does not forget any query operation.

Example

The type SECRETARY is query distinctive, since two objects of sort secr are different whenever they differ in at least one component NAME, AGE, or TYPING SPEED.

The distinction by queries is transitive with respect to parameter passing in parameterized types:

Proposition

If the database scheme DB(PAR) is query distinctive (w.r.t.its parameter type PAR) and ARG is query distinctive (w.r.t. its primitive type OBJ), then the database DB(ARG) is query distinctive (w.r.t. OBJ).

Corollary

If a database specification DB is built only by database abstractions then it is query distinctive if no query operations are forgotten by the signature morphisms used.

Example EMPLOYMENT AGENCY is query distinctive.

Finally, all composed data structures which do not allow to access all their components directly by a single operation are, in general, not query distinctive.

Example

Let the type FILE with parameter OBJ specify files using the updates emptyfile, rest, put, and the query operations get, eof. Then FILE is not query distinctive since files differing in the second, third, ... element cannot be distinguished by the queries get and eof. Note that the procedure get of PASCAL corresponds to the composition get-rest (see section 9): PASCAL-files are not query distinctive.

6.3 KEYS

In a database specification $DB = \langle \Sigma, E, OBJ_1, \dots, OBJ_n \rangle$

with various primitive types OBJ, the queries

leading to a distinguished primitive type OBJio may be sufficient to discriminate the nonprimitive objects:

Definition

A primitive type OBJ_{io} of a database specification $DB = \langle \Sigma, E, OBJ_1, \ldots, OBJ_n \rangle$ is called key if

DB is query distinctive w.r.t.all operations with range in a sort of $OBJ_{\rm in}$.

Example

Let secretaries be uniquely distinguishable by their names. Then NAME may serve as key since

SECRETARY/sig(SECRETARY)>{sage, speed, isadult}

is query distinctive, whereas

SECRETARY/sig(SECRETARY)>{sname}

will, in general, not be query distinctive.

7. IMPLEMENTATION AND BEHAVIOUR EQUIVALENCE OF DATABASES

At the beginning of a database design there stands the requirement analysis formalized as external or conceptual scheme. The conceptual scheme then is the starting point for a joint development towards an internal scheme. This development can be done for example by a stepwise refinement technique (cf. /Ehrig, Fey 81/) or by a series of transformation steps meliorating coherently algorithms and data structures (cf. /Bauer, Wössner 82/). The correctness of such refinement or transformation steps can be expressed and formalized by the notion of implementation. Each subsequent specification is an implementation of the previous one. The algebraic definition of an implementation does not depend on the structure of the specified database but only takes its behaviour into account (cf. /Ausiello et al. 80/).

Definition

A database specification $DB_1 = \langle \Sigma_1, E_1, 0BJ \rangle$ is called an implementation of $DB_2 = \langle \Sigma_2, E_2, 0BJ_2 \rangle$ via a signature morphism

K: OUT → OUT

if $\emptyset \neq Mod(Beh(DB_1,\Sigma_1)/\kappa \subseteq Mod(Beh(DB_2,\Sigma_2)/_{OUT_2})$ where $OUT_i \equiv sig(OBJ_i) \cup OUT(\Sigma_i)$

(see also section 6) .

DB, and DB_ are called behaviourally equivalent if they implement each other. Informally DB, implements DB, if the visible parts of the models of DB, agree (after renaming) with the visible parts of models of DB₂. Therefore the semantics of queries is not affected by an implementation whereas the semantics of updates, in general, is.

When building hierarchies of data abstractions composite objects often share common components. In an implementation step these components are separated : The sum of two products with common subtypes can be implemented by the product of the common subtypes and the sum of the other components.

and

Proposition -

Let
$$OBJ_{=} PROD_{a}(OBJ, OBJ_{+})$$

 $OBJ_{=} PROD_{b}(OBJ, OBJ_{+})$.

