

ON A CERTAIN DISTRIBUTED DATA BASE MODEL

Józef Kubit

The Academy of Economics in Cracow
Department of Computer Science
Poland

ABSTRACT

The paper presents a formal model of distributed data base. Present work is an attempt at defining fundamental notion connected with decompositions of data base model. Distributed data base model consists of family of individual (local) data base models and administrator of distributed data base model. Model is based on a definition of relational schema of data. Data are the object of widely understood measurement.

1. INTRODUCTION

In recent years many different models have been presented describe the data base at the conceptual (mathematical, logical) level. The rapid development of large computer networks led to investigations about distributing data base throughout such a system. Looking at these trends we can identify two basic approaches [Neuhold 1977]:

- 1) Using already existing networks of large computer systems,
- 2) As an alternative, we may build a new data base management system which takes into account the existence of a computer network but does not rely on already existing individual (local) data base systems.

We envision for the next future, that computer networks commonly will encompass both large minicomputers even including intelligent terminals.

The reason for the interest in distributed data base systems is because they provide a solution to very real problems for the geographically distributed organization which needs to preserve a unified information-sharing and processing system.

* The paper has been presented on The International Colloquium COMPCONTROL '79 in Sopron on November, 1979.

In general, the advantage of centralized approach are the disadvantages of a distributed approach and vice versa.

They are, in broad form [Champine 1977]:

Distributed advantages / centralized disadvantages

- communication failsoft capability
- lower communications data rate and cost
- configuration flexibility
- high system performance (fast response and high transition rate)
- modular implementation
- modular up-grade

Generalized advantages / distributed disadvantages

- operations economy
- hardware economy of scale
- unified control
- easy update / retrieval
- compatibility.

Gradually along with a vigorous development of data base systems based on Codd's relational model of data, a new theory of relations has emerged. The birth of this theory must be assigned to Codd's papers [1970] and [1971] where the basic notions were introduced and collected together.

Depending on delimited purposes of realized work the notion of data base model is not used similarly. It concerns for example the works devoted theory of i.s.r systems [Marek and Pawlak 1976].

General comparison of Codd's and Pawlak's and Marek's approach was done by Koczko [1978] who showed that they lead to essentially same consequences.

Many data models have already been proposed, each having its own concepts and terminology. In studying the various models it becomes apparent that they have similarities and differences which are not trivial to analyze [Kerschberg, Klug and Tsichritsis 1976].

The notion of data base model is connected with the notion of information system model. Information system model in [Kubit 1979] consists of four specialized system models. Two of them constitute a certain data base model.

Information system model is component of controlled economic system model.

The purpose of the paper is a presentation of divided - distributed data base model.

In each of distributed data base cases we deal as follows:

The base B is decomposed into smaller entities $B_1, \dots, B_i, \dots, B_n$ and we add new object called administrator (\bar{A}). The administrator of queries. Some of i 's could be empty.

The response for i is the settheoretical union of the responses of B_i to i_i .

Of particular importance may be a formulation of relations occuring between notions of data base model and administrator of data base model.

2, NOTION OF DATA

Let's take an attribute names N . This set will be named alphabet of data names. From elements of this alphabet we'll make a word set N^* .

Let N be a subset of N^* set. N set will be named a set of data names. Data names will be marked by N while for all N there will be

$$N = \langle n_1, \dots, n_j, \dots, n_{r_N} \rangle \quad n_j \in N$$

Denote $\langle n_1, \dots, n_j, \dots, n_{r_N} \rangle$ will be also defined relational schema of data.

Let's take a nonempty family of domains $\{M_i\}_{i \in I}$ and mapping m .

The mapping m need not be one-to-one. This is onto mapping

$$m: N \rightarrow \{M_i\}_{i \in I}$$

According to what was said above for relational schema of data will correspond a finite relation

$$R_N = m(n_1) \times \dots \times m(n_j) \times \dots \times m(n_{r_N})$$

The set of data values corresponding relational schema N will be defined as $\gamma(N), \gamma(N)=P(R_N)$ where $P(R_N) \subseteq R_N$.

By V will be denoted an element of subset $P(R_N), V \in P(R_N)$. Element V will be named data value and

$$V = \langle v_1, \dots, v_j, \dots, v_{r_N} \rangle$$

DEFINITION 1

Data d will be named a pair composed of data name and data value $d = \langle N, V \rangle$.

