



# Indexing for reuse of TV news shots

M. Bertini \*, A. Del Bimbo, P. Pala

Dipartimento Sistemi E Informatica, Universita' di Firenze, Via S Marta 3, 50139 Firenze, Italy

Received 16 November 2000; accepted 16 November 2000

## Abstract

Broadcasters are demonstrating interest in building digital archives of their assets for reuse of archive materials for TV programs or on-line availability. This requires tools for video indexing and retrieval by content. Effective indexing by content of videos is based on the association of high-level information associated with visual data. In this paper a system is presented that enables content-based indexing and browsing of news reports; the annotation of the video stream is fully automated and is based both on visual features extracted from video shots and on textual descriptors extracted from captions and audio tracks. © 2001 Published by Elsevier Science Ltd on behalf of Pattern Recognition Society.

Keywords: Multimedia databases; Video content analysis; Content-based video retrieval; Video shots classification

## 1. Introduction

Broadcasters are demonstrating interest in building large digital archives of their assets for reuse of archive materials for TV programs or on-line availability to other companies and the general public. To satisfy this request there is need of systems that are able to provide efficient management of visual data in terms of storage, transmission, retrieval and browsing. Solutions to storage and transmission issues involve analysis and processing of data streams regardless of their content. Differently, effective retrieval and browsing of images and videos is based on the extraction of content level information associated with visual data and on a compact representation of retrieved shots.

While effective content-based retrieval of images is accomplished by supporting content representation through low-level image features, the same does not apply to content-based retrieval of videos, except for very limited application contexts. Instead, effective retrieval of videos must be based on high-level content descriptors.

Specific knowledge of the application content ease the extraction of high-level descriptors [1].

Recently, news videos have received great attention by the research community. This is motivated by the interest of broadcasters in building digital archives of their assets for reuse of archive materials. On the one hand, reuse of archive materials is identified as one key method of improving production quality by bringing added depth, and historical context, to recent events. On the other hand, the use of stock footage allows to produce faster the news services. An example of the first case is the reuse of shots that show the scene of a crime: they can be reused later to provide the historical context. An example of the second case is the reuse of "generic" shots, e.g. shots that show an airport may be used in a news service about an airport strike. Anyway, it is not possible to reuse all the shots of a news video: the information contained in the speech of an anchorman or in the text and the graphs of a computer graphics shot became obsolete after a short time and can be easily and inexpensively replaced by new shots. An effective reuse of archive materials is possible if the shot description is rich enough to allow retrieval by content and the content has been classified: a thorough description of the contents allows to search the shots that fit into a request of the video producer, while shot classification allows to skip those that cannot

\* Corresponding author. Tel.: +39-055-4796540;

fax: +39-055-4796363.

E-mail addresses: bertini@dsi.unifi.it (M. Bertini), delbimbo@dsi.unifi.it (A. Del Bimbo), pala@dsi.unifi.it (P. Pala).

1 be reused. News have a rather definite structure and do  
 2 not offer a wide variety either of edit effects, which are  
 3 mainly cuts, or of shooting condition (e.g. illumination).  
 4 The definite structure of news is suitable for content anal-  
 5 ysis and has been exploited for automatic classification of  
 6 news sequences in Refs. [2–6]. In all of these systems a  
 7 two stage scene classification scheme is employed. First,  
 8 the video stream is parsed and video shots are extracted.  
 9 Each shot is then classified according to content classes  
 10 such as *newscaster*, *report*, *computer graphics*, *weather*  
 11 *forecast*. The general approach to this type of classifica-  
 12 tion relies on the definition of one or more image tem-  
 13 plates for each content class. To classify a generic shot,  
 14 a *key frame* is extracted and matched against the im-  
 15 age template of every content classes. Other works [6,7]  
 16 deal with the problem of video indexing using informa-  
 17 tion sources like the text of the captions and the audio  
 18 track.

19 This is due to the fact that news videos images have an  
 20 ancillary function with respect to words and video content  
 21 is strongly related to textual and audio information which  
 22 is contained in the audio track.

### 23 1.1. Previous work

24 A method for shot classification based on the syntax  
 25 and the structure of news videos has been proposed in  
 26 Ref. [2]. Shot classification is based on the similarity  
 27 match of frames against a pre-determined set of proto-  
 28 type anchorman images. However, as noted in Refs. [3,5]  
 29 the validity of this approach is limited by the difficulty  
 30 to find a representative set of prototype anchorman im-  
 31 ages. These should account for different cases including  
 32 news editions, change of dresses, modifications of stu-  
 33 dio layout. Furthermore, the method proposed is compu-  
 34 tation intensive since it requires the calculation of simi-  
 35 larity between each frame and prototypical image of the  
 36 anchorman.

37 In order to diminish dependency from the set of sam-  
 38 ple frames in Refs. [3,5] has been proposed to use a  
 39 different approach to the definition of the anchorman  
 40 frame model. In this model each anchorman frame is  
 41 considered a composition of distinctive regions, like the  
 42 shape of the anchorman, the caption of the reporter’s  
 43 name, the graphics that sometimes appear in the top  
 44 third of the frame. A model of the anchorman frame  
 45 is built, which accounts for the spatial distribution of  
 46 basic elements and is independent of the anchorman’s  
 47 sex, apparel and appearance. To determine whether a  
 48 shot contains an anchorman all the frames are compared  
 49 with the model; if they match, they are classified as  
 50 “anchorman”, thus building a set of model images for  
 51 each video. Only the frames of the shots that satisfy the  
 52 similarity criteria according to the spatial model are then  
 53 compared with the model-image set, using a new simi-  
 larity measure. One of the limits of this method is that

if the style of the news changes the database must be 55  
 updated.

