## Data Management for Multimedia Retrieval

Multimedia data require specialized management techniques because the representations of color, time, semantic concepts, and other underlying information can be drastically different from one another. The user's subjective judgment can also have significant impact on what data or features are relevant in a given context. These factors affect both the performance of the retrieval algorithms and their effectiveness. This textbook on multimedia data management techniques offers a unified perspective on retrieval efficiency and effectiveness. It provides a comprehensive treatment, from basic to advanced concepts, that will be useful to readers of different levels, from advanced undergraduate and graduate students to researchers and professionals.

After introducing models for multimedia data (images, video, audio, text, and web) and for their features, such as color, texture, shape, and time, the book presents data structures and algorithms that help store, index, cluster, classify, and access common data representations. The authors also introduce techniques, such as relevance feedback and collaborative filtering, for bridging the "semantic gap" and present the applications of these to emerging topics, including web and social networking.

K. Selçuk Candan is a Professor of Computer Science and Engineering at Arizona State University. He received his Ph.D. in 1997 from the University of Maryland at College Park. Candan has authored more than 140 conference and journal articles, 9 patents, and many book chapters and, among his other scientific positions, has served as program chair for ACM Multimedia Conference'08, the International Conference on Image and Video Retrieval (CIVR'10), and as an organizing committee member for ACM SIG Management of Data Conference (SIGMOD'06). In 2011, he will serve as a general chair for the ACM Multimedia Conference. Since 2005, he has also been serving as an associate editor for the *International Journal on Very Large Data Bases* (*VLDB*).

Maria Luisa Sapino is a Professor in the Department of Computer Science at the University of Torino, where she also earned her Ph.D. There she leads the multimedia and heterogeneous data management group. Her scientific contributions include more than 60 conference and journal papers; her services as chair, organizer, and program committee member in major conferences and workshops on multimedia; and her collaborations with industrial research labs, including the RAI-Crit (Center for Research and Technological Innovation) and Telecom Italia Lab, on multimedia technologies.

# DATA MANAGEMENT FOR

# MULTIMEDIA RETRIEVAL

## K. Selçuk Candan

Arizona State University

## Maria Luisa Sapino

University of Torino

**CAMBRIDGE**
UNIVERSITY PRESS

**CAMBRIDGE**
UNIVERSITY PRESS

# Contents

v

Color plates follow page 38

# Preface

Database and multimedia systems emerged to address the needs of very different application domains. New applications (such as digital libraries, increasingly dynamic and complex web content, and scientific data management), on the other hand, necessitate a common understanding of both of these disciplines. Consequently, as these domains matured over the years, their respective scientific disciplines moved closer. On the media management side, researchers have been concentrating on media-content description and indexing issues as part of the MPEG7 and other standards. On the data management side, commercial database management systems, which once primarily targeted traditional business applications, today focus on media and heterogeneous-data intensive applications, such as digital libraries, integrated database/information-retrieval systems, sensor networks, bioinformatics, e-business applications, and of course the web.

There are three reasons for the heterogeneity inherent in multimedia applications and information management systems. First, the semantics of the information captured in different forms can be drastically different from each other. Second, resource and processing requirements of various media differ substantially. Third, the user and context have significant impacts on what information is relevant and how it should be processed and presented. A key observation, on the other hand, is that rather than being independent, the challenges associated with the semantic, resource, and context-related heterogeneities are highly related and require a common understanding and unified treatment within a multimedia data management system (MDMS). Consequently, internally a multimedia database management system looks and functions differently than a traditional (relational, object-oriented, or even XML) DBMS.

Also acknowledging the fact that web-based systems and rich Internet applications suffer from significant media- and heterogeneity-related hurdles, we see a need for undergraduate and graduate curricula that not only will educate students separately in each individual domain, but also will provide them a common perspective in the underlying disciplines. During the past decade, at our respective institutions, we worked toward realizing curricula that bring media/web and database educations closer. At Arizona State University, in addition to teaching a senior-level

ix

"Multimedia Information Systems" course, one of us (Prof. Candan) introduced a graduate course under the title "Multimedia and Web Databases." This course offers an introduction to features, models (including fuzzy and semistructured) for multimedia and web data, similarity-based retrieval, query processing and optimization for inexact retrieval, advanced indexing, clustering, and search techniques. In short, the course provides a "database" view of media management, storage, and retrieval. It not only educates students in media information management, but also highlights how to design a multimedia-oriented database system, why and how these systems evolve, and how they may change in the near future to accommodate the needs of new applications, such as search engines, web applications, and dynamic information-mashup systems. At the University of Torino, the other author of this book (Prof. Sapino) taught a similar course, but geared toward senior-level undergraduate students, with a deeper focus on media and features.

A major challenge both of us faced with these courses was the lack of an appropriate textbook. Although there are many titles that address different aspects of multimedia information management, content-based information retrieval, and query processing, there is currently no textbook that provides an integrated look at the challenges and technologies underlying a multimedia-oriented DBMS. Consequently, both our courses had to rely heavily on the material we ourselves have been developing over the years. We believe it is time for a textbook that takes an integrated look at these increasingly converging fields of multimedia information retrieval and databases, exhaustively covers existing multimedia database technologies, and provides insights into future research directions that stem from media-rich systems and applications. We wrote this book with the aim of preparing students for research and development in data management technologies that address the needs of rich media-based applications. This book's focus is on algorithms, architectures, and standards that aim at tackling the heterogeneity and dynamicity inherent in real data sources, rich applications, and systems. Thus, instead of focusing on a single or even a handful of media, the book covers fundamental concepts and techniques for modeling, storing, and retrieving heterogeneous multimedia data. It includes material covering semantic, context-based, and performance-related aspects of modeling, storage, querying, and retrieval of heterogeneous, fuzzy, and subjective (multimedia and web) data.

We hope you enjoy this book and find it useful in your studies and your future endeavors involving multimedia.

*K. Selçuk Candan and Maria Luisa Sapino*