Names n_j are the names of specified attributes, whereas parts v_j of data value V are the values of these attributes. Attributes which is defined by n_j includes the definition of scale and that of value unit v_j . Further on will be identified as equivalence class of data $d = \|d\|$.

Data will be obtained in the process of measurement in the broadest sense of the notion. Data measurement notion will be defined by means of data measurement function. Data measurement function will be referred to as isomorphism from real objects (portions of productive factors) algebra, onto data algebra

$$h: A \xrightarrow{\sim} D$$

where: A - universe of algebra of equivalence classes of real objects

$$A = A \langle 0_1, \dots, 0_n \rangle$$

D - universe of data algebra $D = \langle D, 0'_1, \dots, 0'_n \rangle$

(cf. [Kobayashi 1975])

Dated $d \in D$ are data equivalence classes because of defined equivalence relation. Algebras A and D are similar algebras. For these two algebras the following theorem holds true.

Theorem 1

Let f be one-to-one transformation of A_0 set of algebra generators onto D_0 set of similar algebra D generators. If there exists homomorphism h of algebra A into algebra D being an extension f and also homomorphism g of algebra D into algebra A being an extension f^{-1} , then h is isomorphism A onto D and $h^{-1}=g$.

3. SYSTEM OF DATA RETRIEVAL

We assume that I_N denotes a set of attribute names according to the relational schema of data (data name) N . By m_N will be understood restriction of m function to I_N set of attribute names

$$m_N = M \upharpoonright I_N$$

Let \bar{U} be a set of all functions m_N for set N of data names

$$\bar{U} = \{m_N\}_{N \in N}$$

Let I_N^* denote a set of all words over the alphabet I_N of attribute names of data N . The elements of I_N set will be marked as \bar{n} .

By $P(R_N)[\bar{n}]$ we shall define a set of all data value V restrictions to a set of attribute names \bar{n} according to a relational schema of data name N

$$PR_N[\bar{n}] = \{V \bar{n} : V \in P(R_N)\}$$

Let $(R_N)^*$ be understood as a set of all projections $P(R_N)$ for all $\bar{n} \in I_N^*$. By P we shall mark a set of all projections $P(R_N)^*$ for each data name N

$$= \{(R_N)^*\}_{N \in N}$$

We assume the θ_N^j will denote a set of two-argument order relations j -type values domain assigned to the data name N . Two-argument order relations θ_N^j give a set Θ of order relations of data name attribute values which is expressed in the following way

$$\Theta = \{\{\theta_N^j\}_{j \in \{1, \dots, r_N\}}\}_{N \in \mathcal{N}}$$

By Q we shall mark a set of set theory operations [Neuhold 1974] defined for P set.

DEFINITION 2

System of data retrieval for D set of data will be marked septuple

$$S_{DR} = \langle D, N, \bar{N}, \bar{U}, P, \Theta, Q \rangle$$

where:

- D - universe of data algebra
- N - alphabet of data names
- \bar{N} - set of data names
- \bar{U} - set of functions m_N (assigning attribute values domains) for data defined by relational schema from set
- P - set of all projections of data values from data set D
- Θ - set of two-argument order relations of attribute values
- Q - set of set theory operations for set of data values projections.

4. SYSTEM OF CONTROL SYSTEM SETS OF DATA RETRIEVAL

Let set N be alphabet of subset of data names for set D. From elements N of set N we compose of set of all words N^* . Let \hat{N} be a subset of N^* set. \hat{N} will be names a set of names of data subsets. The names of data subsets will be expressed by \hat{N}_k , where for each \hat{N}_k there will be

$$\hat{N}_k = \langle N_1, \dots, N_k, \dots, N_{w_N} \rangle ; \quad N_k \in N$$

The denote \hat{N} will be defined as a name of subset or relational schema of data subsets (sets).