56 A different approach, based on frame statistics, is 57  
 presented in Ref. [8]. The system uses hidden Markov 58  
 models to classify frames of news videos. The classifica- 59  
 tion process takes into account several clues, including 60  
 feature vectors based on difference images, average 61  
 frame color and audio signal. Parameters of the hidden 62  
 Markov model are determined in a training stage using 63  
 a ground-truth database of news videos.

64 The problem of text extraction has been investigated 65  
 by several researchers. A method for the extraction of 66  
 captions and scene texts (e.g. street names or shop names 67  
 in the scene) from movies has been presented in Ref. 68  
 [9]. Techniques for the extraction and OCR of caption 69  
 text for the news video indexing have been examined in 70  
 Ref. [7]. The first problem that must be solved for effec- 71  
 tive text extraction is to determine which frames contain 72  
 captions and the position of the text in the frame. The 73  
 method presented in Ref. [7] is based on the search of 74  
 rectangular regions, composed by elements with sharp 75  
 borders, appearing in sequences of frames; it is also 76  
 based on the assumption that the captions have a high 77  
 contrast on the background.

78 For the purpose of video content annotation, speech 79  
 transcriptions has been used in the CMU Informedia 80  
 project as extremely important source of information 81  
 [10–12].

82 News-on-demand is an application within the Infor- 83  
 media digital video library project [6] that indexes news 84  
 from TV and radio sources and allows the user to retrieve 85  
 news by content. The system creates a time-aligned 86  
 transcript from speech recognition and captions. The 87  
 video data is segmented into news stories using the 88  
 presence of silence and captions as “paragraph” bound- 89  
 aries, while scene breaks and keyframes are identified 90  
 using algorithms based on color histograms. The CMU 91  
 Sphinx-II speech recognition system is used both for 92  
 the speech transcription and for the user interface of the 93  
 content-based retrieval system. There is no shot classifi- 94  
 cation and the speech recognition system uses the whole 95  
 audio track, obtaining variable error rates that depend 96  
 on the audio source [11]. 97

98 Two prototypes for the construction of personalized 99  
 TV news programs have been presented in Ref. [13]. 100  
 The first prototype allows category-based retrieval using 101  
 manual annotation provided by the news producer. The 102  
 second prototype indexes the shot content using teletext 103  
 data that are provided for deaf people by a French TV 104  
 channel. The indexing of news videos uses the video 105  
 parsing system presented in Ref. [14]. This system de- 106  
 tects cuts computing the difference of color histograms 107  
 of consecutive frames. Shots containing anchorman are 108  
 identified by combining shot similarity, person detection 109  
 and the “high variance factor” which accounts for the 110  
 “regular spot presence” of the anchorman shots.

## 1.2. The news indexing and annotation system

In this paper a system for content-based indexing and annotation of news videos is presented. Videos are segmented to identify video shots. On the basis of the first frame of each shot, a statistical analysis is performed to detect which shots recur throughout the video. The shots are thus classified as newscaster shots, and the others are classified as report shots. The content of a generic report shot is described through the use of both visual and textual information and is further classified as computer graphics (*non-realistic*) or *realistic*, in order to improve the reuse of realistic shots, as needed by the broadcasters. Textual information is automatically extracted from textual captions included in the video and from speech associated with the video. Differently from Ref. [6] only anchorman shots are used for speech recognition. A retrieval engine allows the user to search by content and browse through video shots.

This paper is organized as follows: in Section 2, the video segmentation technique used to identify video shots is presented. In Section 3 the shot classification system is presented, and a comparison is carried out with respect to other techniques. In Section 4, video content description is expounded with reference to the extraction of textual information from OCR and speech recognition. Finally, in Section 5 retrieval and browsing examples are provided.

## 2. Video segmentation

In order to perform segmentation of news videos two problems must be dealt with: (i) avoiding incorrect identification of shot changes due to rapid motion or sudden lighting change in the scene (false positives), (ii) identification of sharp shot transitions (cuts) as well as gradual (dissolves, matte). Ref. [15] reports a thorough comparison of video segmentation algorithms. In the following, we concentrate on cuts since they are, by far, the most commonly employed edit effect in news videos. Furthermore, for the purpose of content-based indexing, it is not important to classify the edit effect, but to detect changes of visual content. Table 1 shows the number of sharp and gradual edit effects used in 4 h of news videos of the three most important Italian broadcasters.

The identification of gradual as well as sharp transitions can be performed through a cut detection algorithm,

provided that the video is suitably sub-sampled in time. In fact, gradual transitions become sharp if the video is sub-sampled in the time variable since the difference between consecutive frames increases. The cut detection algorithm is developed following two distinct steps:

*Preliminary cut detection:* Rapid motion in the scene and sudden change in lighting produce a low correlation between contiguous frames especially in case a high temporal sub-sampling rate is adopted. To avoid false cut detection, a metric has been studied which proves highly insensitive to such variations, while being reliable in detecting “true” cuts [16]. Each frame is partitioned into nine sub-frames. Each of these is represented by considering its color histogram in the HSI color space. Actually, to improve independence with respect to lighting conditions, the histogram takes into account only hue  $H$  and saturation  $S$  properties. The HSI color space has been chosen, since as reported in Ref. [17], it is a good compromise between missed detection and computational costs.

Edit effect detection is performed considering the volume of the difference of sub-frame histograms in two consecutive frames. Cuts correspond to zero crossings of the difference of the average values of the difference of the volumes. This method allows edit effect identification also when the frame color statistic remains the same but the position of the color spots is different.

To keep false positive detection low, results of the first pass are refined using a method based on video structure and shot similarity.

