

Introduction to Video Search Engines

David C. Gibbon • Zhu Liu

Introduction to Video Search Engines

 Springer

David C. Gibbon
AT & T Labs Research
200 Laurel Ave.
Middletown NJ 07748
USA
dcg@research.att.com

Zhu Liu
AT & T Labs Research
200 Laurel Ave.
Middletown NJ 07748
USA
zliu@research.att.com

ISBN: 978-3-540-79336-6

e-ISBN: 978-3-540-79337-3

Library of Congress Control Number: 2008932565

ACM Computing Classification (1998): H.2, H.3, H.5, 1.2.10, 1.4

© 2008 Springer Berlin Heidelberg

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilm or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer. Violations are liable to prosecution under the German Copyright Law.

The use of general descriptive names, registered names, trademarks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

Cover design: KünkelLopka, Heidelberg

9 8 7 6 5 4 3 2 1

springer.com

Preface

The evolution of technology has set the stage for the rapid growth of the video Web: broadband Internet access is ubiquitous, and streaming media protocols, systems, and encoding standards are mature. In addition to Web video delivery, users can easily contribute content captured on low cost camera phones and other consumer products. The media and entertainment industry no longer views these developments as a threat to their established business practices, but as an opportunity to provide services for more viewers in a wider range of consumption contexts. The emergence of IPTV and mobile video services offers unprecedented access to an ever growing number of broadcast channels and provides the flexibility to deliver new, more personalized video services. Highly capable portable media players allow us to take this personalized content with us, and to consume it even in places where the network does not reach.

Video search engines enable users to take advantage of these emerging video resources for a wide variety of applications including entertainment, education and communications. However, the task of information extraction from video for retrieval applications is challenging, providing opportunities for innovation. This book aims to first describe the current state of video search engine technology and second to inform those with the requisite technical skills of the opportunities to contribute to the development of this field.

Today's Web search engines have greatly improved the accessibility and therefore the value of the Web. The top portals prominently feature search capabilities and go beyond text search to include image and video search in various forms. A number of smaller companies have begun to offer more sophisticated media search features based on content analysis. Academic research groups have been actively developing algorithms and prototypes in this area for over a decade; incorporating and advancing previously existing constituent technologies.

Most media search systems rely on available metadata or contextual information in text form. Syndication formats such as RSS provide organized access to media sources and include descriptive global metadata. While these information sources are valuable and should be exploited, they are limited because they are typically brief, high level and subjective.

Therefore the current focus of media indexing research is to develop algorithms to exploit the media content itself as much as possible to augment available metadata. In some cases, the media may contain associated text streams such as closed caption or song lyrics. By extracting and operating on these streams, a textual representation of the dialog is obtained and existing text information retrieval methods can then be applied to retrieve relevant media. Speech recognition can be employed to create an approximation of the transcription, and techniques such as video optical character recognition can also be used to generate a textual representation of the media content. Although these technologies are inherently error prone, they have been used with success for indexing applications. Advanced speech retrieval systems use phonetic search to deal with the “out of vocabulary” problem and maintain alternative hypotheses in the form of lattices to boost recall.

Media retrieval that goes beyond the textual media component is more complex because the basic media features are not well defined and may not scale well for large archives. Further, formulating queries may not be as simple as typing a keyword. However, systems have been designed to, for example, retrieve images similar to a given image (query by example) or retrieve images based on a specification of color or shape. For navigating video retrieval results, techniques such as video skimming or mosaicing have been proposed.

The book will have a practical emphasis with the goal of bringing researchers up to date on the state of the art in multimedia search technologies and systems. Part of the presentation will follow a logical flow from content acquisition, analysis to extract index data, data representation, media archival, retrieval and finally rendering results in a Web-based environment. Each of these major functional components will be outlined, and particular emphasis will be given to automated content analysis techniques since this is critical for operating video search engines at scale, and it presents on-going research challenges. To give the readers an understanding of the issues involved, individual media processing algorithms operating on text, audio and video will be addressed including text alignment, case restoration, entity extraction, speech recognition, speaker segmentation, and video shot boundary detection. Additionally, the value of operating on multiple media components simultaneously will be illustrated by examining multimodal processing techniques, e.g. for media segmentation. The role of media segmentation in improving relevance ranking for long-form content will be discussed.

Who should read this book?

The book is intended for senior undergraduates or first-year graduate students in computer science or computer engineering, as well as professionals working in related fields. Although not intended for experts working directly on video search engines, the book will present a refreshing, broad perspective on video search and will have value as a reference tool. The topic of multimedia search spans multiple disciplines so the book will be valuable to experts in the constituent technologies such as speech processing or information retrieval who are looking to broaden their knowledge beyond their current areas of expertise.

