

11

Multimedia Indexing and Retrieval Using Natural Language, Speech and Image Processing Methods

Harris Papageorgiou, Prokopis Prokopidis, Athanassios Protopapas and George Carayannis

11.1 Introduction

Throughout the chapter, we provide details on implementation issues of practical systems for efficient multimedia retrieval. Moreover, we exemplify algorithms and technologies by referring to practices and results of an EC-funded project called Combined IMage and Word Spotting (CIMWOS) [1], developed with the hope that it would be a powerful tool in facilitating common procedures for intelligent indexing and retrieval of audiovisual material. CIMWOS used a multifaceted approach for the location of important segments within multimedia material, employing state-of-the-art algorithms for text, speech and image processing in promoting reuse of audiovisual resources and reducing budgets of new productions.

In the following three sections, we focus on technologies specific to speech, text and image, respectively. These technologies incorporate efficient algorithms for processing and analysing relevant portions from various digital media and thus generating high-level semantic descriptors in the metadata space. After proposing an architecture for the integration of all results of processing, we present indicative evaluation results in the context of CIMWOS. Relevant Content-based Information Retrieval (CBIR) research prototypes and commercial systems are briefly presented in the Section 11.7.

11.2 Audio Content Analysis

Audio content analysis for multimedia indexing and retrieval refers to a number of audio features gathered from acoustic signals. These audio features combined with visual descriptors have effectively been exploited in scene content analysis, video segmentation, classification and summarization [2]. Audio features are extracted in the short-term frame level and the long-term clip level (usually called window). A frame is defined as a group of adjacent samples lasting about 10–40 ms, within which the audio signal is assumed to be stationary. Frame-level features are:

- *Time-domain features*: volume, zero crossing rate, pitch.
- *Frequency-domain features* (based on the fast fourier transform (FFT) of the samples in a frame): spectrum of an audio frame, the ratio of the energy in a frequency subband to the total energy (ERSB), the mel-frequency cepstrum coefficients (MFCC) widely used for speech recognition and speaker recognition.

Clip-level features capture the temporal variation of frame-level features on a longer time scale and include volume-based features like the VSTD-standard deviation of the volume over a clip, the ZSTD proposed by Liu et al. [3] and pitch-based features like the SPR, which is the percentage of frames in a clip that have similar pitch as the previous frames.

Audio tracks are further segmented and classified into a number of classes such as speech, silence, noise, music and environmental sound (recognition of shots, cries, explosions, door slams etc.). State-of-the-art continuous speech recognizers perform on the speech segments, generating speech transcriptions. Producing a transcript of what is being said, and determining who is speaking and for how long are all challenging problems, mainly due to the continuous nature of an audio stream: segments of diverse acoustic and linguistic nature exhibit a variety of problematic acoustic conditions, such as spontaneous speech (as opposed to read or planned speech), limited bandwidth (e.g. telephone interviews), and speech in the presence of noise, music and/or background speakers. Such adverse background conditions lead to significant degradation in the performance of speech recognition systems if appropriate countermeasures are not taken. Likewise, segmentation of the continuous speech stream into homogeneous sections (with respect to acoustic/background conditions, speaker and/or topic) poses serious problems. Successful segmentation, however, forms the basis for further adaptation and processing steps. Adaptations to the varied acoustic properties of the signal or to a particular speaker, and enhancements to the segmentation process, are generally acknowledged [4] as key research areas that will result in rendering indexing systems usable for actual deployment. This is reflected in the amount of effort and the number of projects dedicated to the advancement of the current state-of-the-art in these areas [5].

A speech processing subsystem (SPS) usually comprises a set of technology components responsible for a particular task, arranged in a pipeline. These tasks typically include speaker change detection (SCD), automatic speech recognition (ASR), speaker identification (SID) and speaker clustering (SC). To save processing time, SC and SID can be run in parallel to ASR. Input to the overall system is the audio stream, while the final output is a set of audio descriptors containing speech transcriptions and identified speakers. A typical SPS architecture is depicted in Figure 11.1, while each specific task is presented in detail in the following subsections.

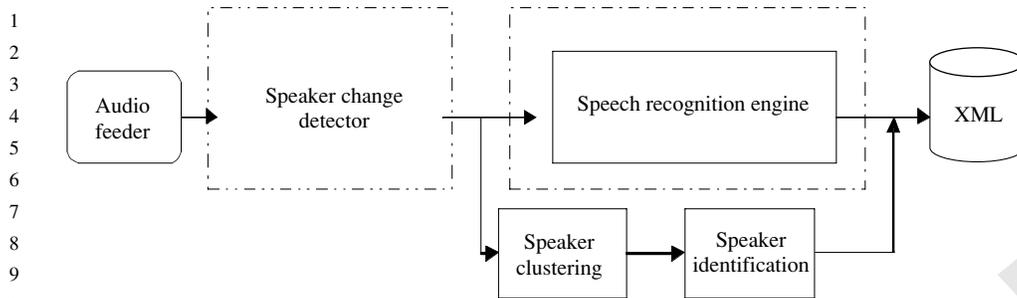


Figure 11.1 Speech processing subsystem

11.2.1 Speaker Change Detection

SCD aims at extracting only the speech portions of the incoming audio stream, and at partitioning those sections into homogeneous sections of speech. The technology used for SCD involves initially running the speech recognizer with a phone-level model (PLD), which is a set of broad phonetic classes of speech sounds like vowels, nasals, obstruents, as well as some non-speech sounds like music, laughter etc.

Using the information produced by the PLD, the SCD module determines, during runtime, likely boundaries of speaker turns based on non-speech portions (e.g. silence, music, background noise) in the signal. Parameters used in this decision are usually empirically determined, as they differ from language to language and even from domain to domain.

11.2.2 Speaker Clustering and Identification

The SC module is responsible for grouping predetermined chunks of audio (the SCD output described above) into a number of clusters. Resulting clusters ideally represent utterances produced by a single speaker. SID uses a predefined set of target models to identify the speaker of a segment or, if identification is not possible, to determine the speaker's gender. The set of target speakers typically comprises anchor speakers of TV stations as well as persons of public interest such as politicians. Current SID technology employs Gaussian mixture models (GMM), trained on corpora of speaker-annotated transcriptions. On the other hand, SC does not use any trained models, but rather a variant of the generalized likelihood ratio criterion at runtime. Evaluation of the CIMWOS SID module on predefined test sets produced 97% precision and 97% recall figures for English, 93% precision and 93% recall for French, and 92% precision and 89% recall for Greek.

11.2.3 Automatic Speech Recognition

In most state-of-the-art speech processing systems [6, 7], ASR is performed via a large vocabulary, speaker-independent, gender-independent, continuous speech recognizer. The recognizers typically use hidden markov models (HMM) with GMM for acoustic modelling (AM) and n-gram models estimated on large corpora for language modelling (LM).