Then SUM(OBJ, OBJ₂) is implemented by PROD(OBJ,SUM'(OBJ', OBJ')) via the signature morphism κ induced by

$$\kappa (mk_1(mk_a(x, y_1)) = mk(x, mk'_1(y_1)) \\ \kappa (mk_2(mk_b(x, y_2)) = mk(x, mk'_2(y_1))$$

Of course the implementation is more economic since the obj-component occurs only once in the product. However, as can be easily seen, both types are behaviourally equivalent.

Example

The type STAFF can be implemented by PROD(NAME, AGE, SUM(TYPING_SPEED, TELEPHONE)) .

A quite different implementation step usually is applied to the set-like database abstraction CORRESPONDENCE. After introducing a key, tables (specified by the type scheme GREX below) may be used to implement the correspondence.

Definition An association is a type scheme

type ASSOCIATION = (KEY, OBJ, σ , PRE, Σ , E) :

<u>data-enrich</u> GREX(KEY, OBJ)/ $_{\sigma}$: PRE <u>by</u> Σ , E endoftype where

Proceedings of the Eighth International Conference on Very Large Data Bases

isin(endt, 1) = false,
preput(g, k, x)
$$\Rightarrow$$

isin(put(g,k,x), 1) = (eqk(k,1) v isin(g,1)),
isin(put(g,k,x), 1) \Rightarrow get(put(g,k,x), 1) =
if eqk(k,1) then x else get(g,1) fi,
predel(cancel(put(g, k, x), 1)) \Rightarrow
cancel(put(g, k, x), 1) =
if eqk(k, 1) then if isin(g, 1)
then cancel(g, 1)
else g fi
else put(cancel(g,1),k,x) fi

endoftype

Proposition

Let OBJ = PROD(KEY, OBJ¹) such that for all key k, finset s there exists at most one

 $o \in obj^1$ with iselem(s, <k, o>). Then CORRESPONDENCE(KEY, OBJ) can be implemented by ASSOCIATION(KEY, OBJ) via the signature morphism κ induced by

```
\kappa(\text{insert}(g, x)) = \text{put}(\kappa(g), \text{sel}_1(x), \kappa)

\kappa(\text{delete}(g, x)) = \text{cancel}(\kappa(g), \text{sel}_1(x))

\kappa(\text{iselem}(g, x)) = \text{isin}(\kappa(g), \text{sel}_1(x))
```

and some particular operations getsome and isinit for κ (choose) and κ (isempty).

The particular choice of getsome such as "get an object with minimal key" often implies that COR-RESPONDENCE and ASSOCIATION are not behaviourally equivalent.

Example

The type CANDIDATES can be implemented by the following type

type CANDIDATES' = (NAME, SECRETARY) :

data-enrich (cand', cinit', csubmit', cget', cwithdraw', iscand', precwithdraw', precsubmit') = CDEV(NAME_SECDETARY)

```
GREX(NAME, SECRETARY)
```

```
by precsubmit'(g, i, x) = -, iscand'(g,i),
precwithdraw'(g, i) = iscand'(g, i)
```

```
by funct (cand') bool cisinit
```

```
funct (cand'c: ¬isinit(c)) secr getsome
```

```
• • •
```

endoftype

via the morphism induced by

 κ (csubmit(g, x)) = csubmit'(κ (g), sname(x), x), κ (iscand(g, x)) = iscand'(κ (g), sname(x)). 8. AN ABSTRACT PROGRAMMING LANGUAGE MODEL

To express and perform complex database operations a database language is needed. The queries and updates of the database scheme serve as its primitive operations. Below we define a schematic implementation of a conceptual database scheme by a simple but adequate programming language. This completes an aim of this paper to propose an abstract programming language as a semantic database model. The programming language should be a scheme language: its semantics should be completely defined relative to an arbitrary database scheme providing the primitive operation symbols. The axioms for the programming language model are simple translations from the axioms of the database specification enriched by the axioms of the programming language constructs.