A basic knowledge of Web application technologies, databases and computer networking issues is assumed. While a basic knowledge of the constituent technologies would be helpful, the intent will be to present these at an introductory level and discuss only the elements applicable to the problem of video search. The book explains the overall process of video content acquisition, indexing and retrieval with browsing, provides overviews of constituent technologies such as information retrieval, Internet video systems, video and multimedia processing to extract index data, and gives examples of existing systems and describes their features. The readers will:

- understand at a basic level all of the technologies used in today's video search engines;
- learn which video indexing techniques are appropriate for a given type of video material and be able to make inferences about which methods will work for new video content types;
- be able to differentiate between proven, practical techniques and those that are speculative, under development, or of narrow applicability;
- be able to determine which topics in video search are of interest to them for further study.

How is this book organized?

The book is divided into three main sections:

- I. Background and fundamentals: Chapter 1 outlines the technology trends which dictate that video material will increasingly be made available on the Web and points out the challenges that video is much more difficult to search than text files, and it is more difficult to browse. Chapter 2 addresses the nature, availability, and attributes of

different sources of video data. Details about available metadata for different types of video (e.g. electronic program guides, transcripts, etc.) are also provided. Chapter 3 reviews Internet video systems, including topics such as bandwidth, compression, random access, streaming, standards, digital rights management, redirector files, etc. Chapter 4 introduces video search engine systems: the process of content acquisition, media processing, building a multimedia database, retrieval, media browsing.

- II. Media processing: To address the challenges, we need to move beyond existing metadata retrieval systems, and analyze the content to extract information for indexing. Chapter 5 gives an overview on automated methods, systems, and algorithms for processing media to extract information for indexing and retrieval purposes. Chapters 6 - 8 discuss the specific media processing technologies that are developed in video, audio, and text domains. Multimodal processing, which is designed to mitigate the error that is inevitable with single modal processing, is discussed in Chapter 9.
- III. Case studies: Chapter 10 reinforces the concepts of video processing through illustrative examples, and provides details about existing solutions. Practical issues are brought to light through presentation of a detailed case study including a system supporting rapid content queries on a 50,000 hour broadcast television archive spanning 10 years and supporting a wide range of streaming media types for different applications. Chapter 11 provides a review of currently deployed Web search engines and identifies a few trends in the field to provide a sense of future directions.

Acknowledgements

This book became possible due to the support and vision of Dr. Behzad Shahraray, the head of the Video and Multimedia Services Research Department at AT&T Labs Research, Dr. Richard Cox and Dr. Lawrence Rabiner, directors of the Speech and Image Processing Research Laboratory, and Prof. Yang Wang at the Electrical Engineering department of Polytechnic University. Our work has been inspired by our colleagues, Lee Begeja, Bernard Renger, Harris Drucker, Andrea Basso and Murat Sarclar. We also benefit from collaborations with Prof. Shi-fu Chang and Eric Zavesky at Columbia University. This book began as a tutorial that the authors gave at the WWW2006 Conference in Edinburgh Scotland as suggested by Robin Chen. All their support and help is greatly appreciated.

Contents

Preface	v
1 Video Search.....	1
1.1 Introduction	1
1.2 Addressing the Opportunity.....	2
1.3 Classification of Web Video Sites.....	5
1.3.1 Content Originators and Traditional Broadcasters	5
1.3.2 Aggregators	6
1.3.3 Download	6
1.3.4 Sharing.....	6
1.3.5 Application Specific	7
1.3.6 Other Video Systems.....	7
1.4 Classification of Video Sources.....	8
1.4.1 Webcams / Security.....	9
1.4.2 Video Telephony / Teleconferencing	9
1.4.3 Industrial / Academic / Medical	9
1.4.4 User Generated Content.....	10
1.4.5 Public Access and Government (PEG) Content	10
1.4.6 Enterprise Content	10
1.4.7 Rushes, Raw Footage	11
1.4.8 News	11
1.4.9 Advertising	11
1.4.10 Episodic TV Programming	11
1.4.11 Feature Films	12
1.4.12 Content Value.....	12
1.5 Challenges of Video Search.....	13
1.5.1 Acquisition	14
1.5.2 Media File Formats.....	15
1.5.3 Data Transport.....	16
1.5.4 Browsing.....	16
1.5.5 Duplication	17
1.5.6 Ranking and Indexing.....	17
1.6 Advantages of Video Search over Text.....	18