*Cut detection refinement:* The algorithm described above features a high false positive detection rate in some critical situations, such as: (i) color instabilities due to noise in the digitalization process, (ii) insertion of graphics or other changes of large zones in images, (iii) news shots recorded in critical situations, or news shots featuring sudden lighting changes. Typically, lighting changes are due either to long sequences of flashes like in press conferences, or to sudden camera movements (like *panning* and *zoom*) and free hand takes, like in reports on demonstrations or war actions.

To reduce errors due to multiple and rapid variation of visual contents of the shot, the knowledge of the specific structure of news videos has been considered. In fact, unlike other types of videos, such as commercials and movies, where the editing can reach frantic levels, in news videos the duration of the shots is long enough to let the audience “understand” the subject. Thus, there is always a minimum temporal distance  $\tau_L$  between two consecutive cuts. This rule is adopted to disregard all those cuts that are less distant than  $\tau_L$  seconds from the preceding cut (inter-cut time difference constraint). Furthermore, since cuts identify a change of the video content the key-frames of shots for two consecutive cuts cannot be too similar. This rule is used to disregard all those cuts whose similarity with the preceding cut exceeds a threshold  $\tau_S$  (inter-cut frame similarity constraint).

Table 1  
Shot boundary statistics for news videos

Shots	Cuts	Diss.+ Wipe	Matte+
1797	1702 (94.7%)	95 (5.3%)	

Table 2  
Statistics of all the videos

Shots	Detected shots	False detections	Missed detections
731	765	43 (5.9%)	9 (1.2%)

The performance of the proposed technique has been evaluated with reference to a test database composed of 12 videos from 6 Italian TV channels: RAI 1, 2 e 3, Mediaset Canale 5 and Cecchi Gori TeleMonteCarlo 1, for a total time of 2 h and 42 min.

Table 2 includes the number of video shots, cuts, gradual edit effects, falsely detected edit effects and missed detections. With respect to cut detection based exclusively on color histogram, the use of cut detection refinement results in a 37% improvement in false detection (from 69 to 43).

### 3. Shot classification

The main goals of shot classification are the classification of reusable and not reusable shots, and the indexing of the video. For each video shot, the first frame is used as the key-frame. Video shots are classified into two main classes: anchorman and news reports. Sub-classifications of the anchorman shots (like “weather forecasts”) are obtained considering the speech content as explained in Section 4. Shot classification is a two step process: the first step classifies anchorman and report shots, using a statistical approach and motion features of the anchorman shots, without requiring any model. Then news report shots are processed in order to detect those that contain computer graphics.

Classification of anchorman and computer graphic shots is important since they cannot be reused. Fig. 1 shows an example of reusable shots: the anchorman introduces a report about accidents in the home, then after some realistic shots that show typical house works there is a computer graphic shot that will show some statistics. While the realistic shots are reuseable in an another report that deals about house works, the anchorman and the graphics are not, and will be replaced by another anchorman and by newer computer graphics.

#### 3.1. Classification of anchorman shots

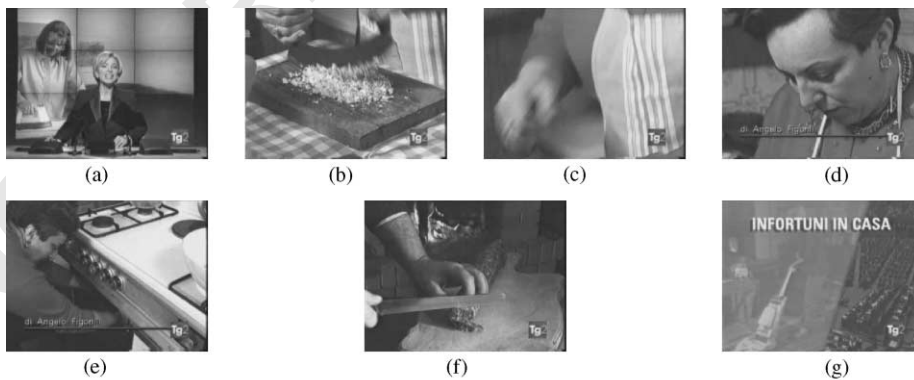
##### 3.1.1. Classification based on statistical features

Shots of the anchorman are repeated at intervals of variable length throughout the video. The first step for the classification of these shots stems from this assumption and is based on the computation, for each video shot  $S_k$ , of its *shot lifetime*  $\mathcal{L}(S_k)$ . The shot lifetime measures the shortest temporal interval that includes all the occurrences of shots with similar visual content, within the video. Given a generic shot  $S_k$  its lifetime is computed by considering the set

$$T_k = \{t_i | \mathcal{S}(S_k, S_i) < \tau_s\},$$

where  $\mathcal{S}(S_k, S_i)$  is a similarity measure applied to key-frames of shots  $S_k$  and  $S_i$ ,  $\tau_s$  a similarity threshold and  $t_i$  is the value of the time variable corresponding to the occurrence of the key-frame of shot  $S_i$ . The lifetime of shot  $S_k$  is defined as  $\mathcal{L}(S_k) = \max(T_k) - \min(T_k)$ .

Shot classification is based on fitting values of  $\mathcal{L}(S_k)$  for all the video shots in a bimodal distribution. This is used to identify a threshold value  $\tau_l$  that is used to classify shots into service and anchorman categories. Particularly, all the shots  $S_k$  so that  $\mathcal{L}(S_k) > \tau_l$  are classified as anchorman shots, where  $\tau_l$  is determined according



Anchorman, report and CG key-frame sequence

Fig. 1. Example of reusable (b, c, d, e, f,) and not reusable(a, g) shots.