The task for the ASR is to convert the information passed on to it by the SCD stage into time-tagged text. This is often done in a multi-stage approach by finding the most likely sequence of words, given the input audio stream. Each level of processing uses increasingly

1 complex acoustic and language models, thereby keeping a balance between model complexity,
 2 search space size and quality of the results. The AMs used by different stages are HMM-
 3 based mixture models with various degrees and kinds of tying (sharing of model parameters).
 4 They form a representation of the acoustic parameters (acoustics, phonetics, environment,
 5 audio source, speaker and channel characteristics, recording equipment etc.). The LM is an
 6 n-gram model with different types of back-off and smoothing. It contains a representation of
 7 the words included in a vocabulary, and of what words and sequences of words are likely to
 8 occur together or follow each other. The recognizer's vocabulary consists of several thousand
 9 wordforms and their pronunciation models. Using wordforms as opposed to words means that
 10 plural and singular forms (or any kind of declinational or conjugational forms) of a word are
 11 considered separate entities in the vocabulary.

12 ASR in CIMWOS uses a multi-stage approach to provide for flexibility and best use of
 13 resources. During *decoding*, in a first step, a *fast-match* is performed using the simplest (and
 14 thus fastest) models. Then a *detailed-match* follows on the results of fast-match. Finally a
 15 *rescoring pass* is applied to the results of the second stage. This staged approach allows for
 16 a dramatical reduction in search space while at the same time recognizer performance and
 17 accuracy are kept high. Evaluation of the CIMWOS ASR module on predefined test sets for
 18 English and french broadcast news (BNs) yielded word error rates (WERs) of 29% and 27%,
 19 respectively. For Greek BNs a 33% WER is reported. Figure 11.2 shows a sample portion of
 20 an XML file automatically created by the speech processing subsystem in CIMWOS.

```

21
22
23 <?xml version="1.0" encoding="UTF-8" ?>
24 <SpeechAnnotation project="CIMWOS">
25 <Header type="SpeechRecognitionMetadata">
26   <Media id="RTBF_20011012_1930_News" xml:lang="fr">
27     ...
28     <Passage id="p77" speaker="male 14" gender="male"
29       mediaRelIncrTimePoint="131864" mediaIncrDuration="1420"
30       xml:lang="fr">
31       <Word id="w3541" mediaRelIncrTimePoint="131864"
32         mediaIncrDuration="13" confidence="0.99">Dans</Word>
33       <Word id="w3542" mediaRelIncrTimePoint="131878"
34         mediaIncrDuration="8" confidence="0.99">le</Word>
35       <Word id="w3543" mediaRelIncrTimePoint="131887"
36         mediaIncrDuration="17" confidence="0.99">hall</Word>
37       <Word id="w3544" mediaRelIncrTimePoint="131905"
38         mediaIncrDuration="10" confidence="0.99">de</Word>
39       <Word id="w3545" mediaRelIncrTimePoint="131916"
40         mediaIncrDuration="48" confidence="0.99">départ</Word>
41       <Word id="w3546" mediaRelIncrTimePoint="131965"
42         mediaIncrDuration="5" confidence="0.99">de</Word>
43       <Word id="w3547" mediaRelIncrTimePoint="131971"
44         mediaIncrDuration="4" confidence="0.99">1'</Word>
45       <Word id="w3548" mediaRelIncrTimePoint="131976"
46         mediaIncrDuration="108" confidence="0.99">aéroport</Word> ...
47     </Passage>
48     ...
49   </Media>
50 </SpeechAnnotation>
  
```

Figure 11.2 Segment of speech processing generated metadata

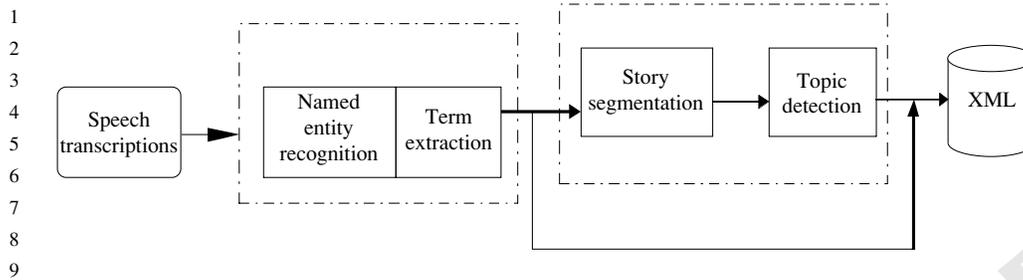


Figure 11.3 Text processing subsystem

11.3 Text Processing Subsystem

Following audio processing, text processing tools are applied on the textual data (speech transcriptions) produced by the speech processing subsystem, in an attempt to enable a retrieval system with the ability to answer questions like what topic a passage is about, or which organizations are mentioned. Named entity detection (NED), term extraction (TE), story segmentation (SD) and topic detection (TD) are tasks that can be included in a text processing pipeline like the one depicted in Figure 11.3.

11.3.1 Named Entity Detection and Term Extraction

The task of the NED module in the context of a multimedia retrieval system is to identify named locations, persons and organizations, as well as dates, percentages and monetary amounts, in the speech transcriptions automatically produced by the SPS. NED involves an initial preprocessing phase where sentence boundary identification, and part-of-speech tagging are performed on textual input. State-of-the-art NED implementations involve lookup modules that match lists of NEs and trigger words against the text, hand-crafted and automatically generated pattern grammars, maximum entropy modelling, HMM models, decision-tree techniques, SVM classifiers etc. [8].

As an example, the NED module for the Greek language in CIMWOS [9] is a combination of list-based matching and parsing with a grammar of rules compiled into finite-state transducers. In a collection of TV news broadcasts, this module identified 320 NEs, compared to the 555 identified by human annotators using a suitable annotation tool (Figure 11.4). There were 235 correct guesses, 85 false positives and 320 missed instances. Although promising precision scores were obtained for persons and locations, (Figure 11.5), lower recall in all NE types is due to:

- greater out-of-vocabulary (OOV) rate in the ASR for Greek;
- missing proper names in the vocabulary of the ASR engine;
- the diversity of domains found in broadcasts, leading to a large number of names and NED indicators that were not incorporated in the NED module's resources.

The TE task involves the identification of single or multi-word indicative keywords (index terms) in the output of the speech processing system. Systems for automatic term extraction using both linguistic and statistical modelling are reported in the literature [10]. Linguistic