As an example we define a functional language (/cf. Buneman, Frankel 80/), in the FP-style of Backus to obtain recursion complete query and update languages.

Throughout this section for simplicity let DB = $<\Sigma$, E, OBJ > be a database specification with exactly one nonprimitive sort db and one primitive type OBJ with primitive sort <u>pr</u>.

8.1 THE BASIC DATABASE LANGUAGE

The basic database language B_DB_L is based on the types DB for the database specification and SEQU(OBJ) specifying sequences of objects. It comprises the sorts <u>query</u> and <u>upd</u> for query resp. update operations.

type $B_DB_L = (DB, OBJ)$:

sort query, upa

Then the database scheme DB is "lifted" by the signature mapping α : sig(DB) \rightarrow sig(B_DB_L) which relates sorts

and function symbols

 $\alpha(\underline{funct}(\underline{s}, \ldots, \underline{s}_n) \underline{s} f) = \underline{funct} \alpha(\underline{s}) \overline{f}$.

Depending on its range s, a function symbol f yields an object \vec{T} either of sort query or upd. Remembering the syntactic classification in section 4 we thus get by definition of α the following correspondence:

holding query f of DB	₩	<u>funct</u> query F
object query f of OBJ	₽	funct query 7
holding update f of DB	↔	funct upd T
object update f of OBJ	↔	funct query Ŧ

Example

In B_DB_L(CANDIDATES, SECRETARY) the operations csubmit, cwithdraw, iscand, and precsubmit of the database CANDIDATES are translated into

funct upd csubmit, cwithdraw

funct query iscand, precsubmit .

8.2 A FUNCTIONAL DATABASE LANGUAGE

The basic database language is enriched by the .anguage constructs $\boldsymbol{\Sigma}$ and axioms \boldsymbol{E} describing a functional programming language.

type $DB_L = (DB, OBJ, ID)$:

data (NAT + ID + CONFIGURATION(DB, OBJ))+
enrich B_DB_L by Σ, Ε

endoftype

Functional programs are constructed using the conditional if . then . else . fi, the functional composition ..., a tuple constructor [.,.] with corresponding selector functions sel and (possibly recursive) function definitions recv and function calls call. For this purpose a type ID of identifiers is assumed. Thus for $v \in$ {query, upd} the signature Σ comprises

funct (query, v, v) v if . then . else . fi ,
funct (v, v') upd ... (upd $\in \{v, v'\}$),
funct (query, query) query ..., [.,.] ,
funct (nat) query sel ,
funct (id, v) v recv ,
funct (id) v call .

In the signature above function symbols with different arities are overloaded; for example there are actually four compositions

funct	(<u>query</u> , <u>upd</u>) <u>upd</u> .•. ,
funct	(upd, query) upd .o.,
funct	(upd, upd) upd .o.,
funct	(query, query) query .o.

three of which yield updates; they may change the state of DB. Note that (recursive) definitions of queries resp. updates have to be discriminated.

To define the semantics let a configuration consist of a database scheme and a sequence of input/ output data:

type CONFIGURATION = (DB, OBJ) :

