

Computer Science
and Applied Mathematics

**DATA BASE ORGANIZATION FOR
DATA MANAGEMENT**

Sakti P. Ghosh

BIBLIOTHEQUE DU CERIST

Data Base Organization for Data Management

SAKTI P. GHOSH

*IBM Research Laboratory
San Jose, California*



ACADEMIC PRESS New York San Francisco London 1977
A Subsidiary of Harcourt Brace Jovanovich, Publishers

BIBLIOTHEQUE DU CERIST

Computer Science and Applied Mathematics

A SERIES OF MONOGRAPHS AND TEXTBOOKS

Editor

Werner Rheinboldt

University of Maryland

HANS P. KÜNZI, H. G. TZSCHACH, and C. A. ZEHNDER. Numerical Methods of Mathematical Optimization: With ALGOL and FORTRAN Programs, Corrected and Augmented Edition

AZRIEL ROSENFELD. Picture Processing by Computer

JAMES ORTEGA AND WERNER RHEINBOLDT. Iterative Solution of Nonlinear Equations in Several Variables

AZARIA PAZ. Introduction to Probabilistic Automata

DAVID YOUNG. Iterative Solution of Large Linear Systems

ANN YASUHARA. Recursive Function Theory and Logic

JAMES M. ORTEGA. Numerical Analysis: A Second Course

G. W. STEWART. Introduction to Matrix Computations

CHIN-LIANG CHANG AND RICHARD CHAR-TUNG LEE. Symbolic Logic and Mechanical Theorem Proving

C. C. GOTTLIEB AND A. BORODIN. Social Issues in Computing

ERWIN ENGELER. Introduction to the Theory of Computation

F. W. J. OLVER. Asymptotics and Special Functions

DIONYSIOS C. TSICHRITZIS AND PHILIP A. BERNSTEIN. Operating Systems

ROBERT R. KORFHAGE. Discrete Computational Structures

PHILIP J. DAVIS AND PHILIP RABINOWITZ. Methods of Numerical Integration

A. T. BERZTISS. Data Structures: Theory and Practice, Second Edition

N. CHRISTOPHIDES. Graph Theory: An Algorithmic Approach

ALBERT NIJENHUIS AND HERBERT S. WILF. Combinatorial Algorithms

AZRIEL ROSENFELD AND AVINASH C. KAK. Digital Picture Processing

SAKTI P. GHOSH. Data Base Organization for Data Management

DIONYSIOS C. TSICHRITZIS AND FREDERICK H. LOCHOVSKY. Data Base Management Systems

Contents

<i>Preface</i>	ix
<i>Acknowledgments</i>	xi
1 Data Structures	
1.0 INTRODUCTION	1
1.1 BASIC MATHEMATICAL CONCEPTS	3
1.2 FINITE GEOMETRY	16
1.3 PRIMITIVES	29
1.4 LOGICAL RELATION	33
1.5 ENTITY SET MODEL	40
1.6 RELATIONAL MODEL	44
1.7 GRAPH STRUCTURE MODEL	49
EXERCISES	58
REFERENCES	60
2 Queries and Query Languages	
2.0 INTRODUCTION	62
2.1 PREDICATES	63
2.2 SIMPLE QUERIES	68
2.3 QUERIES BASED ON HIERARCHICAL STRUCTURES	74
2.4 QUERIES BASED ON RELATIONAL ALGEBRA	77
2.5 QUERIES WITH LOGICAL STRUCTURES	80
2.6 SET REPRESENTATION QUERY LANGUAGE	85
EXERCISES	92
REFERENCES	93

3 Searching on One Field

3.0	INTRODUCTION	94
3.1	SERIAL SEARCH	96
3.2	SORTING	99
3.3	BINARY TREE SEARCH	103
3.4	INDEX SEQUENTIAL SEARCH	111
3.5	OVERFLOW MANAGEMENT	122
	EXERCISES	126
	REFERENCES	127

4 Key to Address Transformation

4.0	INTRODUCTION	128
4.1	RANDOM TRANSFORMATION	132
4.2	DIVISION TRANSFORMATION	135
4.3	RADIX TRANSFORMATION	142
4.4	POLYNOMIAL TRANSFORMATION	144
4.5	RANK TRANSFORMATION	147
4.6	OTHER TRANSFORMATIONS	151
4.7	OVERFLOW HANDLING	156
4.8	COMPARISON OF DIFFERENT KATS	158
	EXERCISES	165
	REFERENCES	166