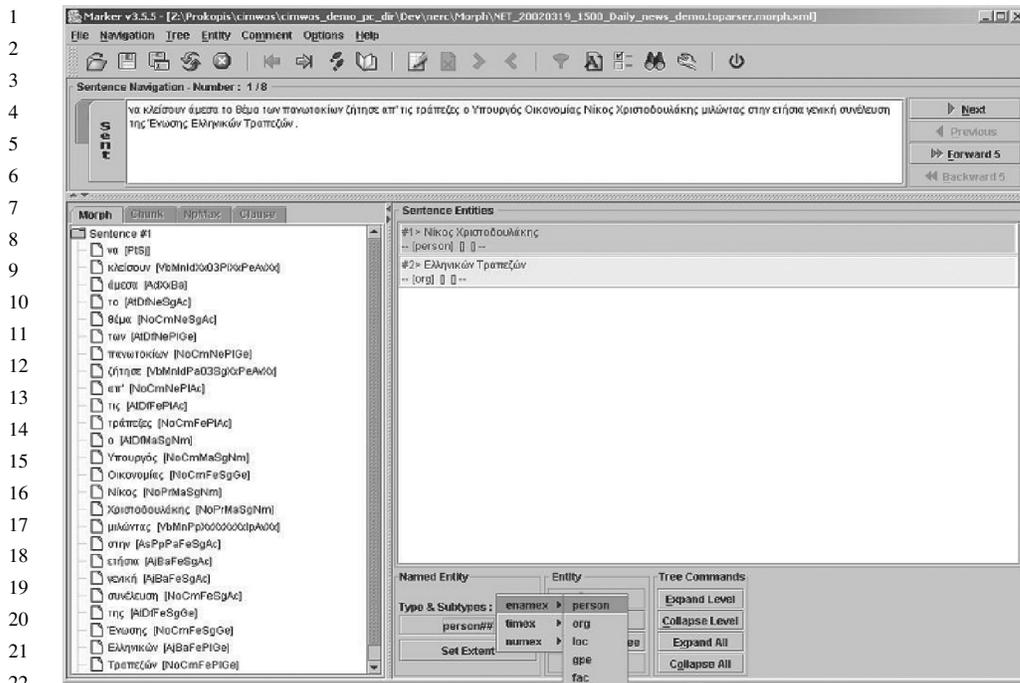


Figure 11.4 Named entity annotation tool

Named entities evaluation

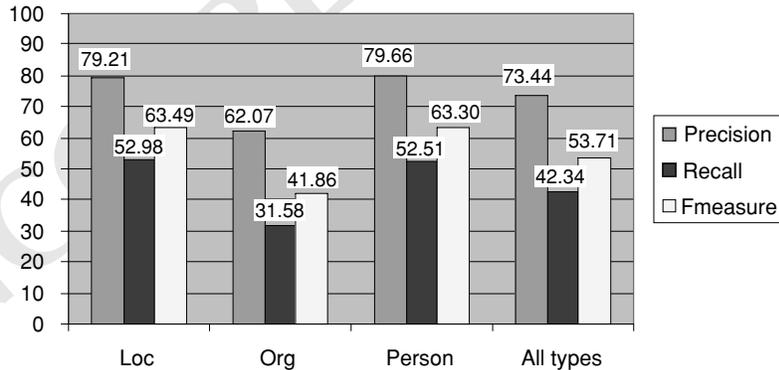


Figure 11.5 Performance of NED module for each NE type

processing is usually performed through an augmented term grammar, the results of which are statistically filtered using frequency-based scores.

The term extractor in CIMWOS was used to automatically identify terms in a testing corpus of manually annotated transcriptions, scoring a 60.28% recall and a 34.80% precision. To

1 obtain higher precision scores, stricter statistical filtering needs to be applied. Nevertheless, a
2 relaxed statistical filtering like the one currently used produces higher recall measures, a factor
3 which plays a more important role in an information retrieval context.

6 *11.3.2 Story Detection and Topic Classification*

7 Story detection and topic classification modules often employ the same set of models, trained
8 on an annotated corpus of stories and their manually associated topics. The basis of these
9 technologies is a generative, mixture-based hidden Markov model that includes one state per
10 topic, as well as one state modelling general language; that is, words not specific to any topic.
11 Each state models a distribution of words given the particular topic. After emitting a single
12 word, the model re-enters the beginning state and the next word is generated. At the end of a
13 story the final state is reached. Detection is performed running the resulting models on a sliding
14 window of fixed size, thereby noting the change in topic-specific words as the window moves
15 on. The result of this phase is a set of 'stable regions' in which topics change only slightly, or
16 not at all. Building on the story boundaries, sections of text are classified according to a set of
17 topic models.

18 In the Informedia project at Carnegie Mellon University [11] each video document is par-
19 titioned into story units based on text, image and audio metadata jointly. Silence periods are
20 identified and subsequently aligned with the nearest shot break indicating the boundary of
21 a story unit. Text markers such as punctuation are taken into consideration in case closed
22 caption is available. In the MAESTRO system [12] developed by SRI, topic boundaries are
23 detected based on prosodic information extracted from the speech waveforms (pause and pitch
24 patterns) combined with word usage statistics. The MITRE Broadcast News Navigator (BNN)
25 incorporates a story segmentation module based on finite state machines [13].

28 **11.4 Image Processing Subsystem**

29 Image processing and understanding is an important field of research where close ties between
30 academic and commercial communities have been established. In the multimedia content
31 analysis framework, a typical image processing subsystem (IPS) like the one presented here
32 consists of modules performing video segmentation and keyframe extraction, face detection
33 and face identification, object identification and video text detection and recognition.

36 *11.4.1 Video Segmentation and Keyframe Extraction*

37 The segmentation of video sequences is a prerequisite for a variety of image processing appli-
38 cations. Video streams consist of many individual images, called frames, generally considered
39 as the smallest unit to be concerned with when segmenting a video. An uninterrupted video
40 stream generated by one camera is called a shot (for example, a camera following an aeroplane,
41 or a fixed camera focusing on an anchorperson), while a shot cut is the point at which shots
42 change within a video sequence. The video segmentation task involves detecting shot cuts and
43 thus partitioning raw material into shots. Under the assumption that frames within a shot have
44 a high degree of similarity for each shot, a few representative frames are selected, referred
45 to as keyframes. Each keyframe represents a part of the shot called subshot. The subdivision
46

1 of a shot into subshots occurs when, for example, there is an abrupt camera movement or
2 zoom operation, or when the content of the scene is highly dynamic so that a single keyframe
3 no longer suffices to describe the content of the entire shot. Keyframes contain most of the
4 static information present in a shot, so that subsequent modules like face detection and object
5 identification can focus on keyframes only without scanning the whole video sequence.

6 In order to detect shot cuts and select keyframes, different methods have been developed for
7 measuring the differences between consecutive frames and applying adaptive thresholding on
8 motion and texture cues. The performance of the CIMWOS video segmentation module was
9 evaluated on some of the news broadcasts available in the project's collection. A recall rate of
10 99% and a precision of 97% were observed in the case of abrupt shot cuts. For smooth shot
11 cuts, performance was lower and highly dependent on the video content, yet still comparable
12 to the state of the art.

13 The performance of the keyframe selection is much harder to quantify, as there is no objec-
14 tive groundtruth available: what constitutes an appropriate keyframe is a matter of semantic
15 interpretation by humans, and cannot be resolved on the basis of low-level image cues. Never-
16 theless, judging from comments by test users, the CIMWOS keyframe extraction module was
17 indeed able to reduce the number of redundant keyframes while keeping important information
18 available.

19
20

21 *11.4.2 Face Detection and Face Identification*

22
23 Given an arbitrary image, the goal of face detection is to determine whether or not there are any
24 faces in the image, and, if present, to return the face dimensions in the keyframe. Face detection
25 is a challenging task, since several factors influence the appearance of the face in the image.
26 These include identity, pose (frontal, half-profile, profile), presence or absence of facial features
27 such as beards, moustaches and glasses, facial expression, occlusion and imaging conditions.
28 Face detection has been, and still is, a very active research area within the computer vision
29 community. It is one of the few attempts to recognize from a set of images a class of objects
30 for which there is a great deal of within-class variability. It is also one of the few classes of
31 objects for which this variability has been captured using large training sets of images.