1.6.1 Applications.....	18
1.6.2 Metadata	19
1.7 Metadata vs. Content	19
1.7.1 Content-based retrieval.....	19
1.8 Conclusion	20
References	21
2 Video Data Sources and Applications.....	23
2.1 Introduction	23
2.1.1 Evolution of Digital Media Metadata.....	23
2.1.2 Consumer Video Metadata	24
2.1.3 Metadata Loss.....	24
2.1.4 Metadata Standards	25
2.1.5 Dublin Core	26
2.1.6 MPEG-7.....	27
2.1.7 MPEG-21.....	27
2.2 Essential Media Metadata.....	29
2.2.1 Embed Global Metadata	29
2.2.2 Elementary Metadata.....	29
2.3 Metadata for Personal Media Collections.....	31
2.3.1 Consumer Media Libraries	31
2.3.2 UPnP Forum	33
2.3.3 MP3 ID3	33
2.3.4 3GP / QuickTime / MP4.....	34
2.3.5 Metadata Services.....	34
2.3.6 Content Identification.....	36
2.3.7 Recorded Television.....	37
2.4 Media Syndication: RSS Content Description	39
2.4.1 Content Syndication	39
2.4.2 Media Enclosures	39
2.4.3 Podcasts	41
2.4.4 RSS for Content Ingest.....	42
2.4.5 MediaRSS.....	43
2.5 Metadata for Broadcast Television.....	43
2.5.1 Electronic Programming Guide (EPG).....	44
2.5.2 Extended Data Service (XDS).....	46
2.5.3 Program and System Identifier Protocol (PSIP).....	47
2.6 Metadata for Video on Demand	47
2.6.1 Introduction	47
2.6.2 Cable Labs	49
2.7 Production Metadata.....	50
2.8 Timed Text Formats	51

2.8.1 Introduction	51
2.8.2 Synchronization Precision and Resolution	52
2.8.3 Transcripts	53
2.8.4 Closed Captions.....	54
2.8.5 Synchronized Accessible Media Interchange	55
2.8.6 Metadata from Social Sources	55
2.8.7 Metadata Issues.....	55
2.9 Conclusion	56
References	56
3 Internet Video	59
3.1 Introduction	59
3.2 Digital Video	59
3.2.1 Aspect Ratio	59
3.2.2 Luminance and Chrominance Resolution.....	61
3.2.3 Video Compression	62
3.3 Internet Protocol Media Systems.....	66
3.3.1 Transport.....	66
3.3.2 Searching VoD vs. Live.....	67
3.3.3 IPTV	68
3.3.4 Rights Management.....	70
3.3.5 Redirector Files	70
3.3.6 Layered Encoding.....	73
3.3.7 Illustrated Audio	73
3.4 Media Captioning	74
3.5 Conclusion	75
References	76
4 Video Search Engine Systems.....	77
4.1 Introduction	77
4.2 Content Acquisition	78
4.2.1 Metadata Normalization	78
4.2.2 User Contributed.....	79
4.2.3 Syndicated Contribution.....	80
4.2.4 Broadcast Acquisition.....	81
4.3 Content Processing	82
4.3.1 Asset Management	82
4.4 Retrieval.....	84
4.5 User Perspectives.....	85
4.5.1 Interaction States	85
4.5.2 Granularity of Search Results Representation	87
4.6 Factors Concerning Scalability.....	88

- 4.6.1 Introduction 88
- 4.6.2 Acquisition 89
- 4.6.3 Processing..... 89
- 4.6.4 Storage..... 90
- 4.6.5 Retrieval 91
- 4.7 Retrieval Interfaces..... 92
- 4.8 Typical System Features..... 93
- 4.9 Conclusion..... 94
- References 94

- 5 Media Processing 97**
 - 5.1 Introduction 97
 - 5.2 Feature Extraction..... 99
 - 5.3 Media Segmentation 100
 - 5.4 Clustering, Structure Generation 101
 - 5.5 Real-Time Processing..... 103
 - 5.6 Systems Issues and Architectures..... 103
 - 5.7 Conclusion 104
 - References 105

- 6 Video Processing 107**
 - 6.1 Introduction 107
 - 6.2 Shot Boundary Determination 108
 - 6.2.1 Feature Extraction 110
 - 6.2.2 Shot Boundary Detectors..... 111
 - 6.2.3 Fusion of Detector Results 117
 - 6.2.4 Evaluation Results 117
 - 6.3 Representative Image Selection..... 118
 - 6.4 Face Detection 121
 - 6.5 Face Recognition 126
 - 6.6 Video Optical Character Recognition..... 129
 - 6.7 Concept Detection 131
 - 6.7.1 Color Feature 133
 - 6.7.2 Texture Feature..... 133
 - 6.7.3 Edge Feature..... 135
 - 6.8 Video Browsing..... 135
 - 6.9 Conclusion 140
 - References 141