(conf, <.,.>, db, io) = TPROD(DB, SEQU(OBJ))
endoftype

The following axioms of the semantic function funct (v, conf) conf apply $v \in \{\text{query}, \text{upd}\}$ taken as left to right term rewrite rules describe an innermost (call-by-value) text substitution machine. The auxiliary function subst denotes the usual substitution, that is subst(x, f, g) replaces all free occurrences of call(x) in f by g; its straightforward definition by structural induction is omitted. apply(f_{j} <db,<o₁,...,o_k>>)=< f(db,o₁,...,o_k), <>> for holding updates funct(db,obj,...,obj,)db f, apply(\overline{f} <db,<o,,...,o,>)= <db,f(db,o,,...,o,)>> for holding queries funct(db,obj,...,ubj,)obj, f apply(T, db, $(o_1, ..., o_k) = (db, (f(o_1, ..., o_k)))$ for object operations $funct(obj_1, \dots, obj_k)obj_{k+1}$ f $io(apply(c,\pi)) = \langle true \rangle \rightarrow$ apply(if c then v, else v, fi, π) = apply(v, π) $io(apply(c, \pi)) = \langle false \rangle \Rightarrow$ apply(if c then v, else v, fi, π) = apply(v, π), j≤k⇒ apply(sel(j), < db, < $o_1, \ldots, o_k >>$) = < db, < $o_j >>$, apply($[q_1, q_2], \pi$) = $<db(\pi)$, io(apply(q, π))& io(apply(q, π))>, apply($v_1 \circ v_2, \pi$) = apply($v_1, apply(v_2, \pi)$), apply(recupd(f, v), π) = apply(subst(f, v, recupd(f, v)), π), $\pi_2 \approx apply(subst(f, v, recquery(f, v)), \pi_2) \rightarrow$ apply(recquery(f, v), π_1) = $\langle db(\pi_1), io(\pi_2) \rangle$ For brevity the axioms for context conditions like j > k ➡ $D(apply(sel(j), < db, < o_1, ..., o_{\nu} >>)) = false$ have been omitted. These axioms also specify that basic operations can only be applied according to their arity.

The functional database language has the following semantic properties:

Theorem

The type DB L is weakly sufficiently complete and admits minimally defined behaviour models.

Proof See /Broy, Wirsing 80/.

Note that the minimally defined models correspond to least fixed points and provide a mathematical semantics for DB_L . The following example shows some properties of these models.

Example

In DB_L = DB_L(CANDIDATES, SECRETARY, ID) the following equations hold (cf. section 3.2.3):

apply(precsubmit, π) = apply(not • iscand, π) apply(precwithdraw, π) = apply(iscand, π)

The operation apply is the only operation in DB_L with range in the primitive type CONFIGU-RATION; therefore all contexts of primitive sort for functional programs have the form apply(c [x], π). An induction on the minimally defined models of DB_L shows that precsubmit and not \circ iscand as well as precwithdraw and iscand are not distinguishable by contexts of primitive sort. Thus

precsubmit = not • iscand

precwithdraw = iscand

holds in all minimally defined behaviour models M of DB_L. Furthermore nonterminating queries or updates are undefined, for example

 $M \models D(recquery(x, call(x))) = false$.

Note that a database specification is implemented by its database language.

Proposition

For every database specification DB^{1} the type $DB_{L}(DB^{1})$ is an implementation of DB^{1} via the signature morphism κ , for example induced by

$$\kappa(f(db,o_1,\ldots,o_k)) = io(apply(\overline{f}, >)$$

for all holding queries f of DB.

The following example shows how to use the functional database language.

Example

In DB_L \equiv DB_L(EMPLOYEES, EMPLOYEE, ID) the query

works (cf. section 3.2.4) can be expressed by recquery(w, body) ,

where

body = if eisempty
 then false
 else if eq o [sel(1), secr o echoose]
 then true
 else call(w) o efire o secr o echoose
 fi

Then in DB_L it is provable that

where works is the operation specified in section 3.2.4.

Proposition

The application of recquery does not change the state of the database, that is

DB \models apply(recquery(f,q), <db₁, i₁>) = <db₂, i₂>

$$db_1 = db_2$$

From the full database languages an update language and a query language can be derived.

