

A GENETIC ALGORITHM-BASED APPROACH TO KNOWLEDGE-ASSISTED VIDEO ANALYSIS

N. Voisine², S. Dasiopoulou^{1,2}, F. Precioso², V. Mezaris^{1,2}, I. Kompatsiaris², M. G. Strintzis^{1,2}

¹Information Processing Laboratory, Electrical and Computer Engineering Department, Aristotle University of Thessaloniki, Greece

²Informatics and Telematics Institute/Centre for Research and Technology Hellas, Thessaloniki, Greece

ABSTRACT

Efficient video content management and exploitation requires extraction of the underlying semantics, a non-trivial task associating low-level features of the image domain and high-level semantic descriptions. In this paper, a knowledge-assisted approach for extracting semantics of domain-specific video content is presented. Domain knowledge considers both low-level features (color, motion, shape) and spatial behavior (topological and directional information). First a segmentation algorithm generates a set of over-segmented, homogenous atom-regions. A genetic algorithm is applied, after the preprocessing step, in order to optimize the scene interpretation according to the knowledge of the specific domain. The proposed approach was tested on the Tennis and Formula One domains and shows promising results.

1. INTRODUCTION

The recent advances in computing technologies have made available vast amounts of digital video content, leading to ever-increasing flow of audiovisual information. This results in a growing demand for efficient segmentation/analysis methods for extracting semantic information from such content, since the acquisition of higher-level information in terms of meaning is the key enabling factor for the management and exploitation of video content. However, due to the possible different interpretations and intended uses, the ambiguity that is inherent in visual information renders the development of faster hardware or the evolution of classic segmentation algorithms insufficient. This difficulty [1] in mapping semantic concepts as perceived by humans into a set of automatically extracted low-level image features, can be alleviated for a particular application domain by means of domain specific knowledge. Different approaches have been used for the implementation of particular parts of the domain knowledge such as formal knowledge representation theories, semantic web technologies, Dynamic Belief networks etc. For example, in [2], semantic web technologies are used, while in [3] *a priori* knowledge representation models are used as a knowledge base that assists semantic-based classification and clustering. In [4], an object ontology coupled with a relevance feedback mechanism is introduced, in [5], semantic entities, in the context of the MPEG-7 standard, are used for knowledge-assisted video analysis and object detection, while in [6], associating low-level representations and high-level semantics is formulated as a probabilistic pattern recognition problem.

This work was supported by the EU projects SCHEMA “Network of Excellence in Content-Based Semantic Scene Analysis and Information Retrieval” (IST-2001-32795) and aceMedia “Integrating knowledge, semantics and content for user centered intelligent media services” (FP6-001765).

In this paper a knowledge-assisted, domain-specific video analysis framework is presented, using a genetic algorithm to support efficient object localization and identification. The localization and identification of the domain salient objects are prerequisites for extracting accurate semantic information. An initial segmentation generates automatically a set of atom-regions and subsequently their low-level descriptors are extracted. Analysis may then be performed by using the necessary processing tools and by relating high-level symbolic representations included in the ontology to visual features extracted from the signal domain. Additionally, the genetic algorithm decides how the atom-regions should be merged in order to form objects in compliance with the object models defined in the domain ontology. Following this approach, the detection of the important objects depends largely on the knowledge base of the system and as a result it can be easily applied to different domains provided that the knowledge base is enriched with the respective domain knowledge.

The remainder of the paper is structured as follows: section 2 considers domain ontology development, section ?? contains a presentation of the employed segmentation and descriptor extraction algorithms, while in section 4 the implementation of the genetic algorithm is discussed. Experimental results are presented in section 5 and finally, conclusions are drawn in section 6.

2. DOMAIN KNOWLEDGE

The knowledge about the examined domain is encoded in the form of an ontology. The developed ontology includes the objects that need to be detected, their visual features and their spatiotemporal interrelations. These descriptions provide the system with the required knowledge to find the optimal interpretation for each of the examined video scenes, i.e. the optimal set of mappings among the available atom-regions and the corresponding domain-specific semantic definitions. To account for objects of no interest that may be present in a particular domain and for atom-regions that fail to comply with any of the object models included in the ontology, the concept of unknown object is introduced. In addition, support is provided for the definition of associations between low-level descriptions and the algorithms to be employed for their extraction. In the following, a brief description of the main classes is presented.