32 Recent experiments by neuroscientists and psychologists show that face recognition is a
33 dedicated process in the human brain. This may have encouraged the view that artificial face
34 identification systems should also be face-specific. Automatic face identification is a chal-
35 lenging task and has recently received significant attention. Rapidly expanding research in
36 the field is based on recent developments in technologies such as neural networks, wavelet
37 analysis and machine vision. Face identification has a large potential for commercialization
38 in applications involving authentication, security system access, and advanced video surveil-
39 lance. Nevertheless, in spite of expanding research, results are far from perfect, especially in
40 uncontrolled environments, because of lighting conditions and variations, different facial ex-
41 pressions, background changes and occlusion problems. One of the most important challenges
42 in face recognition is to distinguish between intra-personal variation (in the appearance of a
43 single individual due to different facial expressions, lighting etc.) and extra-personal variation
44 (between different individuals).

45 In CIMWOS, the face detection and face identification modules associate detected faces
46 occurring in video streams with names [14] (Figure 11.6). Both modules are based on support



18 **Figure 11.6** Face detection results

19
20
21 vector machine models trained on an extensive database of facial images with a large variation
22 in pose and lighting conditions. Additionally, a semantic base consisting of important persons
23 that should be identified has been constructed. During identification, images extracted from
24 keyframes are compared to each model. At the decision stage, the scores resulting from the
25 comparison are used either to identify a face or to reject it as 'unknown'.

26 In order to be able to compare the CIMWOS FD performance with other systems' results
27 reported in the literature, the module was evaluated against the MIT/CMU dataset [15], where
28 groundtruth information for frontal face detection is readily available. A correct detection rate
29 of 90% (with 1 false detection in 10. 000 windows) was obtained. A qualitative evaluation on
30 the CIMWOS news broadcasts database shows that similar results are obtained in that context,
31 if the extremely small or out-of-plane rotated faces are ignored.

32 33 34 *11.4.3 Video Text Detection and Recognition*

35
36 Text recognition in images and video aims at integrating advanced optical character recognition
37 (OCR) technologies and text-based search, and is now recognized as a key component in the
38 development of advanced video and image annotation and retrieval systems. Unlike low-level
39 image features (such as colour, texture or shape), text usually conveys direct semantic informa-
40 tion on the content of the video, like a player's or speaker's name, location and date of an event
41 etc. However, text characters contained in video are usually of low resolution, of any colour
42 or greyscale value (not always white), embedded in complex background. Experiments show
43 that applying conventional OCR technology directly to video text leads to a poor recognition
44 rate. Therefore, efficient location and segmentation of text characters from the background is
45 necessary to fill the gap between images or video documents and the input of a standard OCR
46 system.



Figure 11.7 Video text detection and recognition

In CIMWOS, the video text detection and recognition module [16, 17] is based on a statistical framework using state-of-the-art machine learning tools and image processing methods (Figure 11.7). It consists of four modules:

- Text detection, aiming at roughly and quickly finding image blocks that may contain a single line of text characters.
- Text verification, using a support vector machine model, to remove false alarms.
- Text segmentation, attempting to extract pixels from text images belonging to characters with the assumption that they have the same colour/greyscale value. This method uses a Markov random field model and an expectation maximization algorithm for optimization.
- Finally, all hypotheses produced by the segmentation algorithm are processed by the OCR engine. A string selection is made based on a confidence value, computed on the basis of character recognition reliability and a bigram language model.

The CIMWOS video text detector was evaluated on 40 minutes of video consisting of 247 text strings, 2899 characters and 548 words. Characters were correctly extracted and recognized in 94% of all cases, while a 92% recognition rate was reached as far as words are concerned. Precision showed that more than 95% of extracted characters correspond to the characters in the groundtruth. In the 2899 characters, 7% of characters were scene text characters and 93% were captions.

11.4.4 Object Identification

The object identification (OI) problem can be defined as the task of identifying a non-deformable man-made object in a 'recognition view', given knowledge accumulated from a set of previously seen 'learning views'. Until recently, visual object recognition was limited

1 to planar (flat) objects, seen from an unknown viewpoint. The methods typically computed
2 numerical geometric invariants from combinations of easily extractable image points and lines.
3 These invariants were used as a model of the planar object for subsequent recognition. Some
4 systems went beyond planar objects, but imposed limits on the possible viewpoints or on the
5 nature of the object. Probably the only approach capable of dealing with general 3D objects
6 and viewpoints is the 'appearance-based' one, which nevertheless requires a very large amount
7 of example views and has fundamental problems in dealing with cluttered recognition views
8 and occlusion [18]. A system capable of dealing with general 3D objects, from general viewing
9 directions, is yet to be proposed.

10 Since 1998, a small number of systems have emerged that seem to have the potential for
11 reaching the general goal. These are all based on the concept of 'region', defined as a small,
12 closed area on the object's surface. In CIMWOS, the object surface is decomposed into a
13 large number of regions automatically extracted from the images [19]. These regions are
14 extracted from several example views (or frames of a movie) and both their spatial and temporal
15 relationships are observed and incorporated in a model. This model can be gradually learned
16 as new example views (or video streams) are acquired. The power of such a proposal consists
17 fundamentally in two points: first, the regions themselves embed many small, local pieces of
18 the object at the pixel level, and can reliably be put in correspondence along the example views.
19 Even in the case of occlusion or clutter in the recognition view, a subset of the object's regions
20 will still be present. The second strong point is that the model captures the spatio-temporal
21 order inherent in the set of individual regions, and requires it to be present in the recognition
22 view. In this way the model can reliably and quickly accumulate evidence about the identity
23 of the object in the recognition view, even with only a small number of recognized regions.

24 Thanks to the good degree of viewpoint invariance of the regions, and to the strong model
25 and learning approach developed, the object recognition module in CIMWOS copes with 3D
26 objects of general shape, requiring only a limited number of learning views [20]. Moreover, it
27 can recognize objects from a wide range of previously unseen viewpoints, in possibly cluttered,
28 partially occluded, views. Based on experiments conducted on the project's database (searching
29 for 13 objects in 325 keyframes, on average, for each object), a 83.7% precision and a 96.4%
30 recall were obtained.

31

32

33 **11.5 Integration Architecture**

34

35 Due to the prevalence of XML as a medium for data exchange between applications, all
36 processing output in the three modalities (audio, image and text) can converge to a textual XML
37 metadata annotation document following standard MPEG-7 descriptors. These annotations
38 can be further processed, merged, synchronized and loaded into the multimedia database. A
39 merging component amalgamates various XML annotations to create a self-contained object
40 compliant with a database scheme.

41 The CIMWOS architecture follows an N-tier scenario by integrating a data services layer
42 (storage and retrieval of metadata), a business services layer incorporating all remote multi-
43 media processors (audio, video and text intelligent engines) and a user services layer which
44 basically includes the user interface (UI) and web access forms. Web services (SOAP) could
45 be used as the integration protocol to access heterogeneous and loosely-coupled distributed
46 technologies. The CIMWOS architecture is depicted in Figure 11.8.