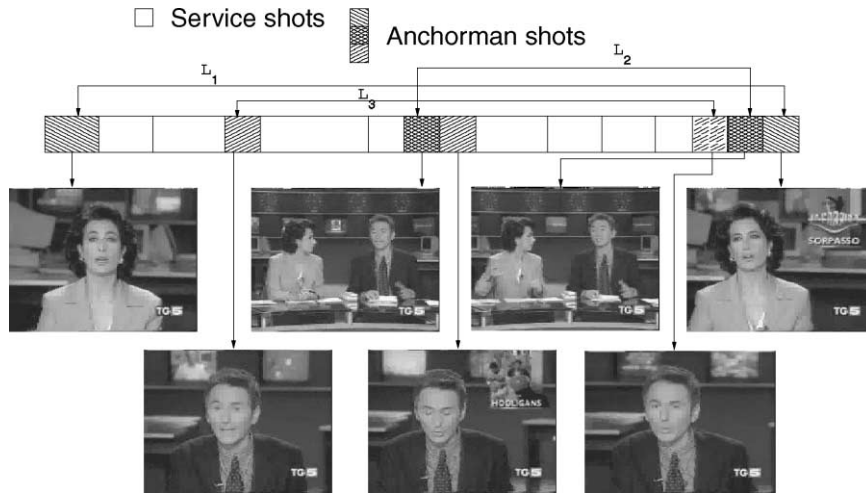


Fig. 2. Lifetime of anchorman shots.

1 to the statistics of the test database, and set to 4.5 s. Re-  
 2 remaining shots are classified as news service shots.

3 This classification method does not rely on any  
 4 pre-defined model of the anchorman shots; rather it is  
 5 based on the time structure of news videos. Fig. 2 shows  
 6 lifetimes for three different types of anchorman shots  
 7 identified in a news video.

### 3.1.2. Classification based on motion features

9 Shot classification based on statistical features can  
 10 sometimes lead to the erroneous classification of some  
 11 news service shots as anchorman shots. This occurs  
 12 mainly in correspondence to interviews and reports. In  
 13 fact, in

14 *Interviews:* The camera alternatively takes shots of the  
 15 interviewer and the interviewed people. Erroneous clas-  
 16 sification of interview shots have been discussed in Ref.  
 17 [8].

18 *In reports:* A reporter describes the content shown  
 19 in some shots; at the end of every shot (or series of  
 20 shots) there is the shot of a new reporter describing the  
 21 next series; this structure replicates the whole structure  
 22 of news video. An example is shown in Fig. 3 where the  
 23 recurrence of shots (a), (c) and (g) leads to erroneous  
 24 classification of these shots as anchorman shots.

25 To avoid these errors, the preliminary classifica-  
 26 tion based on statistical feature is refined considering  
 27 motion features of the anchorman shots. Classifica-  
 28 tion refinement stems from the assumption that in an  
 29 anchorman shot, both the camera and the anchorman  
 30 are almost motionless. In contrast, for both interview  
 31 and news service shots, background objects and cam-  
 32 era movements—persons and vehicles, free-hand shots,  
 33 camera panning and zooming—cause relevant motion  
 components throughout the shot.

Classification refinement is performed by computing  
 an index of the *quantity of motion*  $\mathcal{Q}_S$ , for each possible  
 anchorman shot. The algorithm for the analysis of this  
 index takes into account the frame to frame difference  
 between the shot key-frame  $f_1$  and subsequent frames  
 $f_i$  in the shot according to

$$\mathcal{Q}_S = \sum_{f_i \in \mathcal{S}} D_i$$

with

$$D_i = \sum_{xy} d_{RGB}(f_1(x, y), f_i(x, y)), \quad (1)$$

$$d_{RGB}(f_1(x, y), f_i(x, y))$$

$$= \begin{cases} 0 & \text{if } \|f_1(x, y) - f_i(x, y)\| < \tau_{RGB}, \\ 1 & \text{if } \|f_1(x, y) - f_i(x, y)\| \geq \tau_{RGB}. \end{cases} \quad (2)$$

To enhance sensitivity to motion the shot is  
 sub-sampled in time, and the frames are compared to  
 the key-frame  $f_1$ . Only those shots whose  $\mathcal{Q}_S$  does not  
 exceed a threshold  $\tau_Q$  are definitely classified as an-  
 chorman shots. By using this classification refinement,  
 false anchorman shots shown in Fig. 3 are eliminated.  
 In fact, shots (a), (c) and (g) feature a relevant mo-  
 tion component on account of camera zooming and  
 panning and movement of people and objects in the  
 background.

### 3.2. Classification of computer graphics shots

Shots classified as containing news report are pro-  
 cesses in order to detect whether they contain computer  
 graphics. Fig. 4 shows an example of computer graph-  
 ics shot. Usually, those type of shots show information  
 about money change rates, economic indexes and other