We say that \hat{N} is a lattice with respect of the operation \cup (join) and \cap (meet) if the following equations hold (axioms of lattice) (cf. [Kuratowski and Mostowski 1976]).

$$\begin{array}{ll}
 (1) & \hat{N} \cup \hat{N} = \hat{N} & \hat{N} \cap \hat{N} = \hat{N} \\
 (2) & \hat{N} \cup \hat{N}' = \hat{N}' \cup \hat{N} & \hat{N} \cap \hat{N}' = \hat{N}' \cap \hat{N} \\
 (3) & \hat{N} \cup (\hat{N}' \cup \hat{N}'') = (\hat{N} \cup \hat{N}') \cup \hat{N}'' & \hat{N} \cap \hat{N}' \cap \hat{N}'' = \hat{N} \cap \hat{N}' \cap \hat{N}'' \\
 (4) & \hat{N} \cap (\hat{N} \cup \hat{N}') = \hat{N} & \hat{N} \cup (\hat{N} \cap \hat{N}') = \hat{N}
 \end{array}$$

We call a lattice distributive if

$$(5) \quad \hat{N} \cap (\hat{N}' \cup \hat{N}'') = (\hat{N} \cap \hat{N}') \cup (\hat{N} \cap \hat{N}'') \quad \hat{N} \cup (\hat{N}' \cap \hat{N}'') = (\hat{N} \cup \hat{N}') \cap (\hat{N} \cup \hat{N}'')$$

We introduce an order relation between elements of lattice:

$$(6) \quad \hat{N} \leq \hat{N}' \equiv \hat{N} \cup \hat{N}' = \hat{N}'$$

or, equivalently,

$$(7) \quad \hat{N} \leq \hat{N}' \equiv \hat{N} \cap \hat{N}' = \hat{N}$$

Similarly we define the elements o and i as the elements satisfying conditions

$$(8) \quad \hat{N} \cup o = \hat{N}, \quad \hat{N} \cap i = \hat{N}$$

It is easy to show o is the smallest element in \hat{N} and that i is the largest, namely for every $\hat{N} \in \hat{N}$

$$(9) \quad o \leq \hat{N} \leq i$$

Let's take a set I of control system queries. We assume that for every query correspond a finite nonempty set $\hat{N}_i \in \hat{N}$, such a way that

$$\bigcup_{i \in I} \hat{N}_i = \hat{N}$$

and for all $i, j \in I$, $i \neq j$ $\hat{N}_i \cap \hat{N}_j = \emptyset$

Instead of speaking about the partition $\{\hat{N}_i\}_{i \in I}$ we may consider the equivalence relation R_I on $\hat{N} \times \hat{N}$ that the family of its equivalence classes is indexed by the set I and

$$\hat{N} / R_I = \{\hat{N}_i\}_{i \in I}$$

DEFINITION 3

A system of control system sets of data retrieval is a sextuple

$$S_{\text{CSSDR}} = \langle D, \hat{N}, U, \Omega, R_I, U \rangle \quad (\text{cf. [Raš 1978]})$$

where:

- D - data set
- $\langle \hat{N}, U, \Omega \rangle$ - lattice of names of data subsets
- R_I - an equivalence relation in and I is set of control systems queries
- $U: \hat{N} \rightarrow 2^D$ - monotonic function satisfying the following conditions:

$$(\forall \hat{N}, \hat{N}' \in \hat{N} (\hat{N} \leq \hat{N}') = U(\hat{N}) \subseteq U(\hat{N}'))$$

$$(\forall d \in D) (\exists \hat{N} \in \hat{N} (d \in U(\hat{N})))$$

DEFINITION 4

Let S_{CSSDR} be a system of control system sets of data retrieval.

Let $D' \subseteq D$.

- a) D' is said to be describe within S_{CSSDR} iff there is $\hat{N} \in \hat{N}$ such that $U(\hat{N})=D'$.
- b) $B(S_{\text{CSSDR}})$ is the family of all subsets of D describable within S_{CSSDR}

DEFINITION 5

A system is selective iff for all $d \in D$, $U(\hat{N}_d) = \{d\}$

Theorem 2 (cf. [Lipski and Marek 1975])

A system S_{CSSDR} is selective iff $B(S_{\text{CSSDR}}) = 2^D$ (recall that we consider only the case when I (set of control system queries), and consequently D , are finite).

P r o o f:

If $B(S_{\text{CSSDR}}) = 2^D$ then obviously S_{CSSDR} is selective. If S_{CSSDR} is selective then all $d \in D$ have different descriptions, hence D is finite. For any $D' \subseteq D$ we have then

$$D' = U\left(\sum_{d \in D} \hat{N}_d\right),$$

i.e. D' is describable.