Class **Object** is the superclass of all objects to be detected during the analysis process. When the ontology is enriched with the domain specific information it is subclassed to the corresponding domain salient objects. Class **Object Interrelation Description** describes the objects spatiotemporal behavior, while **Low-Level Description** refers to the set of their representative low-level visual features. Since real-world objects tend to have multiple different

instantiations, it follows that each object prototype instance can be associated with more than one spatial description and respectively multiple low-level representations. Different classes have been defined to account for the different types of low-level features (color, shape, motion etc.). These are further subclassed to reflect the different ways to represent such a feature (e.g. color information could be represented by any of the color descriptors standardized by MPEG-7, the distribution models of the respective color space etc.) The actual values that comprise the low-level descriptors (e.g. the DC value elements, color space) are under the **Low-Level Descriptor Parameter** class.

Providing domain-specific spatiotemporal information proves to be particularly useful for the identification of specific objects, since it allows discrimination of objects with similar low-level characteristics as well as of objects whose low-level features alone are not adequate for their identification. The implemented spatial relations consider two-dimensional, binary relations, defined between regions with connected boundaries. In the current implementation the included spatial relations are the four relative directional relations, i.e. right, left, above, below, and the eight topological relations resulting from the 9-intersection model as described in earlier works on spatial relations representation and reasoning [7, 8]. More qualitative spatial relations such as near, far, adjacent, between etc. can be defined combining the ones previously mentioned. Symmetricity and transitivity properties as well as the inverse of each of the defined spatial relations are specified. Consequently, more complex spatial descriptions can be inferred, reducing at the same time the number of required explicit descriptions. The used low-level descriptors are the MPEG-7 Dominant Color and Region Shape descriptors, the motion norm of the averaged global motion-compensated block motion vectors for each region blocks and the ratio between a region's area and the square of its perimeter (compactness).

Enriching the ontology with domain specific knowledge results in populating the ontology with appropriate instances, i.e. prototypes for the objects to be detected. The presented system interprets the provided information (i.e. the low-level and spatial relation descriptions) as a conjunctive normal form clause, with one conjunct for each description category. Each conjunct is the disjunction of the respective category descriptors associated with the particular prototype instance. As will be explained in section 4, fuzzy matching criteria are incorporated in the fitness function used to determine the plausibility of each interpretation. The followed approach proves twofold advantageous as it also allows to tackle the inevitable loss of objects connectivity in the 2D image plane: over-segmented atom-regions belonging to the same object class are appropriately merged to form a single instance of the respective concept.

3. PREPROCESSING

Under the proposed framework, a set of over-segmented atom-regions is generated by combining the resulting color and motion segmentation masks of the preprocessing step.

3.1. Color and motion initial segmentation

The color segmentation is based on the extraction of up to eight dominant colors in the frame, as proposed in the MPEG-7 Dominant Color descriptor [9]. We then initialize a simple K-means

algorithm with the extracted colors, as detailed in [10]. The motion segmentation is based on a two step algorithm.

We first apply the segmentation methodology of [4], considering a block matching approach [11], in order to obtain a coarse but very fast segmentation. Indeed, an iterative rejection scheme [12] based on the bilinear motion model is used to effect foreground/background segmentation. Following that, meaningful foreground spatiotemporal objects are formed by initially examining the temporal consistency of the output of iterative rejection, clustering the resulting foreground macroblocks to connected regions and finally performing region tracking. Furthermore, this first step provides an efficient estimation of the 8 parameters of the bilinear camera motion model.

As a second step, we consider the previous motion segmentation as initialisation of a region-based motion segmentation algorithm using smoothing spline active contours [13]. Smoothing splines offer a robust active contour implementation to overcome the problem of considering noisy data that working with mpeg streams implies. Hence, we improve the accuracy of the first step motion segmentation. Furthermore, the contour defining extracted moving regions is given by a parametric equation which allows a fast computation for geometric curve features as perimeter, area, or moments, involved in the low-level feature descriptor extraction.

Finally, we merge the results of color and motion segmentation by giving priority to color information. That is to say, if a motion-based segmentation region consists of two or more color-based segmented atom-regions, we split this region according to the color segmentation. Then we apply a region-based smoothing spline active contour segmenting homogeneous regions of the color segmentation mask in order to provide the parametric contour equation of each color regions.

