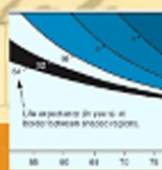Michael **LESK**

# Understanding
# Digital Libraries

**SECOND EDITION**

# Understanding
# Digital Libraries

# The Morgan Kaufmann Series in Multimedia Information and Systems

*Series Editor*, Edward A. Fox, Virginia Polytechnic University

*Understanding Digital Libraries, Second Edition*
Michael Lesk

*Bioinformatics: Managing Scientific Data*
Zoe Lacroix and Terence Critchlow

*How to Build a Digital Library*
Ian H. Witten and David Bainbridge

*Digital Watermarking*
Ingemar J. Cox, Matthew L. Miller, and Jeffrey A. Bloom

*Readings in Multimedia Computing and Networking*
Edited by Kevin Jeffay and HongJiang Zhang

*Introduction to Data Compression, Second Edition*
Khalid Sayood

*Multimedia Servers: Applications, Environments, and Design*
Dinkar Sitaram and Asit Dan

*Managing Gigabytes: Compressing and Indexing Documents and Images,*
*Second Edition*
Ian H. Witten, Alistair Moffat, and Timothy C. Bell

*Digital Compression for Multimedia: Principles and Standards*
Jerry D. Gibson, Toby Berger, Tom Lookabaugh, Dave Lindbergh, and Richard L. Baker

*Readings in Information Retrieval*
Edited by Karen Sparck Jones and Peter Willett

# Understanding
# Digital Libraries

## Second Edition

Michael Lesk

Again, this book is dedicated to the late Gerard Salton (1927–1995), a great information retrieval researcher who introduced many algorithms in the 1960s that are only now becoming commercially recognized, whose students have taught and studied in many universities, and who in 1961 first introduced me to computers, programming, and information retrieval.

This Page Intentionally Left Blank

# Contents

This Page Intentionally Left Blank

# Figures

This Page Intentionally Left Blank

# Tables

# Figure Credits

Figure 1.2    Based on data from Leach, S. 1986. "The Growth Rate of Major Academic Libraries: Rider and Purdue Reviewed." College and Research Libraries 37 (Nov) and from Rider, F. 1944. *The Scholar and the Future of the Research Library.* New York: Hadham Press.

Figure 1.3    Based on data from Meadows, J (1993) "Too Much of a Good Thing?" in H. Woodward and S. Pilling, eds. *The International Serials Industry*, Aldershot, Hampsire: Gower Publishing.

Figure 1.5    Courtesy of Mrs. Joan Adria D'Amico, the executor of the estate of Mary Crimi, wife of Alfred D. Crimi.

Figure 1.6    Courtesy of the Public Records Office, London FD 850/234.

Figure 1.7    Based on data from several sources including "Birth of a Chip." 1996. *Byte* (Dec) and "The Personal Computing Industry by the Numbers." 1997 *PC Magazine* (March).

Figure 1.8    Based on data from several sources including Pugh, E. 1995. *Building IBM: Shaping an Industry and Its Technology.* Cambridge, MA: MIT Press.

Figure 1.9    Adapted from Hammann, D. 1983. "Computers in Physics: An Overview." *Physics Today* 36(5).

Figure 2.2    Courtesy of the *London Daily Mail.*

| | |
|---|---|
| Figure 2.3 | Courtesy of the *American Chemical Society*. |
| Figure 3.1 | Courtesy Ann B. Lesk. |
| Figure 3.2 | 'Flying' scanner by Bob Kobres, University of Georgia Libraries. Reprinted with permission. |
| Figure 3.3 | (a) 4DigitalBooks.com, (b) www.kirtas-tech.com. Reprinted with permission. |
| Figure 3.7 | From DjVuZone.com, Leon Bottou. Reprinted with permission. |
| Figure 3.8 | The World. *Burney London Daily Journal.* 8.3.1728, PB.MIC.BUR.823B. Used by permission of the British Library. |
| Figure 3.9 | Screenshot of heads. Used by permission of the British Library. |
| Figure 3.10 | Screenshot of heads. Used by permission of the British Library. |
| Figure 3.11 | *The Journal of Physical Chemistry*, Vol.92 No.26, p.7162 (Fig.2). Copyright 1988 American Chemical Society. Reprinted with permission. |
| Figure 3.13 | *The Journal of Analytical Chemistry*, 63(17), 1697-1702. © Copyright 1991 American Chemical Society. Reprinted with permission. |
| Figure 3.14 | *The Journal of Organic Chemistry*, Vol. 56, No. 26, p.7270 (Fig.1) © Copyright 1991 American Chemical Society. Reprinted with permission. |
| Figure 3.15 | (a) Courtesy of the Harvard University Library, (b) Courtesy of the French National Library, Francois Mitterand, Dominique Perrault, Architect, Alain Goustard, Photographer. |
| Figure 4.2 | (a) Courtesy of Anne B. Lesk. |
| Figure 4.4 | Courtesy of www.bberger.net/rwb/gamma.html. |
| Figure 4.5 | Courtesy of EE/CS University of California, Berkeley. |
| Figure 4.6 | Courtesy of EE/CS University of California, Berkeley. |
| Figure 4.7 | "Automatic Image Annotation and Retrieval using CrossMedia Relevance Models by Jeon, Lavrenko and Manmatha," *SIGIR'03*. Reprinted with permission. |
| Figure 4.8 | Courtesy of the Department of Engineering Science, University of Oxford. |
| Figure 4.9 | Courtesy of Department of Engineering Science, University of Oxford. |