5, A DATA BASE MODEL,

DEFINITION 6

Data base model of productive object is a pair

$$DB = \langle S_{DR}, S_{CSSDR} \rangle$$

where:

- | | |
|-------------|---|
| S_{DR} | - system of data retrieval |
| S_{CSSDR} | - system of control system sets of data retrieval |

The following properties of data base model presented above can be distinguished:

1. Model is based on a definition of relational schema of data as proposed by Codd [1971].
2. Data are the object of widely understood measurement.
3. For the presented model i.s.r, system theory can be applied.
4. Identification of data base follows the identification of information system for productive object.
5. Presented data base model can be used in the construction of distributed data base model.

Remark

Presented data base model is a proposal resulting from necessity for complex analysis of problems connected with its modelling. It concerns primarily the contents of data base.

6. DISTRIBUTED DATA BASE MODEL AND ITS PROPERTIES

DEFINITION 7

System of storage and retrieval of documents (cf. [Marek and Pawlak 1976]) is quadruple

$$I = \langle X, B, R_J, V \rangle$$

where:

- X - set of documents
- B - set of descriptors
- R_J - equivalence relation on B of finite index
- V - maps B into 2^X ($V: B \rightarrow 2^X$) and satisfied the following two conditions:

- 1) if $a R_J b$ $a \neq b$, then $V(a) \cap V(b) = \emptyset$
- 2) $\cup \{V(b) : b R_J a\} = X$ for each $b \in B$.

DEFINITION 8

Let $S_{CSSDR} = \langle D, \hat{M}, U, \cap, R_I, U \rangle$ be system of control system sets of data retrieval and I be system of storage and retrieval of documents. System I coverage S_{CSSDR} IFF

- 1) there exists an one-to-one function Ψ such that $\Psi: J \rightarrow I$
- 2) there exists a function

$$\phi : \bigcup_{j \in J} \hat{N}_{\Psi(j)} \rightarrow B \text{ such that}$$

$$V(b) = \cup \{U(\hat{N}) : \hat{N} \in \phi^{-1*} \{b\}\} \quad \text{and}$$

$$\hat{N} \in \hat{N}_{\Psi(j)} \Rightarrow \phi(\hat{N}) \in B_j$$

The intuition which is connected with this following: I is a "presystem" classifying the names of data subsets in according to some documents from set of documents X.

DEFINITION 9

Let S_{CSSDR} be system of control system sets of data retrieval and I be system of storage and retrieval of documents. We say that I strongly covers S_{CSSDR} iff I covers S_{CSSDR} and

$$1) (\forall b)_B \text{ card}(\phi^{-1} * \{b\}) \geq 2$$

and

$$2) \forall i \in I \text{ card}(\hat{N}_i) \leq 2$$

Let K be set of indexes of existing individual (local) systems of control system sets of data retrieval.

DEFINITION 10

Let S_{CSSDR} be system of control system sets of data retrieval. Let $\{I_k\}_{k \in K}$ be a partition of the set I . An induced family $\{S_{CSSDR}_k\}_{k \in K}$ of individual (local) systems of control system sets of data retrieval is formed as follows:

$$S_{CSSDR}_k = \langle D_k, \hat{N}^{(k)}, U_k, \cap_k, R_{I_k}, U_k \rangle$$

$$1) N^{(k)} = \bigcup_{i \in I_k} \hat{N}_i$$

$$2) R_{I_k} = R_I \cap (\hat{N}^{(k)} \times \hat{N}^{(k)})$$

$$3) U_k = U \uparrow \hat{N}^{(k)}; \quad \cap_k = \uparrow \hat{N}^{(k)}$$

$$4) U_k = U \uparrow \hat{N}^{(k)}$$

$$5) D_k = \hat{N}^{(k)} \in \hat{N}^k U_k(\hat{N}^{(k)})$$

By S_{DR_k} will be understood restriction of system of data retrieval S_{DR} to set of data D_k

$$S_{DR_k} = \langle D_k, N, N^{(k)}, U_k, P_k, \Theta_k, Q_k \rangle$$

where:

- D_k - set of data
- N - alphabet of data names
- $N^{(k)}$ - set of names of data set D_k
- U_k - set of functions m_N (assigning attribute values domains) for data defined by relational schema from $N^{(k)}$ set
- P_k - set of all projections of data values for data set D_k
- Θ_k - set of all two-argument order relations of data name attribute values for $N^{(k)}$ set of data names
- Q_k - set of set theory operations for P_k set of data values projections.