3.2. low-level feature extraction descriptors

The low-level descriptors defined in section 2 are extracted for each atom-region as follows. We compute the Dominant Color descriptor applying the MPEG-7 eXperimentation Model (XM) [9]. The region motion feature, based on the aforementioned motion segmentation algorithm, is defined by the norm of the average global-motion-compensated motion vectors evaluated on the blocks belonging to the atom-region considered. To extract the compactness descriptor, we compute the area and the perimeter of each region using a fast algorithm, proposed in [14], based on spline properties of the parametric contour description.

An other good spline contour property allows us to cope with one of the main issue of the region shape descriptor extraction. Indeed, the Angular Radial Transform involved in the MPEG-7 region shape computation considers the evaluation for each regions of specific normalized central moments. These moments are defined considering that each region is included into the unit disk. This normalization process can be efficiently managed using the spline structure affine invariance, since a spline curve subjected to an affine transformation is still a spline curve whose parameters (the control points) are obtained by subjecting the original spline curve control points to that affine transformation. Thus we compute the first geometric moments of the considered region (its area, its centroid,...) using fast algorithms [14] in order to evaluate the parameters of the affine transformation corresponding to the normalized contour and then extract the MPEG-7 region shape descriptor.

4. GENETIC ALGORITHM

As previously mentioned, the initially applied color and motion segmentation algorithms, result in a set of over-segmented atom-regions. Assuming for a single image N_R atom regions and a domain ontology of N_O objects, there are $N_R^{N_O}$ possible scene interpretations. To overcome the computational time constraints of testing all possible configurations, a genetic algorithm is used [15]. Genetic algorithms (GAs) have been widely applied in many fields involving optimization problems, as they proved to outperform other traditional methods. They build on the principles of evolution via natural selection: an initial population of individuals (chromosomes encoding the possible solutions) is created and by iterative application of the genetic operators (selection, crossover, mutation) the optimal, according to the defined fitness function, solution is reached.

In our framework, each individual represents a possible interpretation of the examined scene, i.e. the labelling of all atom-regions either as one of the considered domain objects or as unknown. An object instantiation is identified by its corresponding concept and an identifier used to differentiate instances of the same concept. The domain ontology contains information about the maximum allowed number of detected instances for each object. In order to reduce the search space, the initial population is generated by allowing each gene to associate the corresponding atom-region only with those objects that the particular atom-region is most likely to represent. For example in the domain of Tennis a green atom-region may be interpreted as a Field, Wall or Unknown object but not as Ball or Player. Therefore, for each individual included in the initial population, the corresponding gene is associated with one of the three aforementioned object concepts (instead of the available N_O). The set of plausible candidates for each atom-region is estimated according to the low-level descriptions included in the domain ontology.

The following functions are defined to estimate the degree of matching between R_i and o_j , in terms of low-level visual and spatial features respectively:

- the interpretation function $\mathcal{I}_M(g_i) \equiv \mathcal{I}_M(R_i, om_j)$, assuming that g_i associates region R_i with object o_j having model om_j , to provide an estimation of the degree of matching between an object model om_j and a region R_i . $\mathcal{I}_M(R_i, om_j)$ is calculated using the descriptor distance functions realized in the MPEG-7 XM and is subsequently normalized so that $\mathcal{I}_M(R_i, om_j)$ belongs to $[0, 1]$.
- the interpretation function \mathcal{I}_R , which provides an estimation of the degree to which a relation \mathcal{R} holds between two atom-regions.

Since each individual represents the scene interpretation, the Fitness function has to consider the above defined low-level visual and spatial matching estimations for all atom-regions. As a consequence the employed Fitness function is defined as follows:

$$Fitness(G) = \sum_{g_i} \mathcal{I}_M(g_i) + \sum_k \sum_{(g_i, g_j), g_i \mathcal{R}_k g_j} \mathcal{I}_{\mathcal{R}_k}(g_i, g_j)$$

where $\mathcal{I}_M(g_i)$ is the estimation function of gene g_i regarding low-level visual similarity and $\mathcal{I}_{\mathcal{R}_k}(g_i, g_j)$ is the estimation function of spatial similarity between g_i and g_j in terms of \mathcal{R}_k . It follows from the above definitions that the optimal solution is the one that maximizes the Fitness function. This process elegantly handles the

merging of atom-regions: any neighboring such regions belonging to the same object according to the generated optimal solution are simply merged. In our implementation, the following genetic operators were used: roulette wheel selection, in which individuals are given a probability of being selected that is directly proportionate to their fitness and uniform crossover, where genes of the parent chromosomes are randomly copied.