Definition

For a database language $DB_L \equiv DB_L (DB, ID, SEQU(OBJ))$ a query language Q_L may be given by

$$\begin{bmatrix} data-enrich DB_L & by & (v \in \{query, upd\}) \\ \\ \underline{funct}(v, conf) & \underline{seq}(obj) & qapply, \\ \\ qapply(v, \pi) &= io(apply(v, \pi))]/\\ \hline \\ \hline \\ OUERY \end{bmatrix}$$

endoftype

where

$$\overline{QUERY} = \{ recquery(f, v) | f \in \underline{id} \quad x \ v \in W(DB_L) \}$$
$$U \{ \underline{funct} (\underline{db}, \underline{conf}) \ \underline{sequ(obj)} \ qapply \} .$$

An update language UPD_L may be given by

type UPD_L ≡

$$\frac{[\text{data-enrich DB}_L \text{ by}}{\text{funct } (v, \text{ conf}) \text{ db updapply }, (v \in \{\text{query, upd}\})}$$
$$\frac{(v \in \{\text{query, upd}\})}{(v \neq \pi)} = \frac{(v \neq \pi)}{(v \neq \pi)}$$

endoftype

where

UPDATE = sig(DB_L) \cup {updapply } {apply } .

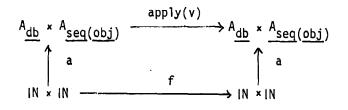
Thus programs of the query language do not change the state of the database, whereas programs of the update language do.

9. RECURSION COMPLETENESS

It has been advocated (/Chandra, Harel 80/) that the language to manipulate the database should have full computational power, i.e. that every partial recursive function can be computed.

Definition

A sig(DB_L)-algebra A is called recursion complete, if there exists a surjective recursive mapping a : $|N \times |N \rightarrow A_{db} \times A_{obj}$ such that for every partial recursive function f : $|N \times |N \rightarrow |N \times |N$ there exists $v \in W(DB)$ of sort <u>query</u> or <u>update</u> such that $A \models apply(v, a(n, m)) = a(f(n, m))$, that is the diagram



commutes. A is called query resp. update recursion complete, if analogously for every partial recursive $f : |N \times |N \rightarrow |N$ there exists $q \in W(DB)$ of sort query such that

 $A \models qapply(q, a(n, m)) = a(f(n, m))$ resp. an $u \in W(DB)$ of sort <u>upd</u> such that $A \models updapply(u, a(n, m)) = a(f(n, m))$

holds.

Then for every database DB^1 , $DB \perp (DB^1)$ is called query resp. update recursion complete if all minimally defined behaviour models of $DB \perp (DB^1)$ are query resp. update recursion complete.

Proposition

DB L(FINSET) is query and update recursion complete.

DB L is weakly sufficiently complete. Therefore the semantic operation apply - and thus also qapply - is a partial recursive function in every minimally defined model of Q L(FINSET). Conversely the recursion completeness implies that every partial recursive relation can be simulated using qapply by

$$qapply(q) : \underline{obj} \times \dots \times \underline{obj} \rightarrow \{a(0), a(1)\}.$$

Both properties together show that Q L(FINSET) is complete in the sense of /Chandra, Harel 80/.

Thus our notion of recursion completeness generalizes the completeness of Chandra, Harel ; it also generalizes their "extended completeness" since FINSET is parameterized with arbitrary object types.

10. CONCLUDING REMARKS

Algebraic types seem to be an interesting tool for the specification and analysis of database schemes. Such a specification may be considered as a stepwise development process where requirement analysis, external, conceptual, and internal schemes mark significant levels:

The transition from the requirement analysis to an external scheme can be seen as a refinement process where in each step the requirements are precized. The conceptual scheme is a solution of an "abstract domain equation" where the views of the external scheme have to be integrated. Then a joint development of data structures and algorithms leads to an internal scheme. In this step, for example keys are introduced and recursion can be removed.

The formal correctness of these transformations can be uniformly defined by an algebraic notion of implementation and thus algebraic methods can be used to support the specification of databases.

On the other hand for a sound database specification at each level it is necessary to analyse the properties of the algebraic types. Such an analysis may comprise the structure and behaviour of admissible models, the completeness of the specification, the complexity of the operations and the redundancy of information. This gives a guideline for the further development and a feedback with the informal ideas in mind.