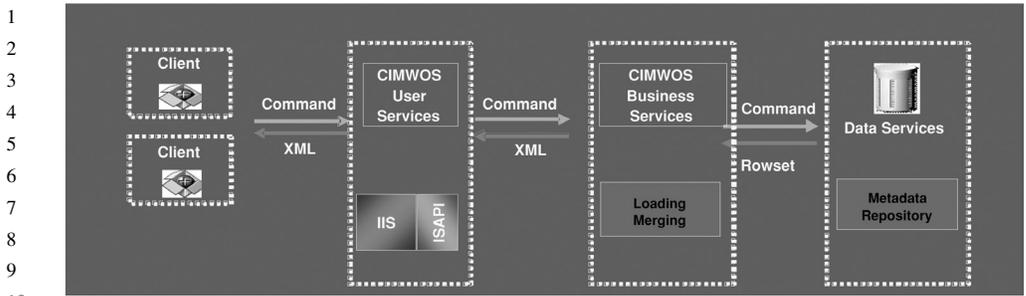


Figure 11.8 The CIMWOS architecture

11.5.1 Indexing, Search and Retrieval

The basic approach to indexing and retrieval is to apply speech, image and text technologies to automatically produce textual metadata and then to use information retrieval techniques on these metadata. The CIMWOS retrieval engine [21] is based on a weighted Boolean model. After query submission, the retrieval engine is invoked with the set of criteria combined flexibly with standard Boolean operators. A matching operation computes the similarity between the query and each passage to determine which ones contain the given set of combined query terms. The calculated similarity measure takes into consideration three different factors: metadata-level weights based on the overall precision of each multimedia processor, value-level statistical confidence measures produced by the engines and $tf*idf$ scores for all textual elements. Finally, passages are ranked based on the result of similarity computation.

The basic retrieval unit is the passage (Figure 11.9), which has the role of a document in a traditional system. Passages are defined on the speech transcriptions—that is, on the

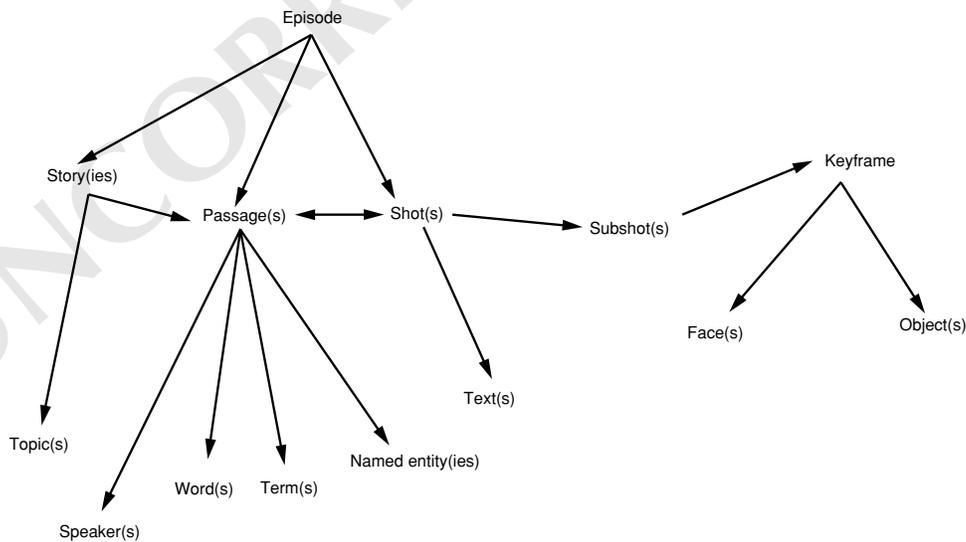


Figure 11.9 CIMWOS indexing schema

1 post-processed output of the ASR—therefore they correspond to some level of audio seg-
2 mentation. Other processing streams, that is, text and video, will typically result in segments
3 (stories and shots, respectively) not aligned with passages. This lack of correspondence with
4 the original multimedia object has implications for both indexing and retrieval. The solution
5 given must be based on the eventual use of the material, therefore the needs of the users should
6 be taken into account. Referring to the example of a user searching for Putin co-occurring with
7 Clinton on a particular crisis topic, it is necessary for a search and retrieval system to define the
8 temporal extent over which the existence of both annotations counts as co-occurrence, mean-
9 ing that the relevant time segment is to be retrieved. For news broadcasts, a news story is the
10 appropriate conceptual unit. However, in the merged annotations one has access to processed
11 outputs from the various streams, which do not map onto the desired unit in a well-defined
12 manner. Thus the selection of passages is made to satisfy both practical constraints (availabil-
13 ity, typical duration) and use requirements (to contain all relevant information and nothing
14 irrelevant, as much as possible).

15 The passage is indexed on a set of textual features: words, terms, named entities, speakers
16 and topics. Each passage is linked to one or multiple shots, and each shot is indexed on another
17 set of textual features: faces, objects and video text. By linking shots to passages, each passage
18 is assigned a broader set of features to be used for retrieval. Passages are represented as sets
19 of features and retrieval is based on computed similarity in the feature space.

20 A video clip can take a long time to be transferred, e.g. from the digital video library to the
21 user. In addition, it takes a long time for someone to determine whether a clip meets his or her
22 needs. Returning half an hour of video when only one minute is relevant is much worse than
23 returning a complete book when only one chapter is needed. Since the time to scan a video
24 cannot be dramatically shorter than the real time of the video, it is important to give users only
25 the material they need.

26 Multimedia information retrieval systems typically provide various visualized summaries of
27 video content where the user can preview a specific segment before downloading it. In Infor-
28 media, information summaries can be displayed at varying detail, both visually and textually.
29 Text summaries are displayed for each news story through topics and titles. Visual summaries
30 are given through thumbnails, filmstrips and dynamic video skims [11]. In CIMWOS, while
31 skimming a passage, the end user can view its associated metadata, the transcribed speech
32 and results of all processing components, as well as a representative sequence of thumbnails.
33 These thumbnails can act as indicative cues on how relevant the segment content is to the user
34 query. Each thumbnail corresponds to a keyframe that is recognized and extracted by the video
35 segmentation module. Results can also include more pieces of bibliographical information,
36 such as the title of the video to which the passage belongs and the duration of the passage.
37 Finally, the user can play the passage via streaming.

38

39

40

11.6 Evaluation

41 Apart from the evaluation of each module contributing metadata annotations, assessing the
42 overall system response to user queries is of course of extreme importance for a multimedia
43 information retrieval (MIR) system. Suggestions and guidelines on evaluation procedures have
44 been put under testing in the framework of international contests. For example, the TREC 2002
45 Video track [22] included three sessions focusing on shot boundary detection, extraction of
46 segments containing semantic features like ‘People’, ‘Indoors Location’ etc., and search for
particular topics from a handcrafted list.

1 In what follows, a brief overview of the CIMWOS retrieval exercise is given. The CIMWOS
2 system and its overall performance has been evaluated in two distinct phases. During the
3 first phase, we tested video search and retrieval of passages relevant to a particular topic. We
4 repeated the retrieval task during the second phase, this time on video material for which
5 boundaries of relevant stories for each topic had been previously identified by human annota-
6 tors. Our testing material consisted of Greek news broadcasts, produced by state and private
7 TV networks between March 2002 and July 2003. For the first phase we used 15 videos
8 (henceforth, Collection A) of a total duration of approximately 18 hours, captured in BETA
9 SP and transcribed in MPEG-2 format. During the second phase we used 15 news broadcasts
10 (henceforth, Collection B) that amounted to approximately 17 hours of video, captured via
11 standard PC TV cards in MPEG-2 format.