DEFINITION 11

Individual (local) data base model DB_k is a pair

$$DB_k = \langle S_{DR_k}, S_{CSSDR_k} \rangle$$

- where:
- S_{DR_k} - individual (local) system of data retrieval
 - S_{CSSDR_k} - individual (local) system of control system sets of data retrieval

DEFINITION 12

Divided data base model is a quadruple

$$DDB = \langle \{DB_k\}_{k \in K}, I, \Psi, \phi \rangle$$

where: $\{DB_k\}_{k \in K}$ - family of individual (local) data base models
 I - system of storage and retrieval of documents
(covering system)
 Ψ, ϕ - covering functions

Administrator of divided data base model will be named triple

$$\bar{A} = \langle I, \Psi, \phi \rangle .$$

Final remark

Presented divided data base modes is a proposal resulting from necessity for analysis of very real problems for the geographically distributed organization of large data bases for an enterprise.

7. LITERATURE

1. Champine, G.A.: "Six approaches to distributed data bases". Datamation 5(1977) pp. 69-72.
2. Codd, E.F.: "A relational model of data for large shared data banks". Comm. ACM, Vol.13, No.6, June 1970, pp.377-387.
3. Codd, E.F.: " Further normalization of the data base relational model", Couran Comp. Sc. Symp. 6: Data base systems, Prentice-Hall, N.J., May 1971, pp. 65-98.
4. Kerschberg, L. - Klug, A. - Tsichritzis, D.: "A taxonomy of data models"; System for Large Data Bases; Lockemann, P.C.-Neuhold, E.J. (eds.), North-Holland Publishing Company, 1976.

5. Kobayashi, I.: "Information and information processing structure"; Information Systems 1(1975), pp.39-49.
6. Koczkodaj, W.W.: "A relational model of data and its connections with the i.s.r. systems"; ICS PAS Reports 306. Warsaw 1978.
7. Kubit, J.: "On a certain information system model"; 5th Symposium of algorithms-ALGORITHMS'79, April 1979, Strebske Pleso, pp.309-315.
8. Kuratowski, K. - Mostowski, A.: "Set theory" ; North-Holland Publishing Company, Polish Scientific Publishers, Warsaw 1976.
9. Lipski, W. - Marek, W.: "On information storage and retrieval systems"; CC PAS Report 200, Warsaw 1975.
10. Marek, W. - Pawlak, Z.: "Information storage and retrieval systems - Mathematical foundations"; Theoretical Computer Science 2(1976), pp. 331-354.
11. Marek, W. - Rode - Babczenko, I.: "A decomposition of informational systems"; CC PAS Reports 212, Warsaw 1975.
12. Neuhold, E.J.: "Formal properties of data bases"; University of Stuttgart, Mathematical Centre Tracts 63, 1974, pp.121-177.
13. Neuhold, E.J. - Biller, H.: "POREL: A distributed data base on an inhomogeneous computer network"; Third International Conference on Very Large Data Bases; Tokyo, Japan, October 1977, pp. 380-389.
14. Raś, Z.: "Algebraic foundations of information storage and retrieval system II"; (in Polish), ICP PAS Reports 339, Warsaw 1978.

ÖSSZEFOGLALÁS

Egy osztott adatbázis modell

József Kubit

Osztott adatbázisok egy formális modelljét mutatjuk be, a felbontásaival kapcsolatos alapvető jelöléseket definiáljuk. Lokális adatbázis modellek egy családját és az osztott adatbázis modell adminisztrátorát értjük osztott adatbázis modellel. A modell a relációs adatséma egy definícióján alapul.

Об одной модели распределенных баз данных

Йожеф Кубит

В данной работе представляется некоторая формальная модель распределенных баз данных, и предлагаются определения основных понятий связанных с ее декомпозицией. Модель распределенной базы данных состоит из администратора и семейства локальных моделей. Предложенная модель основана на некотором определении реляционной схемы данных.