When working with and reasoning about database specification a specification language seems to be needed to express flexibly various operations with algebraic types. But it would be very uneconomic (and to difficult for non-experts) to write down arbitrarily complex axiom systems. Therefore complementary to a specification language elaborated database abstractions may build the "skeleton" of a database-specification for which standard solutions and high level implementations are available.

ACKNOWLEDGEMENT

We gratefully acknowledge valuable discussions with Prof. G. Ausiello, C. Batini, Prof. F.L. Bauer, M. Brodie, Prof. E. Neuhold, and A. Pettorossi. Thanks go to R. Hyerle for reading a draft.

This research has been partially sponsored by the Sonderforschungsbereich 49, Programmiertechnik, München, and the Consiglio Nazionale delle Ricerche, Roma.

/ALBANO et al. 81/ A. Albano, M.E. Occhiuto, R. Orsini: GALILEO: A conceptual language for database applications. A preliminary definition. CNR-PFI-DATAID Report No 13, November 1981 /AUSIELLO et al. 80/ G. Ausiello, C. Batini, M. Moscarini: Conceptual relations between databases transformed under join and projection. In: P. Dembinski (ed.): 10th Conf. on Mathematical Foundations of Computer Science. LNCS 88 (1980) /BAUER, WUSSNER 82/ F.L. Bauer, H. Wössner: Algorithmic language and program development. Berlin-Heidelberg-New York: Springer (1982) /BAUER et al. 81/ F.L. Bauer, M. Broy, W. Dosch, R. Gnatz, F. Geiselbrechtinger, W. Hesse, B. Krieg-Brückner, A. Laut, T.A. Matzner, B. Möller, H. Partsch, P. Pepper, K. Samelson, M. Wirsing, H. Wössner: Report on a wide spectrum language for program specification and development. Report TUM-18104, Institut für Informatik, Technische Universität München (1981) /BJØRNER 80/ D. Bjørner: Formalization of Data Base Models. In: D. Bjørner (ed.): Proc. Advanced Course on Abstract Software Specifications, LNCS 86 (1980) /BRODIE 81/ M.L. Brodie: Association: A data base abstraction for semantic modelling. In P.P. Chen (ed.): Entity-Relationship Approach to Information Modelling and Analysis, ER Institute, Los Angeles, Oct. 1981, p. 583-608 /BRODIE, ZILLES 81/ M.L. Brodie, A.N. Zilles (eds.): Proc. Workshop on Data Abstraction, Data-Bases and Conceptual Modelling. SIGPLAN Notices 16:1 (1981) /BROY, WIRSING 80/ M. Broy, M. Wirsing: Algebraic definition of a functional programming language and its se-mantic models. Report TUM-18008, Institut für Informatik, Technische Universität München, 1980. Also in RAIRO (to appear) /BUNEMANN, FRANKEL 80/ 0.P. Bunemann, R.E. Frankel: FQL - a functional query language. Proc. ACM SIGMOD, May 1979 /BURSTALL, GOGUEN 80/ R.M. Burstall, J.A. Goguen: The semantics of CLEAR: a specification language. In D. Bjørner (ed.): Proc. Advanced Course on Abstract Software Specification. LNCS 86 (1980)