12 The overall retrieval of the system was tested in both ‘interactive’ and ‘manual’ search
13 modes. In interactive search, users are familiar with the test collection, and have full access
14 to multiple interim search results. On the other hand, in the case of manual search, users with
15 knowledge of the query interface but no direct or indirect knowledge of the search test set
16 or search results are given the chance to translate each topic to what they believe to be the
17 most effective query for the system being tested. In each phase, a user familiar with the test
18 collection was responsible for generating a set of topics that s/he judged to be of interest, opting
19 for topics that were represented in more than one news broadcast of the collection. After the
20 list was finalized, the user had to manually locate all video sequences relevant to each topic,
21 thus allowing developers to associate start and end timecodes with each topic.

22 During the search task, users were given the chance to translate each topic of the list to what
23 they believed to be the most effective queries for the system being tested, using combinations
24 of criteria like Terms, Named Entities and/or ASR Text. Users formed five queries, on average,
25 for each topic. Using the HTML interface (Figure 11.10), they were able to browse the results
26 in order to intuitively assess their overall satisfaction with the database results. Nevertheless,
27 results were also saved in an XML format as in Figure 11.11, where the topic ‘earthquake
28 prediction’ has been translated into a query comprising the word *πρόγνωση* (*prediction*) and
29 a name (the name of a Greek geology professor).

30 In their initial reactions when searching for a topic in the videos of Collection A, users
31 reported that passages returned by the system were too short and fragmentary. The explanation
32 was that passages were based on automatic segmentation by the ASR engine, based on hints
33 from speaker changes. A different approach was taken while testing on Collection B, in order
34 to produce more intuitive results for end users.

35 A user was again responsible for manually identifying relevant segments in each of the
36 videos of the collection. These segments corresponded to the stories of each news broadcast.
37 Following that, we aligned the ASR transcription with the manually identified stories. Each
38 video segment was also assigned a description, derived from a list of topics created during
39 the training phase of the Greek ASR module. Thus, while it was possible to have passages
40 consisting of just one shot during the first phase, this counterintuitive phenomenon was absent
41 in the semi-automatically annotated material of Collection B.

42 The XML files corresponding to each user’s queries were assembled and tested against
43 groundtruth data. In Table 11.1, we show results for both video collections. We also tested the
44 system’s response when we did not take into account results that scored less than 60% in the
45 interface’s ranking system. This filtering had a negative effect on the system’s recall in the case
46 of Collection B. Nevertheless, it significantly increased the system’s precision with data from
both sets.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

The screenshot displays the CIMWOS HTML interface for a video segment. It is divided into several sections:

- Source Video Info:**

Creator	RTBF
Description	RTBF News Broadcast of 2003-08-05, 18:30
Creation Date	2003-08-05T18:30+00:00
Creation Location	Brussels
Rights Owner	RTBF
Total Video Duration	0:00:0
Availability Type	on demand delivery
Availability Start Date	2003-08-05
Availability End Date	2006-08-05
- Video Segment Info:**

Title	News of 2003-08-05 (English) Journal Télévisé 2003-08-05 (FRAN01)
Video Segment Duration	02:35 sec
Text	Douleur après le sang de pri pour plus attachant un peu moins pour les chroniqueurs de chaleur ici d'additionnement mieux vaut ne pas trop il pense ... très chaud pour très froid ... le travailleur défend pas . Pour l'un ils changent d' autre la grille-chauffe ... le bon boulot ... comment la climatisée . Au DVD avaient leur va se mettre au soleil le royume prend tout son futur . Quand même en des tares aujourd'hui apporter bonnet schtroupe en volume nous avait déjà de chercher des vêtements traite le nombre d' enfants français ont trouvé souvent au nouveau chef de tee-shirts autour d'un magasin de changer de ce que va dire nous regarde . Tant au nord au sud de que nous garde 13 de disant nous de leur métier ... que tous les conditions mêmes nouvelles détails créatives de il des surges à est un ancien balancier qui a connu aussi la la promesse des rêves de véritables passions les Dites ...
Named Entities	DVD soleil marc
Topics	Amerique du Nord Culture Economie
Video Text	Plan Tonello Marc Tonello RTBF Carré Responsable stock-surgelés Châteaux DVD RTBF 411 RTBF RTBF RTBF Dehors RTBF Simon Usage RTBF D Bardet.Parcoursage RTBF RTBF RTBF
Format	Streaming Video (Requires MS Media Player)
- Thumbnail:** A grid of small video frames showing various scenes from the broadcast.

Figure 11.10 CIMWOS HTML interface

11.7 Related Systems

Many multimedia indexing and retrieval systems have been developed over the past decade. Most of these systems aim at dynamically processing, indexing, organizing and archiving digital content [23, 24]. A brief overview of the most relevant systems is given below.

News on demand (NOD) systems have made their appearance in the past decade, integrating image, speech and language technology, and taking advantage of cross-modal analysis in order to monitor news from TV, radio and text sources and provide personalized newscasts [13].

MIT's Photobook [25] was an early example that allowed retrieval from textures, shapes and human faces. Chabot [26], developed by UC Berkeley, is another early system. It provided a combination of text-based and colour-based access to a collection of photographs. The system was later renamed Cypress, and is now known as the DWR picture retrieval system [27], and has been incorporated into the Berkeley Digital Library project. Berkeley has continued similar research with its Blobworld software [28], which has one of the more sophisticated region segmentation functionalities. Noteworthy are recent developments at Berkeley, which allow finding of specific object classes like horses and naked people. VisualSEEK [29] was the first of a series of systems developed at Columbia University. A strong point was that queries were quite flexible in the spatial layout specifications that they could accept. WebSEEK [30] was a further development, aimed at facilitating queries over the web. Emphasis is on colour, but in another prototype called VideoQ [31], motion was added. Also, relative layouts of regions with specific colours or textures can be specified. This allows for some primitive

```

1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <criteria Language="Greek" Text="πρόγνωση"
4     Named_Entities="Παπαζάχος"
5     Topics="γεωλογία,μετεωρολογία"
6     Logical_Operator="OR" />
7   <passage score="100%" number="1"
8     media_title="News of 2002-05-24"
9     passage_id="p44"startpoint_in_sec="529,05"
10    duration_in_sec="15,76">
11     <content>αναστάτωση αλλά και ανησυχίες
12       προκαλεί στην κοινή γνώμη δημοσιοποίηση
13       έκθεσης του Βασίλη του καθηγητή Βασίλη
14       Παπαζάχου για μεσοπρόθεσμη πρόγνωση
15       ισχυρός σεισμός σε τέσσερις περιοχές της
16       Ελλάδας ο κύριος Παπαζάχος με δηλώσεις
17       του ο υπερασπίστηκε την επιστημονική
18       πέραση.</content>
19     <speaker>male 22</speaker>
20     <media id="NET_20020524_1800_News"
21       creation_date="2002-05-24T18:00+02:00"
22       creator="NET" rights_owner="NET"
23       availability_type="ondemand delivery">
24       <mediatitle>News of 2002-05-24
25         </mediatitle>
26       <source duration="1:74:52,48" />
27     </media>
28     <thumbnails number="1" />
29   </passage>
30   ...
31 </results>