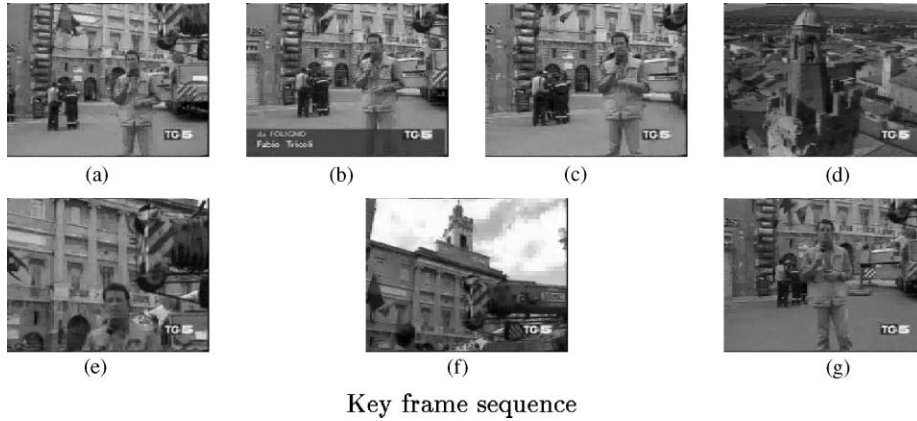
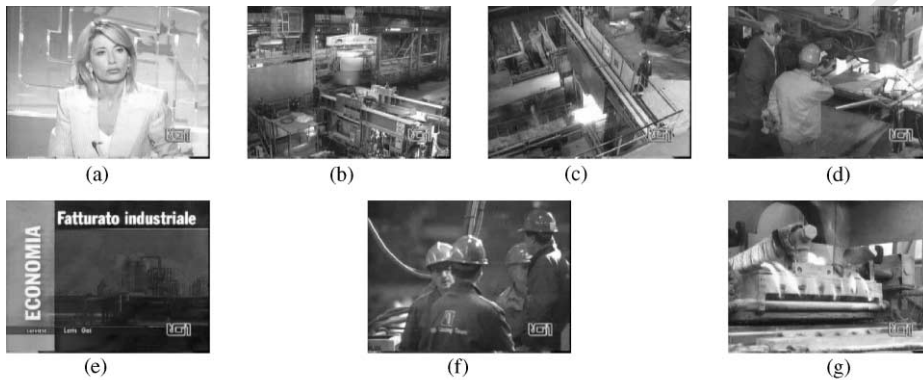


Fig. 3. Example of “false positive” class.



Anchorman, report and CG key-frame sequence

Fig. 4. Example of reusable (report) and not reusable (anchorman and CG) shots.

1 graphs. They are not reusable due to the fact that the  
 2 information they convey is subject to fast changes, and  
 3 can be inexpensively replaced. The shot represented by  
 4 key-frame (e) in Figs. 4 and 5 shows the sales of Febru-  
 5 ary 2000 compared to those of February 1999, and has  
 6 little reuse value. Unlike it, the shots that show workers  
 7 in a factory can be reused in other reports. Unlike the  
 8 anchorman shots it is possible to use neither the struc-  
 9 ture of the video nor the layout, as a hint to detect the  
 10 computer graphic shots.

11 The features used to classify those shots are based on  
 12 statistical parameters and motion features. The first step  
 13 calculates an index of the *quantity of motion*  $\mathcal{Q}_{CG}$  divid-  
 14 ing each shot into sub-shots and taking into account the  
 15 frame to frame difference between the sub-shot key-frame  
 16  $f_i$  and  $f_{(i+1)}$  according to the previous equation.

17 To reduce possible misclassification of still images the  
 18 preliminary classification is refined analyzing the color  
 19 histogram in the HSI color space. The histogram takes  
 20 into account only the H and S components, and calcula-  
 21 tes two indexes:  $N_{bin}$  is the number of histogram bins

whose value is higher than a  $\tau_{bin}$  percentage of frame  
 22 pixels;  $N_{pix}$  is the percentage of pixels represented by a  
 23 selection of the biggest bins of the histogram.  $N_{bin}$  and  
 24  $N_{pix}$  are calculated for each key-frame of the sub-shots  
 25 and are summed; if one of these values exceed a thresh-  
 26 old they are discarded.  $N_{bin}$  and  $N_{pix}$  take into account the  
 27 fact that computer graphics shots present a more “com-  
 28 pact” color histogram than realistic shots, with a low con-  
 29 trast background that allows higher quality legibility of  
 30 text and graphics. Table 4 reports the performance of the  
 31 computer graphics shots classification.

32 The algorithm takes into account the feasibility of the  
 33 presence of small motion in the CG shot, due to moving  
 34 text and graphics; An example is shown in Fig. 5.

### 3.3. Performance evaluation

35 The shot classification algorithms have been tested on  
 36 a test database of news video. To verify the robustness  
 37 of the classification process the database includes news  
 38 videos of different broadcasters, featuring different styles  
 39



Fig. 5. Example of computer graphics text.

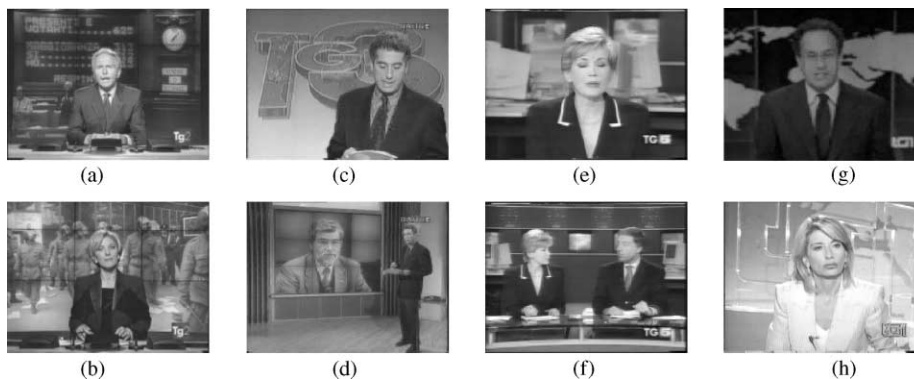


Fig. 6. Example of anchorman shots' styles.

1 and layouts, both for the anchorman shots and for the  
2 computer graphic shots.