5 Algebraic Filing Schemes

5.0	INTRODUCTION	168
5.1	TIME-SPACE TRADE-OFF	171
5.2	BALANCED FILING SCHEMES FOR BINARY ATTRIBUTES	173
5.3	BALANCED FILING SCHEMES FOR MULTIPLE-VALUED ATTRIBUTES	189
5.4	ASYMMETRIC FILING SCHEMES	202
	EXERCISES	210
	REFERENCES	210

6 Consecutive Retrieval Property

6.0	INTRODUCTION	212
6.1	C-R PROPERTY FOR BINARY ATTRIBUTES	220
6.2	C-R PROPERTY FOR MULTIPLE-VALUED ATTRIBUTES	233
6.3	GRAPH THEORETIC APPROACH TO C-R PROPERTY	244
6.4	THE C-R PROPERTY WITH REDUNDANCY	251
	EXERCISES	266
	REFERENCES	267

7 Organization on Drum Storage

7.0	INTRODUCTION	269
7.1	BINARY SEARCH ON A DRUM	273
7.2	CONSECUTIVE STORAGE ON DRUMS	289
7.3	DRUM SCHEDULING	307
	EXERCISES	314
	REFERENCES	316

8 Access Path Retrieval

8.0	INTRODUCTION	317
8.1	SIMPLE ACCESS PATHS	319
8.2	MULTIPLE ACCESS PATHS	326
8.3	GENERAL ACCESS PATHS	328
8.4	SEARCH PATH ALGORITHMS	352
	EXERCISES	368
	REFERENCES	369

<i>Author Index</i>	371
<i>Subject Index</i>	373

Preface

In the past few years the area of data management has become an extremely important one in computer science. Large volumes of non-numerical data are stored in large data banks and are processed by complex queries from time to time as the need arises. The techniques that are used for numerical information processing cannot be easily applied to nonnumerical information because of the complex logical structures inherent in them. Most of the techniques for handling nonnumerical information processing have been developed by practitioners, and very few formal descriptions are available in written form. Researchers associated with prestigious universities and large computer manufacturers have been doing basic research to understand these techniques at a fundamental level. Their works are scattered in many journals, technical reports, proceedings of meetings, etc. This presents some very difficult problems for an instructor putting together a course in fundamentals of data management organization. The problem is more difficult for students. I have lost count of the number of times I have been asked by a student for reprints and references and by teachers to give seminars and lectures in their classes. There is no textbook on the market dealing with the theories of data base organization and usable at a graduate level. There is no book on the market covering such basic concepts as theories of data description models, logical structures of queries, combinational query sets, quadratic residue transformations, balanced filing schemes, the consecutive retrieval property, and organization on drums. These topics can be found in journal articles, which are on a level that good researchers can understand. They are not written for students. In view of the great importance of this area of knowledge in computer science, I felt that this vacuum should not be allowed to remain any more.

The textbooks that have been written on data management are directed toward undergraduate students at elementary or senior level and concentrate on either algorithmic aspects or system aspects of data management. In this book I have tried to put the basic theories and techniques of data management in the foreground, ignoring the systems details of information management. Chapters 1–4 are directed toward senior-level undergraduate students and Chapters 5–8 are directed toward graduate students.

The eight chapters of the book cover the following subjects: (1) data structure, (2) queries and query languages, (3) searching on one field, (4) key to address transformation, (5) algebraic filing schemes, (6) consecutive retrieval property, (7) organization on drum storage, and (8) access path retrieval.

Chapter 1 covers some primitive concepts of data description, basic mathematical ideas relevant to the materials covered in the book, and some data description models, namely, entity, relational, and graph theoretic models.

Chapter 2 covers various types of queries, their parametric representation, simple queries, complex queries, logical structured queries, queries based on a relational algebra set representation query language, etc.

Chapter 3 contains techniques of searching on one field; namely, serial search, sorting, binary tree search, index sequential search, and hierarchical access methods.

Chapter 4 deals with techniques used in key to address transformation, the overflow problems, and comparison of different key to address transformation methods.

Chapter 5 contains filing schemes for answering queries based on multiple values of multiple attributes. In this chapter properties of finite geometries are extensively used for construction of filing schemes.

Chapter 6 deals with file organizations that need no redundant storage and at the same time have minimum access time.

Chapter 7 discusses techniques of organizing records on drum storage; namely, binary search on drums, consecutive storage on drums, scheduling on drums, etc.

Chapter 8 deals with types of access paths used for retrieval of data and definition of logical structures. It contains linear lists, logical access paths, search path algorithms, etc.