```

Figure 11.11 Query results in XML

forms of object detection, e.g. finding the American flag on the basis of the alternation of red and white stripes, with a blue region at the top left corner. The MARS project at the University of Illinois [32] put emphasis on control by the user, and introduced extensive relevance feedback mechanisms. Feature weights are set dynamically. An example of European

Table 11.1 Retrieval results on Greek video collections

	Precision	Recall	F-measure
Collection A	34.75	57.75	43.39
Collection A + 60% filter	45.78	53.52	49.35
Collection B	44.78	50.24	47.36
Collection B + 60% filter	64.96	37.07	47.20

1 CBIR technology is the Surfimage system from INRIA [33]. It has a similar philosophy to
2 the MARS system, using multiple types of image features which can be combined in different
3 ways, and offering sophisticated relevance feedback. The RSIA system co-developed by ETH
4 Zurich and DLR [34] was dedicated to search in large satellite image databases. It supports
5 remote access through the web and is best on texture features, learned as the user indicates
6 positive and negative examples of the kind of regions searched for. The Netra system [35] uses
7 similar cues to many of the foregoing systems, i.e. colour, texture, shape and spatial location
8 information, but as in Blobworld a more sophisticated segmentation function is included.
9 Synapse [36] is a rather different kind of system, in that it bases its similarity judgements on
10 whole image matching. Query By Image Content [37] is a component of the IBM DB2 database
11 product, performing image retrieval based on example images and selected colour and texture
12 patterns.

13 Convera [38] provides scene change detection and can send audio tracks to any SAPI-
14 compliant recognition engine. It also provides a closed captioning/teletext extractor and an
15 SMPTE timecode [39] reading module, but these are more 'decoders' than true 'recognition
16 engines', as this information is already available in digital form. On its part, Virage [40] claims
17 to be more comprehensive, including face recognition, on-screen text recognition, speech and
18 speaker, as well as intelligent segmentation.

19
20

21 **11.8 Conclusion**

22 Multimedia information retrieval systems address real needs of multimedia producers,
23 archivists and others who need to monitor and/or index audiovisual content. By utilizing vast
24 amounts of information accumulated in audio and video, and by employing state-of-the-art
25 speech, image and text processing technologies, an MIR system can become an invaluable
26 assistant in efficient reuse of multimedia resources, and in reducing new production expenses.

27 These systems open new possibilities and provide enabling technology for novel application
28 areas and services. Use of markup languages (MPEG-7, MPEG-21) for annotation of audiovi-
29 sual content can promote standardization and unification of access procedures to large archives
30 and multiple repositories, thus assisting the rise of multimedia digital libraries, and improving
31 the working conditions and productivity of people involved in media and television, video,
32 news broadcasting, show business, advertisement and any organization that produces, markets
33 and/or broadcasts video and audio programmes.

34
35

36 **Acknowledgements**

37 Project CIMWOS was supported by shared-cost research and technological development con-
38 tract IST-1999-12203 with the European Commission. Project work was done at the Institute
39 for Language & Speech Processing, Athens, Greece (coordination, integration, text process-
40 ing, search and retrieval), Katholieke Universiteit Leuven, Belgium (video segmentation, face
41 detection), Eidgenoessische Technische Hochschule Zurich, Switzerland (scene and object
42 recognition), Sail Labs Technology AG, Vienna, Austria (audio processing, speech recogni-
43 tion, text processing), Canal+ Belgique, Brussels, Belgium (user requirements and validation)
44 and Institut Dalle Molle d'Intelligence Artificielle Perceptive, Martigny, Switzerland (speech
45 recognition, video text detection and recognition, face recognition).
46

References

- 1 [1] H. Papageorgiou, A. Protopapas, CIMWOS: a multimedia, multimodal and multilingual indexing and retrieval
2 system. In E. Izquierdo (ed.) *Digital Media Processing for Multimedia Interactive Services, Proceedings of the*
3 *4th European Workshop on Image Analysis for Multimedia Interactive Services*, Queen Mary, University of
4 London, 9–11 April 2003, pp. 563–568. World Scientific, Singapore, 2003.
- 5 [2] Y. Wang, Z. Liu, J.-C. Huang, Multimedia content analysis: using both audio and visual clues. *IEEE Signal*
6 *Processing Magazine*, **17**(6), 12–36, 2000.
- 7 [3] Z. Liu, Y. Wang, T. Chen, Audio feature extraction and analysis for scene segmentation and classification. *Journal*
8 *of VLSI Signal Processing Systems for Signal, Image, and Video Technology*, **20**(1), 61–79, 1998.
- 9 [4] J.-L. Gauvain, R. De Mori, L. Lamel, Advances in large vocabulary speech recognition. *Computer Speech and*
10 *Language*, **16**(1), 1–3, 2002.
- 11 [5] J.-L. Gauvain, L. Lamel, G. Adda, Audio partitioning and transcription for broadcast data indexation. *Multimedia*
12 *Tools and Applications*, **14**, 187–200, 2001.
- 13 [6] F. Kubala, J. Davenport, H. Jin, D. Liu, T. Leek, S. Matsoukas, D. Miller, L. Nguyen, F. Richardson,
14 R. Schwartz, J. Makhoul, The 1997 BBN BYBLOS System applied to Broadcast News Transcription. In
15 *Proceedings of the DARPA Broadcast News Transcription and Understanding Workshop*, Lansdowne, VA,
16 8–11 February 1998. National Institute of Standards and Technology, Gaithersburg, MD, 1998. Available at:
17 <http://www.nist.gov/speech/publications/darpa98/pdf/eng110.pdf>.
- 18 [7] J.L. Gauvain, L. Lamel, G. Adda, Transcribing broadcast news for audio and video indexing. *Communications*
19 *of the ACM*, **43**(2), 64–70, 2000.
- 20 [8] E.F. Tjong Kim Sang, Introduction to the CoNLL-2002 shared task: language-independent named entity recognition.
21 In D. Roth, A. van den Bosch (eds) *CoNLL-2002, Sixth Conference on Natural Language Learning*, Taipei,
22 Taiwan, 31 August –1 September 2002. Available at: <http://cnts.uia.ac.be/conll2002/proceedings.html>.
- 23 [9] I. Demiros, S. Boutsis, V. Giouli, M. Liakata, H. Papageorgiou, S. Piperidis, Named entity recognition in Greek
24 texts. In *Proceedings of the Second International Conference on Language Resources and Evaluation, LREC2000*,
25 Athens, Greece, 31 May–2 June 2000, pp. 1223–1228. ELRA, 2000.
- 26 [10] C. Jacquemin, E. Tzoukermann, NLP for term variant extraction: synergy between morphology, lexicon and
27 syntax. In T. Strzalkowski (ed.) *Natural Language Information Retrieval*. Kluwer, Dordrecht, 1999.
- 28 [11] A. Hauptmann, R. Smith, Text, speech, and vision for video segmentation: the Informedia Project. In *AAAI*
29 *Fall 1995 Symposium on Computational Models for Integrating Language and Vision*. American Association for
30 Artificial Intelligence, 1995.
- 31 [12] SRI Maestro Team, MAESTRO: conductor of multimedia analysis technologies. *Communications of the ACM*,
32 **43**(2), 57–63, 2000.
- 33 [13] A. Merlino, D. Morey, M. Maybury, Broadcast news navigation using story segments. In *Proceedings of the Fifth*
34 *ACM International Conference on Multimedia*, Seattle, WA, 9–13 November 1997, pp. 381–391. ACM, Seattle,
35 WA, 1997.
- 36 [14] F. Cardinaux, C. Sanderson, S. Marcel, Comparison of MLP and GMM classifiers for face verification on
37 XM2VTS. In J. Kittler, M.S. Nixon (eds) *Proceedings of the 4th International Conference on Audio- and Video-*
38 *Based Biometric Person Authentication (AVBPA)*, Guildford, 9–11 June 2003. Springer-Verlag, Berlin, 2003.
- 39 [15] The Combined MIT/CMU Test Set with Ground Truth for Frontal Face Detection, [http://vasc.ri.cmu.edu/idb/html/](http://vasc.ri.cmu.edu/idb/html/face/frontal_images/index.html)
40 [face/frontal_images/index.html](http://vasc.ri.cmu.edu/idb/html/face/frontal_images/index.html).
- 41 [16] D. Chen, H. Bourlard, J.-Ph. Thiran, Text identification in complex background using SVM. In *Proceedings of*
42 *the International Conference on Computer Vision and Pattern Recognition*, Kauai, Hawaii, December 2001, vol.
43 2, pp. 621–626. IEEE, 2001.
- 44 [17] J.M. Odobez, D. Chen, Robust video text segmentation and recognition with multiple hypotheses. In *Proceedings*
45 *of the International Conference on Image Processing, ICIP 2002*. IEEE, 2002.
- 46 [18] V. Ferrari, T. Tuytelaars, L. Van Gool, Wide-baseline multiple-view correspondences. In *Proceedings of the*
International Conference on Computer Vision and Pattern Recognition ICIP 2003, Madison, WI, June 2003.
IEEE, 2003.
- [19] T. Tuytelaars, A. Zaatri, L. Van Gool, H. Van Brussel, Automatic object recognition as part of an integrated
supervisory control system. In *Proceedings of IEEE Conference on Robotics and Automation, ICRA00*, San
Francisco, CA, April 2000, pp. 3707–3712. IEEE, 2000.
- [20] V. Ferrari, T. Tuytelaars, L. Van Gool, Real-time affine region tracking and coplanar grouping. In *Proceedings*
of the International Conference on Computer Vision and Pattern Recognition, Kauai, Hawaii, December 2001,
vol. 2. IEEE, 2001.