3 A short analysis of the most recurrent styles for an-  
4 chorman shots is reported below:

- 5 • The anchorman shots are taken using a fixed position  
6 camera while there's an image in the background that  
7 shows the subject of the incoming service (Fig. 6(a)  
8 and (b)). This style is adopted for all the editions of  
9 RAI TG2.
- 10 • The anchorman can be either standing or seated, with  
11 camera movements and edit effects. The background is  
12 usually static (photos or logos). An example is shown  
13 in Fig. 6(c) and (d). For example, this style is used  
14 in the evening edition of RAI TG3.
- 15 • Two anchorman alternate each other. This is shown  
16 in Fig. 6(e) and (f). Background is almost fixed or  
17 the movement is in small regions. This is used in the  
18 evening edition of Mediaset TG5 and in some CNN  
19 editions.
- 20 • There is a more or less uniform background, some  
21 camera movements, limited number of anchor-  
22 man shots, for example front view and 3/4 view  
23 (Fig. 6(g) and (h)).

24 The results on the test set used in Section 2 are reported  
25 in Table 3. The use of the motion feature reduces the  
number of false detection errors from 14 to 3.

Table 3  
Results of the shots classification process

Anchorman shots	Detected anchorman shots	False detections	Missed detections
66	67	3 (4.5%)	2 (3%)

Table 4  
Results of the CG shots classification process

Shots	CG shots	False detections	Missed detections
318	15	10	2

Missed detections occurred with type (d) shots when  
the background contains motion. False detection occurred  
in the presence of an interview which is similar to types  
(a) and (b). To improve false detection in (b) analysis  
was restricted to the central part of the frame, according  
to the broadcaster's style.

The test set used for the computer graphic classifica-  
tion is a sub-set of the one used for the anchorman clas-  
sification (see Table 4).

Missed detections of computer graphics are due to fast  
action, like fast moving text, while the false detection  
occurred in the presence of still images, or shots that

27  
29  
31  
33  
35  
37

1 featured very little motion, with low contrast that lead to  
 2 color histogram distributions similar to that non-realistic  
 3 shots.

#### 4. Video content representation

5 To support effective video retrieval by content,  
 6 high-level information must be extracted from videos  
 7 and used to perform shot sub-classification based on their  
 8 content. Additional information to shot classification is  
 9 extracted from text captions and anchorman speech.

##### 4.1. Text recognition

11 In news videos, text captions are used to show several  
 12 information about the shot being broadcasted, such as  
 13 the site where the action takes place (in service shots)  
 14 and the names of the people shown in the video (both in  
 15 anchorman and service shots).

16 Extraction of text information from video captions has  
 17 been performed by integrating a traditional OCR within  
 18 our system. The OCR engine cannot be supplied with  
 19 raw video frames: a pre-filtering phase is required. This  
 20 phase includes two distinct steps: caption identification  
 21 and text/background separation.

22 *Caption identification:* If a shot includes a caption,  
 23 it is not guaranteed that the caption is present in the  
 24 first frames of the shot. Sometimes the caption appears  
 25 in the middle of the shot and disappears after the last  
 26 frame of the shot. Identification of frames including a  
 27 caption is based on the fact that captions are always  
 28 used in combination with graphic elements that improve  
 29 text readability. These graphic elements follow differ-  
 30 ent styles and may include opaque backgrounds and col-  
 31 ored lines (Fig. 7). Captions are always located in the  
 32 lower part of the frame. Caption identification is based  
 33 on the matching of a pre-defined model of the graphic  
 34 elements with shapes extracted in the lower part of the  
 35 frames. The model accounts for the presence of horizon-  
 36 tal stripes/long lines either colored or opaque, according  
 37 to the different broadcasters' styles.

38 *Text/background separation:* Text separation is com-  
 39 plicated by the presence of captions featuring a poor

40 text/background contrast (Fig. 7(c) and (d)). This sort of  
 41 problem is dealt with by using a text/background separa-  
 42 tion method that exploits persistence of patterns over con-  
 43 tiguous frames. This method is based on the assumption  
 44 that for the entire display of a caption all the pixels cor-  
 45 responding to the text have more or less the same value. On  
 46 the other hand, the value of the pixels in the background  
 47 changes. Text/background separation is performed by  
 48 highlighting the pixels the value of which is almost con-  
 49 stant. Captions usually display over two or more consec-  
 50 utive shots. A critical instance of text/background separa-  
 51 tion occurs when a caption without an opaque back-  
 52 ground appears over a static scene (e.g. a photo or a  
 53 painting) and is displayed only for the duration of a sin-  
 54 gle shot. This is indeed a rare condition that we did not  
 55 encounter in our test sequences.

56 Let us assume that  $\{f_0, \dots, f_k\}$  is a sequence of frames  
 57 that has been identified as including a caption. A new  
 58 sequence  $\{\hat{f}_0, \dots, \hat{f}_k\}$  is computed as follows:

$$\hat{f}_0(i, j) = 0,$$

$$\hat{f}_k(i, j) = \begin{cases} \min(255, \hat{f}_{k-1}(i, j) + \Delta) & \text{if } f_k(i, j) = f_{k-1}(i, j), \\ \max(0, \hat{f}_{k-1}(i, j) - \Delta) & \text{if } f_k(i, j) \neq f_{k-1}(i, j), \end{cases}$$

59 where  $f_k(i, j)$  is the gray level value of pixel  $(i, j)$  in  
 60 frame  $k$  and  $\Delta$  a pre-defined incremental step. In this  
 61 way, the sequence  $\{\hat{f}_0, \dots, \hat{f}_k\}$  is characterized by the  
 62 text caption that gradually fades in Fig. 8. This method  
 63 has proven to be robust even in those cases where the se-  
 64 quence of frames includes several captions that are sep-  
 65 arated by editing effects such as dissolves and cuts.

