



MORGAN KAUFMANN PUBLISHERS, INC.

MANAGING GIGABYTES

COMPRESSING AND INDEXING DOCUMENTS AND IMAGES

Ian H. Witten, Alistair Moffat, & Timothy C. Bell

Second edition; Spring 1999; 500 pages; cloth; 1-55860-570-3; \$59.95

As today's information explosion generates greater and greater volumes of raw data, the challenge of storing and retrieving this information in the most efficient manner continues to grow, whether the data is stored on a local disk or distributed over the World-Wide Web.

Managing Gigabytes helps you to meet this challenge by showing how to capitalize on new methods of compressing and accessing data, enabling you to store information more efficiently and locate specific items more quickly and cost-effectively than ever before. It uniquely covers fully-tested techniques for both text and image compression and shows how to construct a tailor-made electronic index for accessing text, scanned documents, and images.

This book largely avoids extensive theoretical and mathematical discussions, making it accessible to curious laypersons who seek a clear, uncomplicated understanding of this new technology. Real, large-scale problems are illustrated, and the technical material is sprinkled with anecdotes and background information. The new edition is updated with information about recent standards and discoveries. An instructor's supplement will be available.

Managing Gigabytes provides current and comprehensive tools and techniques that will help professionals and academics to work more confidently and effectively in today's increasingly paperless society.

NEW FEATURES

All chapters have been thoroughly updated and new references added to the recent literature;

Up-to-the-minute information on new text compression algorithms is included such as block sorting, approximate arithmetic coding, and fast Huffman coding;

New sections on context-based index compression and on distributed querying, and two new data structures for fast indexing;

Significant new material on image coding like descriptions of de facto standards that are in widespread use on the Web (GIF and PNG), information on CALIC, a high-performance lossless scheme, the new proposed "JPEG Lossless" international standard, and the new JBIG2 international standard for compressing document images;

New information on the Internet and World-Wide Web, digital libraries, Web search engines, and agent-based information retrieval;

New Appendix describing the New Zealand Digital Library, a large-scale web-based application incorporating many of the key ideas presented in the book;

An Instructor's Supplement will be available that contains review and test questions.

TABLE OF CONTENTS

PREFACE

1. OVERVIEW

- 1.1 Document databases
- 1.2 Compression
- 1.3 Indexes
- 1.4 Document images
- 1.5 The MG system
- 1.6 Further reading

2. TEXT COMPRESSION

- 2.1 Models
- 2.2 Adaptive models
- 2.3 Huffman Coding
- 2.4 Arithmetic coding
- 2.5 Symbolwise models
- 2.6 Dictionary models
- 2.7 Synchronization
- 2.8 Performance comparisons
- 2.9 Further reading

3. INDEXING

- 3.1 Sample document collections
- 3.2 Inverted file indexing
- 3.3 Inverted file compression
- 3.4 Performance of index compression methods
- 3.5 Signature files and bitmaps
- 3.6 Comparison of indexing methods
- 3.7 Case folding, stemming, and stop words
- 3.8 Further reading

4. QUERYING

- 4.1 Accessing the lexicon
- 4.2 Partially specified query terms
- 4.3 Boolean query processing
- 4.4 Ranking and information retrieval
- 4.5 Evaluating retrieval effectiveness
- 4.6 Implementation of the cosine measure
- 4.7 Interactive retrieval
- 4.8 Distributed retrieval
- 4.9 Further reading

5. INDEX CONSTRUCTION

- 5.1 Memory-based inversion
- 5.2 Sort-based inversion
- 5.3 Exploiting index compression
- 5.4 Compressed in-memory inversion
- 5.5 Comparison of inversion methods
- 5.6 Constructing signature files and bitmaps
- 5.7 Dynamic collections
- 5.8 Further reading

6. IMAGE COMPRESSION

- 6.1 Types of image
- 6.2 The CCITT fax standard for bilevel images
- 6.3 Context-based compression of bi-level images
- 6.4 JBIG: A standard for bilevel images
- 6.5 Lossless compression of continuous-tone images
- 6.6 JPEG: A standard for continuous-tone images
- 6.7 Progressive transmission of images
- 6.8 Summary of image compression techniques
- 6.9 Further reading

7. TEXTUAL IMAGES

- 7.1 The idea of textual image compression
- 7.2 Lossy and lossless compression
- 7.3 Extracting marks
- 7.4 Template matching
- 7.5 From marks to symbols
- 7.6 Coding the components of a textual image
- 7.7 Performance: lossy and lossless modes
- 7.8 System considerations
- 7.9 JBIG2: A standard for textual image compression
- 7.10 Further reading

8. MIXED TEXT AND IMAGES

- 8.1 Orientation
- 8.2 Segmentation

- 8.3 Classification
- 8.4 Further reading

9. IMPLEMENTATION

- 9.1 Text compression
- 9.2 Text compression performance
- 9.3 Images and textual images
- 9.4 Index construction
- 9.5 Index compression
- 9.6 Query processing
- 9.7 Further reading

10. THE INFORMATION EXPLOSION

- 10.1 Two millennia of information
- 10.2 The Internet: a global information resource
- 10.3 The paper problem
- 10.4 Coping with the information explosion
- 10.5 Digital libraries
- 10.6 Managing gigabytes better
- 10.7 Small is beautiful
- 10.8 Personal information support for life
- 10.9 Further reading

A. GUIDE TO THE MG SYSTEM

- A.1 Installing the MG system
- A.2 A sample storage and retrieval session
- A.3 Database creation
- A.4 Querying an indexed document collection
- A.5 Nontextual files
- A.6 Image compression programs

B. GUIDE TO THE NZDL

- B.1 What's in the NZDL?
- B.2 How the NZDL works
- B.3 Implications
- B.4 Further reading

REFERENCES

INDEX

ABOUT THE AUTHORS

IAN H. WITTEN is Professor of Computer Science at the University of Waikato in New Zealand. He is a Fellow of the ACM and of the Royal Society of New Zealand, and a member of professional computing, information retrieval, and engineering associations in the UK, USA, Canada, and New Zealand. He is co-author of *The Reactive Keyboard* (1992) and *Text Compression* (1990), as well as many journal articles and conference papers.

ALISTAIR MOFFAT is Associate Professor of Computer Science at the University of Melbourne. He is the author of numerous peer-reviewed papers, which have explored such areas as algorithms and data structures for text and image compression, self-adjusting data structures for dictionaries and priority queues, and algorithms for adaptive searching and sorting.

TIMOTHY C. BELL is a Senior Lecturer in Computer Science at the University of Canterbury. He is co-author of *Text Compression* (1990) and the author of a number of papers covering text and image compression, as well as computers and music, and computer education.



Order from Morgan Kaufmann Publishers

Mail: Harcourt Brace & Company Australia, Locked Bag 16
Marrickville, NSW 2204, Australia
Phone: 61-2-9517-8999, **Fax:** 61-2-9517-2249

Mail: Harcourt Brace & Co., Attn. Order Fulfillment Dept.,
6277 Sea Harbor Drive, Orlando, FL 32887,
Phone: 800-745-7323, 407-345-3800, **Fax:** 800-874-6418, 407-345-4060,
Email: orders@mkp.com, **Visit Morgan Kaufmann on the Web:** www.mkp.com