5. EXPERIMENTAL RESULTS

The proposed approach was tested on a variety of Formula One and Tennis domain MPEG-2 videos. As illustrated in 1 and 2, the system output is a segmentation mask outlining the semantic interpretation, i.e. a mask where different colors representing the objects defined in the ontology are assigned to each of the produced regions. The objects of interest included in each domain ontology along with their low-level models and spatial relations are illustrated in table 1. In both domains, the low-level descriptors values included in the corresponding knowledge base were extracted from a training set of manually annotated images.

The time required for performing the previously described tests was between 5 and 10 seconds per frame, excluding the process of motion information extraction via block matching for which efficient and inexpensive hardware implementations exist [11]. More specifically, the time to perform pixel-level segmentation was about 2 seconds, while the time required by the genetic algorithm to reach an optimal solution varied depending on the number of atom-regions and the number of spatial relations. The extraction of the low-level and spatial descriptions is performed before the application of the genetic algorithm. In general, the proposed approach proved to produce satisfactory results as long as the initial color-based segmentation did not segment two objects as one atom-region.

Concept	Visual models	Spatial relations
Road	$DC_{road}^1 \vee DC_{road}^2 \vee DC_{road}^3$	Road ADJ Grass,Sand
Car	$MOV_{car}^1 \wedge CPS_{car}^1$	Car INC Road
Sand	$DC_{sand}^1 \vee DC_{sand}^2$	Sand ADJ Grass, Road
Grass	$DC_{grass}^1 \vee DC_{grass}^2 \vee DC_{grass}^3$	Grass ADJ Road,Sand
Field	$DC_{field}^1 \vee DC_{field}^2 \vee DC_{field}^3$	Field ADJ Wall
Player	MOV_{Player}^1	Player INC Field
Line	$DC_{line}^1 \wedge CPS_{line}^1$	Line INC Field
Ball	$DC_{Ball}^1 \wedge CPS_{Ball}^1$	Ball INC Field
Wall	$DC_{Wall}^1 \vee DC_{Wall}^2 \vee DC_{Wall}^3$	Wall ADJ Field

Table 1. Formula One and Tennis domain definitions (*dominant color descriptor (DC), motion descriptor (MOV), compactness descriptor (CPS), adjacency relation (ADJ), and inclusion relation (INC)*)

6. CONCLUSIONS

In this paper, a knowledge-assisted domain-specific video analysis approach, which exploits the fuzzy inference capabilities of a genetic algorithm, is presented. Domain knowledge includes both low-level visual descriptors and spatial interrelations, and is encoded in the form of an ontology. The genetic algorithm provides a fundamentally different framework compared to knowledge-based systems using formal rules. By encoding the object models defined in the ontology in the form of constraints (fitness function definition), a global optimal interpretation of the examined scene

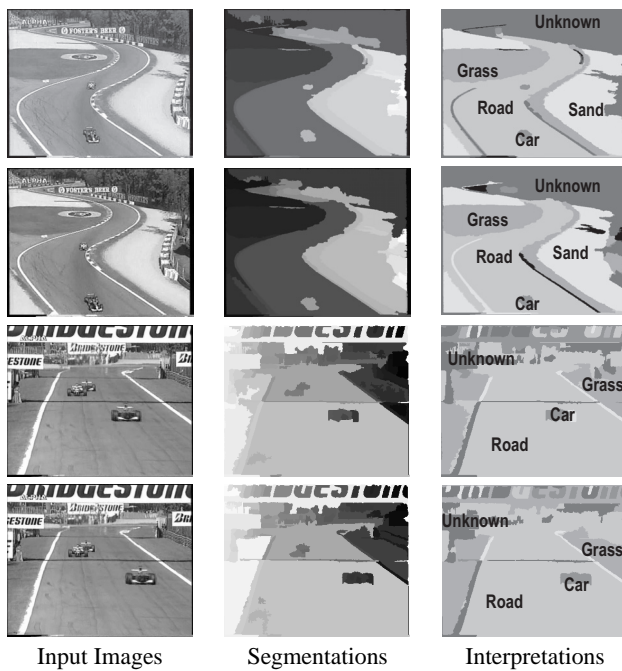


Fig. 1. Formula One domain results

is reached. The developed domain ontology provides a flexible conceptualization that allows the easy addition of new low-level and spatiotemporal descriptors, i.e. supports different abstraction levels.