66 Finally, the sequence  $\{\hat{f}_0, \dots, \hat{f}_k\}$  is processed in or-  
 67 der to extract some frames that are used to feed the OCR  
 68 engine. For this purpose, the correlation  $C(\hat{f}_{k-1}, \hat{f}_k)$  is  
 69 computed for every pair of contiguous frames. Frames  
 70 characterized by local maxima of the correlation function  
 71 are passed to the OCR engine.

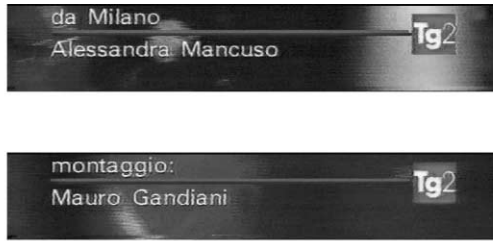
72 The graphic elements like the line and the TG2 logo  
 73 are removed since they interfere with OCR processing.

74 *The OCR engine:* To increase separation between the  
 75 single characters, and ease their segmentation, on the part  
 76 of the OCR program, *thresholding* is applied to images  
 77 extracted from the previous step.



Fig. 7. Different styles of captions.





da Milano  
Alessandra Mancuso

montaggio:  
Mauro Gandiani

Fig. 8. Character extraction.

Two OCR programs have been tested: (i) TextBridge OCR (Windows commercial OCR package), (ii) SOCR (open source OCR developed by the University of Waikato, New Zealand, <http://www.socr.org>): the most recent complete version is 0.1 and recognition rate changes according to the fonts employed. Results are shown in Table 5.

#### 4.2. Speech recognition

During the anchorman shot, the content of the following news service is summarized. While in the news service the reporter provides detailed information on the topic.

In order to improve the speech recognition rate and extract only relevant information about news service content, the speech recognition engine is fed with the audio track of the anchorman shots only. In fact, as reported in Ref. [12], generally there is not an exact synchronization between speech and objects shown in the news service. Often the content of the shots does not correspond to the reporter's description, consequently the association of the audio track content with the corresponding shots may lead to erroneous results. Furthermore, the audio track of news services is typically disturbed by background noise and sometimes includes speech transmitted through low-quality telephone or satellite links.

*Speech recognition engine:* The speech recognition engine used is IBM ViaVoice 98 that features speaker independence and continuous speech processing.

The speech recognition engine is based on a hidden Markov model of the language and uses the following resources: (i) a language thesaurus that can be customized

Table 5  
OCR and speech recognition results

Speech recognition			OCR		
Anchor. shots			All audio track	TextBridge	SOCR
Not trained	Trained	Trained			
57%	84%	52%	87%		~ 60%

and enhanced, (ii) a customizable model of word usage, (iii) word pronunciation models. A database of audio tracks was used to train the words usage model. The database included the audio tracks of anchorman of several broadcasters. Sentences corresponding to speech were manually transcribed. Their content covered different topics such as sports news, politics, chronicles and gossip.

The speech recognition rate was measured on a test database that did not include any of the audio tracks used for training. Results are shown in Table 5.

Words extracted by the speech recognition engine were filtered in order to wipe out all utility words (articles, pronouns, conjunctions and prepositions—this accounts approximately for 50% of the speech in latin languages). Remaining words are used to describe the content of the following news services.

#### 5. Video retrieval

Techniques for video segmentation, shot classification and shot content description presented in the previous sections have been integrated into a system for content-based retrieval of TV news. At archiving time, news videos are automatically processed in order to extract content descriptors for each video shot. The content descriptor of a generic shot includes:

*Shot type identifier:* This can be either *anchorman*, *news service* or *computer graphics*.

*TV broadcaster identifier.*

*Broadcast date and time.*

*Visual shot descriptor.* This is the key-frame of the shot.

*Textual shot descriptor.* This is the set of words extracted from shot captions and from speech recognition of the previous anchorman shot. Manual annotation can be added.

At retrieval time, the system supports video querying and browsing. To reduce the effect of the errors of the OCR programs, the retrieval system uses the AGREP approximate text search that allows to find words that

(a)



(b)



(c)

Fig. 9. (a) Specification of a fuzzy search for “President Clinton”. (b) shows result page with shots that match both the keywords, (c) shows the last page with shots that match only the keyword “President”.

1 contain errors. Queries formulated according to TV  
 2 broadcaster, date, time, content and any Boolean combination of these are supported. One or more words  
 3 can be input by the user. These are matched against textual shot descriptors of database videos through the use of  
 4 a thesaurus so as to support exact word and synonym matching. Matched shots are presented to the user for  
 5 browsing. For each matched shot all the information stored in its content descriptor is shown.

6 In Fig. 9(a) a sample query by content is shown.  
 7 The user enters a Boolean combination of the words

8 ‘President’ or ‘Clinton’ to search for shots with similar  
 9 content. Retrieval results are shown in Fig. 9(b). The  
 10 query also retrieves shots classified with the Italian word  
 11 “Presidente” (speech transcription and manual annotation), since the “fuzzy” search method is used. Retrieved  
 12 shots are shown in decreasing order of match. The first  
 13 shots match both query keywords and show news and anchorman shots related to “President Clinton”. The other  
 14 shots retrieved match only the keyword “President” and show “President Milosevic” and “President Scalfaro”.  
 15 The “Previous” and “Next” buttons on the top of the

13  
 15  
 17  
 19  
 21

1 window allow the user to navigate through all the re-  
 2 trieved shots.

3 Selection of a shot key-frame from the output interface  
 4 allows display of the entire shot through a movie player  
 5 application.