/CASANOVA et al. 81/ M.A. Casanova, J.M.V. de Castilho, A.L. Furtado: Properties of conceptual and external database schemes. Technical Report DB108103, Depart. Informática, Rio de Janeiro, Brasil (1981) /CHANDRA, HAREL 80/ A.K. Chandra, D. Harel: Computable queries for relational data bases, Journal of Computer and System Sciences 21, 156-178 (1980) /CODD 70/ E.F. Codd: A relational model of data for large shared data banks. Comm. of the ACM 13:6, 377-387 (1970) /EHRIG, FEY 81/ H. Ehrig, W. Fey: Methodology for the specification of software systems: from formal requirements to algebraic design specifications. In W. Brauer (ed.): Proc. 11th GI-Jahrestagung, München.Informatik-Fachberichte 50, Springer-Verlag (1981) /EHRIG et al. 78/ H. Ehrig, H.-J. Kreowski, H. Weber: Algebraic specification schemes for data base systems. Proc. 4th Int. Conference on Very Large Data Bases, Berlin West (October 1978) /FAGIN 82/ R. Fagin: Armstrong Databases. Proc. 7th IBM Symposium on Mathematical Foundations of Computer Science, Japan (May 1982) /GALLAIRE 81/ H. Gallaire: Impacts of logic on data bases. Proc. 7th Int. Conf. on Very Large Data Bases, Cannes, France (September 1981) /HAMMER 76/ M.M. Hammer: Data abstractions for data bases. Proc. Conf. on Data: Abstraction, Definition and Structure. ACM FDT 8 (2) /HAMMER, BERKOWITZ 80/ M. Hammer, B. Berkowitz: DIAL: A programming language for data intensive applications, SIGMOD 1980 (HAMMER, McLEOD 75/ M.M. Hammer, D.J. McLeod: Semantic integrity in a relational database system. Proc. 1st Int. Conf. on Very Large Data Bases, Framingham, Massachusetts (Sept. 1975) /HUPBACH 81/ U. Hupbach: Modelling data bases by abstract data types. 4th Int. Datenbank Seminar, Schwerin, GDR (August 1981) /LOCKEMANN et al. 79/ P.C. Lockemann, H.C. Mayr, W.H. Weil, W.H. Wohlleber: Data abstraction for data base systems. ACM Transactions on Database Systems,

4:1, 60-75 (1979)

/MAKOWSKY 81/ J.A. Makowsky: Characterizing data base dependencies. In: S. Even, O. Kariv (eds.): Proc. 8th Int. Colloquium on Automata, Languages and Programming. LNCS 115 (1981) /NEUHOLD, OLNHOFF 81/ E.J. Neuhold, Th. Olnhoff: Building data base management systems through formal specifications. In: J. Diaz, I. Ramos (eds.): Formalization of programming concepts. LNCS 107(1981) /PAOLINI 81/ P. Paolini: Abstract data types and data bases. In /Brodie, Zilles 81/ /dos SANTOS et al. 80/ C.S. dos Santos, E.J. Neuhold, A.L. Furtado: A data type approach to the entity-relationship model. In: P.P. Chen (ed.): Int. Conf. of the Entity-Relationship Approach to Systems Analysis and Design, North-Holland Amsterdam 1980, 103-119 /SCHIEL et al. 82/ W. Schiel, A.L. Eurtado, E.J. Neuhold: Towards multi-level and modular conceptual schema specifications. Report 2/82, Institut für Informatik, Universität Stuttgart /SCHMIDT 80/ J.W. Schmidt: Data Type Concepts for Databases. In: M. Atkinson (ed.): Data Design. Infotech State of the Art Report, Series 8, No 4 (1980) /SMITH, SMITH 77/ J.M. Smith, D.C.P. Smith: Database Abstractions: Aggregation and Generalization: ACM Transactions on Database Systems 2:2(1977) /VELOSO et al. 81/ P.A.S. Veloso, J.M.V. de Castilho, A.L. Furtado: Systematic derivation of complementary specifications. Proc. 7th Int. Conference on Very Large Data Bases, Cannes, France (September 1981) /WIRSING 82/ M. Wirsing: Structured Algebraic Specifications. Proc. AFCET Symposium on Mathematics for Computer Science, Paris, March 1982 /WIRSING et al. 80/ M. Wirsing, P. Pepper, H. Partsch, W. Dosch, M. Broy: On hierarchies of abstract data types. TUM-18007. Institut für Informatik, Technische Universität München. Also Acta Informatica (to

appear),