- 1 [21] H. Papageorgiou, A. Protopapas, T. Netoušek, Retrieving video segments based on combined text, speech and
2 image processing. In *Proceedings of the Broadcast Engineering Conference*, April 2003, pp. 177–182. National
3 Association of Broadcasters, 2003.
- 4 [22] A.F. Smeaton, Paul Over, The TREC-2002 Video track report. In E.M. Voorhees, L.P. Buckland (eds) *The*
5 *Eleventh TExt Retrieval Conference (TREC 2002)*, Gaithersburg, MD, 19–22 November 2002. NIST Special
6 Publication 500-251. National Institute of Standards and Technology, Gaithersburg, MD, 2002. Available at:
7 http://trec.nist.gov/pubs/trec11/t11_proceedings.html.
- 8 [23] N. Dimitrova, H.-J. Zhang, B. Shahraray, I. Sezan, T. Huang, A. Zakhor, Applications of video-content analysis
9 and retrieval. *IEEE Multimedia*, **9**(3), 42–55, 2002.
- 10 [24] *IEEE Multimedia*, Special issue: Content-Based Multimedia Indexing and Retrieval, **9**(2), 18–60, 2002.
- 11 [25] A. Pentland, R.W. Picard, S. Sclaroff, Photobook: content-based manipulation of image databases. *International*
12 *Journal of Computer Vision*, **18**(3), 233–254, 1996.
- 13 [26] V.E. Ogle, M. Stonebraker, Chabot: retrieval from a relational database of images. *IEEE Computer*, **28**(9),
14 164–190, 1995.
- 15 [27] DWR photo database, <http://elib.cs.berkeley.edu/photos/dwr/about.html>.
- 16 [28] Blobworld, <http://elib.cs.berkeley.edu/photos/blobworld>.
- 17 [29] J.R. Smith, S.-F. Chang, Tools and techniques for color image retrieval. In I.K. Sethi, R.C. Jain (eds) *Storage*
18 *and Retrieval for Still Image and Video Databases IV*, SPIE vol. 2670, pp. 426–437. SPIE, 1996.
- 19 [30] J.R. Smith and S.-F. Chang, Image and video search engine for the world wide web. In I.K. Sethi, R.C. Jain (eds)
20 *Storage and Retrieval for Image and Video Databases V*, SPIE vol. 3022, pp. 84–95. SPIE, 1997.
- 21 [31] S.F. Chang, W. Chen, H.J. Meng, H. Sundaram, D. Zhong, VideoQ: an automated content-based video search
22 system using visual cues. In *Proceedings of International Multimedia Conference*, Seattle, WA, November 1997,
23 pp. 313–324. Addison-Wesley, Reading, MA, 1997.
- 24 [32] T.S. Huang, S. Mehrotra, K. Ramchandran, Multimedia Analysis and Retrieval System (MARS) project. In P.B.
25 Heidorn, B. Sandore (eds) *Proceedings of the 33rd Annual Clinic on Library Application of Data Processing:*
26 *Digital Image Access and Retrieval*, Urbana, IL, March 1996, pp. 100–117. University of Illinois, 1997.
- 27 [33] C. Nastar, M. Mitschke, C. Meilhac, N. Boujemaa, Surfimage: a flexible content-based image retrieval system.
28 In *Proceedings of International Multimedia Conference*, Bristol, 12–16 September 1998, pp. 339–344. Addison-
29 Wesley, Harlow, 1998.
- 30 [34] K. Seidel, M. Datcu, G. Schwarz, L. Van Gool, Advanced remote sensing information system at ETH Zurich
31 and DLR/DFD Oberpfaffenhofen. In IEEE International Geoscience and Remote Sensing Symposium IGARSS
32 '99, Hamburg, Germany, pp. 2363–2365. IEEE, 1999.
- 33 [35] W.Y. Ma, B.S. Manjunath, NeTra: a toolbox for navigating large image databases. In *Proceedings of International*
34 *Conference on Image Processing*, Santa Barbara, CA, October 1997, vol. 1, pp. 568–571. IEEE, 1997.
- 35 [36] R. Manmatha, S. Ravela, Syntactic characterization of appearance and its application to image retrieval. In B.E.
36 Rogowitz, T.N. Pappas (eds) *Human Vision and Electronic Imaging II*, SPIE vol. 3016, pp. 484–495. SPIE, 1997.
- 37 [37] Query By Image Content, <http://wwwqbic.almaden.ibm.com>.
- 38 [38] Convera, <http://www.convera.com>.
- 39 [39] Society of Motion Picture and Television Engineers, <http://www.smpte.org>.
- 40 [40] Virage, <http://www.virage.com>.
- 41
- 42
- 43
- 44
- 45
- 46