- 7 Audio Processing..... 145**
 - 7.1 Introduction 145
 - 7.2 Audio Signal and Its Representation 146

7.3 Audio Features.....	148
7.3.1 Frame-Level Features	148
7.3.2 Clip-Level Features	154
7.4 Audio Segmentation	156
7.4.1 Speaker Segmentation	157
7.4.2 Audio Scene Segmentation.....	158
7.5 Audio Content Categorization	160
7.5.1 Speaker Recognition.....	160
7.5.2 Audio Scene Detection	162
7.5.3 Music Genre Classification	163
7.6 Speech Recognition	164
7.7 Audio Query and Browsing Techniques.....	166
7.7.1 SpeechLogger	167
7.7.2 Query by Example	171
7.8 Conclusion	172
References	173
8 Text Processing	177
8.1 Introduction	177
8.2 Story Segmentation.....	178
8.2.1 Cue Phrases	178
8.2.2 Cosine Similarity	179
8.2.3 Dynamic Programming.....	181
8.2.4 Topic Classification.....	183
8.3 Named Entity Extraction	183
8.3.1 Rule Based NEE	184
8.3.2 Data Driven NEE.....	185
8.3.3 NEE Tools	186
8.4 Part-of-Speech Tagging.....	187
8.5 Capitalization.....	189
8.5.1 Linguistic Processing Architecture.....	191
8.5.2 Web Document Collection	191
8.5.3 Text Capitalization Algorithm.....	192
8.6 Information Retrieval.....	194
8.6.1 Stemming.....	194
8.6.2 Term Weighting.....	195
8.6.3 Ranking.....	196
8.7 Text Summarization	197
8.7.1 Keyword Extraction.....	199
8.8 Conclusion	201
References	201

9 Multimodal Processing	203
9.1 Introduction	203
9.2 Case Studies.....	205
9.2.1 Closed Caption Alignment	205
9.2.2 Multimodal News Story Segmentation.....	209
9.2.3 Major Cast Detection.....	214
9.3 Conclusion	217
References	217
10 Research Systems	221
10.1 Introduction	221
10.2 Academic and Industrial Research	222
10.3 Early Internet Deployments.....	226
10.3.1 SpeechBot.....	226
10.3.2 StreamSage	227
10.3.3 SingingFish.....	227
10.4 Selected Commercial Systems.....	228
10.4.1 Virage and Convera	228
10.4.2 Nexidia (FastTalk).....	228
10.5 Resources: Datasets, Evaluations, Conferences	229
10.6 Media Monitoring Deployments.....	231
10.7 Case Study: AT&T MIRACLE	232
10.7.1 Introduction	232
10.7.2 System Architecture	232
10.7.3 Collections.....	233
10.7.4 Data Organization.....	235
10.7.5 Acquisition / Ingest.....	236
10.7.6 Content Processing	238
10.7.7 Real-time processing	239
10.7.8 Query Engine.....	239
10.7.9 Applications.....	240
10.7.10 Performance.....	240
10.8 Conclusion	242
References	242
11 Current Trends in Video Search	247
11.1 Introduction	247
11.2 Video Production.....	248
11.2.1 Metadata Retention.....	248
11.2.2 Multiple Distribution Channels	248
11.2.3 Mobisodes and Webisodes	249
11.3 Video Distribution	249

11.3.1 Streaming Protocols.....	250
11.3.2 Electronic Sell Through.....	250
11.3.3 Peer-to-peer Delivery	251
11.3.4 Managed Download.....	251
11.3.5 Syndication	252
11.4 The Video Web and User Interaction	252
11.4.1 Web-Based Editing.....	252
11.4.2 Media Browsing	252
11.4.3 Social Tagging.....	253
11.4.4 Dynamic Interfaces.....	253
11.4.5 Video Blogs (vlogs).....	254
11.4.6 Integrated Collections.....	254
11.5 Television Technology and Consumption	254
11.5.1 Proliferation of Channels.....	255
11.5.2 Live to Time Shifted.....	255
11.5.3 Mobile Consumption	255
11.6 Trends in Media Devices	256
11.6.1 Increased Media Capabilities.....	256
11.6.2 Increasing Accessibility.....	257
11.6.3 DRM.....	257
11.6.4 Home Media Systems.....	257
11.7 Media Processing Research	257
11.8 Deployments	260
11.9 Conclusion	261
References	261
Glossary	265
Index.....	271