## 6. Conclusions

7 This paper presents a system for content-based in-  
 8 dexing and retrieval of news videos. The system fea-  
 9 tures content-based shot classification of anchorman and  
 10 *non-realistic* shots, to allow the reuse of report shots. Ex-  
 11 traction of high-level content descriptors through caption  
 12 OCR and speech recognition. Shot classification is based  
 13 on statistical and motion features of the news video struc-  
 14 ture, so as to provide independence from TV broadcaster  
 15 style.

## References

- 17 [1] A. Del Bimbo, C. Colombo, P. Pala, Semantics in visual  
 18 information retrieval, *IEEE Multimedia* 6 (3) (1999).  
 19 [2] D. Swanberg, C. Shu, R. Jain, Knowledge guided parsing  
 20 in video databases, *Spie* 1908 (13) (1993) 13–24.  
 21 [3] S. Smoliar, H.J. Zhang, Y. Gong, Automatic parsing  
 22 of news video, *Proceedings of the IEEE Conference*  
 23 *on Multimedia Computing and Systems*, May 1994, pp.  
 24 45–54.  
 25 [4] T. Kanade, Y. Nakamura, Semantic analysis for video  
 26 contents extraction spotting by association in news video,  
 27 *ACM Multimedia* 97 (1997).  
 28 [5] H. Zhang, B. Furht, S.W. Smoliar, *Video and Image*  
 29 *Processing in Multimedia Systems*, Kluwer Academic  
 30 Publishers, Dordrecht, 1995.  
 31 [6] A.G. Hauptmann, M.J. Witbrock, *Informedia: news-on-*  
 32 *demand multimedia information acquisition and re-*  
 33 *trieval*, *Intell. Multimedia Informat. Retrieval* (1997)  
 213–239. 35  
 [7] T. Sato, T. Kanade, E.K. Hughes, M.A. Smith, Video ocr  
 for digital news archive, *IEEE International Workshop*  
 on Content-Based Access of Image and Video Databases  
 CAIVD' 98, 1998, pp. 52–60. 37 39  
 [8] S. Eickeler, S. Muller, Content-based video indexing of tv  
 broadcast news using hidden Markov models, *Proceedings*  
 of the IEEE International Conference on Acoustics,  
 Speech, and Signal Processing (ICASSP), March 1999,  
 pp. 2997–3000. 41 43  
 [9] R. Lienhart, Indexing and retrieval of digital video  
 sequences based on automatic text recognition, *Fourth*  
*ACM International Multimedia Conference*. 45 47  
 1996.  
 [10] A.G. Hauptmann, Speech recognition in the informedia  
 digital video library: uses and limitations, *ICTAI 95*, 1995. 49  
 [11] A.G. Hauptmann, H.D. Wactlar, M.J. Witbrock,  
*Informedia: News-on-demand experiments in speech*  
 recognition, *ARPA Speech Recognition Workshop*, 1996. 51 53  
 [12] M.J. Witbrock, A.G. Hauptmann, Speech recognition for  
 a digital video library, *JASIS*, 1996. 55  
 [13] D. Luparello, B. Merialdo, K.T. Lee, J. Roudaire,  
 Automatic construction of personalized tv news programs,  
*Proceedings of ACM Multimedia 99*, 1999, pp. 323–330. 57  
 [14] B. Merialdo, Automatic indexing of tv news, *WIAMIS*  
 '97, June 1997. 59  
 [15] J.S. Boreczky, L.A. Rowe, Comparison of video  
 shot boundary detection techniques, *Technical*  
 Report, Computer Science Division-EECS, University of  
 California Berkeley. 61 63  
 [16] M. Caliani, C. Colombo, A. Del Bimbo, Commercial  
 video retrieval by induced semantics, *IEEE International*  
 Workshop on Content-Based Access of Image and Video  
 Databases CAIVD '98, 1998, pp. 72–80. 65 67  
 [17] U. Gargi, R. Kasturi, An evaluation of color histogram  
 based methods in video indexing, *First International*  
 Workshop on Image Database and Multi-media Search,  
 1996, pp. 75–82. 69 71

**About the Author**—MARCO BERTINI has a research grant and carries out his research activity at the Department of Systems and Informatics at the University of Florence, Italy. He received a MS in electronic engineering from the University of Florence in 1999. His main research interest is content-based indexing and retrieval of videos. 73

**About the Author**—ALBERTO DEL Bimbo graduated in 1978 from the University of Florence, Italy, where he is presently Full Professor of Computer Engineering. Presently he is the Director of the Master in Multimedia at the University of Florence and the Deputy Rector for Research and Innovation Transfer of the University of Florence.

His scientific interests and activities have addressed the subjects of Image Technology and Multimedia, with particular reference to object recognition and image sequence analysis, content-based retrieval for image and video databases and advanced man-machine interaction. Prof. Del Bimbo is the author of over 150 publications, that have appeared in the most distinguished international journals and conference proceedings and is the author of the monography "Visual Information Retrieval" edited by Morgan Kaufman in 1999. From 1996 to 2000, he has been the President (formerly Vice-President) of the Italian Chapter of IAPR, the International Association for Pattern Recognition. He obtained the IAPR fellowship in 2000. He presently serves as Associate Editor of *IEEE Transactions on Multimedia*, *IEEE Transaction on Pattern Analysis and Machine Intelligence*, *Pattern Recognition*, *Pattern Analysis and Applications Journal*, *Journal of Visual Languages and Computing* and *Multimedia Tools and Applications Journal*. Since 1999 he has been Member at Large of the IEEE Publications Board.

**About the Author**—PIETRO PALA is an assistant professor in the Department of Systems and Informatics at the University of Florence, Italy. He received his MS in electronic engineering at the University of Florence in 1994. He received a Ph.D. in information science from the same university in 1998. His current research interests include pattern recognition, image and video retrieval by content, and related applications. 75