7. REFERENCES

- [1] W. Al-Khatib, Y.F. Day, A. Ghafoor, and P.B. Berra. Semantic modeling and knowledge representation in multimedia databases. *IEEE Transactions on Knowledge and Data Engineering*, 11(1):64–80, Jan/Feb 1999.
- [2] J. Hunter, J. Drennan, and S. Little. Realizing the hydrogen economy through semantic web technologies. *IEEE Intelligent Systems Journal - Special Issue on eScience*, 19:40–47, 2004.
- [3] A. Yoshitaka, S. Kishida, M. Hirakawa, and T. Ichikawa. Knowledge-assisted content-based retrieval for multimedia databases. *IEEE Multimedia*, 1(4):12–21, Winter 1994.
- [4] V. Mezaris, I. Kompatsiaris, N.V. Boulgouris, and M.G. Strintzis. Real-time compressed-domain spatiotemporal segmentation and ontologies for video indexing and retrieval. *IEEE Trans. on Circuits and Systems for Video Technology*, 14(5):606–621, May 2004.
- [5] G. Tsechpenakis, G. Akrivas, G. Andreou, G. Stamou, and S.D. Kollias. Knowledge-Assisted Video Analysis and Object Detection. In *Proc. Eumite02*, Algarve, Portugal, September 2002.
- [6] M. Ramesh Naphade, I.V. Kozintsev, and T.S. Huang. A factor graph framework for semantic video indexing. *IEEE Trans. on Circuits and Systems for Video Technology*, 12(1):40–52, Jan. 2002.
- [7] David A. Randell, Zhan Cui, and Anthony G. Cohn. A spatial logic based on regions and connection. In *KR*, pages 165–176, 1992.
- [8] Max J. Egenhofer and Robert D. Franzosa. Point set topological relations. *International Journal of Geographical Information Systems*, 5:161–174, 1991.
- [9] B.S. Manjunath, J.-R. Ohm, V.V. Vasudevan, and A. Yamada. Color and texture descriptors. *IEEE Trans. on Circuits and Systems for Video Technology, special issue on MPEG-7*, 11(6):703–715, June 2001.
- [10] V. Mezaris, I. Kompatsiaris, and M.G. Strintzis. A framework for the efficient segmentation of large-format color images. In *Proc. International Conference on Image Processing*, volume 1, pages 761–764, 2002.
- [11] J.-C. Tuan, T.-S. Chang, and C.-W. Jen. On the data reuse and memory bandwidth analysis for full-search block-matching VLSI architecture. *IEEE Trans. on Circuits and Systems for Video Technology*, 12(1):61–72, January 2002.
- [12] T. Yu and Y. Zhang. Retrieval of video clips using global motion information. *Electronics Letters*, 37(14):893–895, July 2001.
- [13] F. Precioso, M. Barlaud, T. Blu, and M. Unser. Smoothing B-spline active contour for fast and robust image and video segmentation. In *International Conference on Image Processing*, Barcelona, Spain, September 2003.
- [14] M. Jacob, T. Blu, and M. Unser. An exact method for computing the area moments of wavelet and spline curves. *IEEE Trans. Pattern Anal. and Mach. Intell.*, 23(6):633–642, June 2001.
- [15] M. Mitchell. *An introduction to Genetic Algorithms*. MIT Press., 1996.

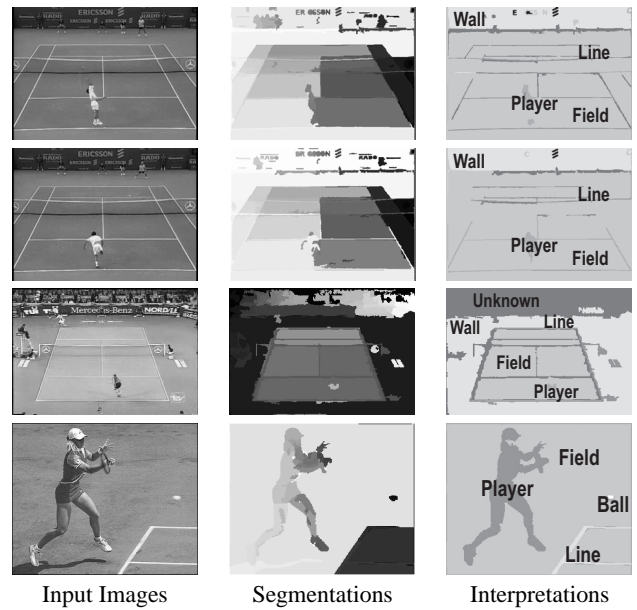


Fig. 2. Tennis domain results