



ISSN: 2350-0328

**International Journal of Advanced Research in Science,
Engineering and Technology**

Vol. 5, Issue 1 , January 2018

Smart sentence framing in video

ADITYA KALE, OM KOLTE, ANKITA KATE, APURVA KAMAJI, SHWETA KAKADE

JSPM RSCOE, S.P.PUNE UNIVERSITY,PUNE,INDIA
JSPM RSCOE, S.P.PUNE UNIVERSITY,PUNE,INDIA
JSPM RSCOE, S.P.PUNE UNIVERSITY,PUNE,INDIA
JSPM RSCOE, S.P.PUNE UNIVERSITY,PUNE,INDIA
JSPM RSCOE, S.P.PUNE UNIVERSITY,PUNE,INDIA

ABSTRACT: Nowadays there is the popularity of online video sharing is high, this is the source of information and entertainment. Therefore video annotation came in picture in past few years. In this paper four-step approach to automatically annotate video shots with sentence. In the first step is to converting the video shots into the continuous frame image. The Second step is to find similar candidate elements of the sentence about frame. The main elements are objects, events, scenes, modifiers. This candidate element are gained by similar images with the video frame in collected image data set. The third step is to select the best candidate element by a weighted scoring algorithm. In fourth step uses the correlation graph algorithm to find out the relation among the best element. This methods are effective to annotate the video with sentence. The weighted scoring algorithm and correlation graph algorithm are more effective in this experiment.

KEYWORDS: Algorithms, video Annotation, image data set, sentence element, YouTube

I. INTRODUCTION

In 21st century with help of social media image have become more and more accessible in the public. But image is not more informative as per the people's Demand. The other informative media needed the video become a more informative then the image, and different technologies help to understand the video content. Conventional approach to video annotation dominantly focuses on administered identification of limited set of concepts. The no of ambiguous meaning will comes in picture when only keyword is provided for searching the video.

In these scenes, generating sentence for better understanding and highlighting the activities happening in the video is annotated with a sentence. It's easier to user understand and flexible search. While the idea of video annotating with the sentence is likely to turn out well, there are so many challenges. It is not easy to collect video as training data set because most of the online video has no label. It is not possible to manually annotate a huge amount of video with sentence. The challenging task is representing the contents of videos with natural language is more complex and multifaceted compare to manual tags; need to take all objects, actions and scenes of events. More important is the correct grammar is another factor should consider of. In this video main focus on unconstrained video data download from online video portals such as YouTube. The content of these videos are very distinct in theme and practical in content, which makes our sentence generation more challenging. experiment are conducted on the unprocessed video shots and NUS-WIDE image data set.

- In this paper main proposal is automatic sentence generation approach for the free style homemade videos downloaded from YouTube. Main task is to collect the series of descriptive vocabularies for the video shot and generated sentence.
- In this approach user annotated data are used for avoiding the costly addition of manual annotated data. While working on project use image with the user generated tags instead of avoiding the problem. This problem occurs because of unlabeled training videos.
- In third phase proposed weighted scoring algorithm and correlation graph algorithm, Weighted scoring algorithm used to verify the Correctness of framed sentence and correlation Graph algorithm to making good sense of that sentence.



International Journal of Advanced Research in Science, Engineering and Technology

Vol. 5, Issue 1 , January 2018

- It makes interesting way to generate the sentence, and it is time consuming.

content-based image retrieval (CBIR), also known as query by image content (QBIC) and content-based visual information retrieval (CBVIR) is the application of computer vision techniques to the image retrieval problem. The term ‘content’ in this context trying to match colors, shapes, textures, tags or any other information that can be derived from the image itself. The underlying search algorithms are so much depending on the application, but result of the images is depending on elements with the provided example. The major features of CBIR.

- Found the match to particular combination of color, texture or shape features (e.g. green stars);
- arrangement of specific types of object (e.g. chairs around a table);
- the drawing of a specific type of event (e.g. a football match);
- the availability of named individuals, locations, or events (e.g. the Queen greeting a crowd);
- corresponding to emotions one might associate with the image (e.g. happiness);

Automatic image annotation also known as automatic image tagging or linguistic indexing is the process by which a computer system uses keywords to query image. This application of annotation techniques is used image retrieval systems to organize and point to the images of interest from a database.

II. CONTENT BASED IMAGE RETRIEVAL (CBIR)

In Content-based image retrieval (CBIR), user has to provide query image as input to search engine instead of providing text query as in case of text based image retrieval. The user search for query image from the storage in hard disk or same kind of image provide by user. Content based image retrieval (CBIR) retrieves images based on features like color, texture and shape. Content-based means that image search will search for the main contents of the image. The ideal or main content of this context might refer color, shape, texture, or any other feature/information that can be extracted from the image itself. Features of query image and images in database are extracted and are stored in feature vector. In CBIR feature vector of all images stored in database trying to compare with feature vector of query image and the image whose feature which matches maximum will be indexed first. Therefore, images will be indexed according to visual content for features like shape, color and texture or any other feature or a combination of set of query image features.

Advantages:

In CBIR there is the automatic retrieval of images based on their visual content, whereas the keyword-based approach requires time-consuming annotation of images in database. CBIR retrieves relevant images fast and does not need of manual annotation of images.

Limitation: the image with the same feature may different users have different tags.

Image Retrieval system	Concept	Advantages	Limitations
Content based Image Retrieval Systems	For this System user has to provide query images as input and system retrieve images similar to query images. Similarity is calculated based on visual content of images	Retrieves relevant images quickly and doesn't need of manual annotation of images	High feature similarity may not always match to semantic similarity

Annotation technique	Advantage	Disadvantage
Machine learning	Improves Sentiment Analysis, Improves Natural Language Processing	Requires a time & clustering, Bayesian network, Costly
Rule learning	Used modularity, Uniformity, Naturalness in rule learning.	possibility of contradictions and Infinite chaining
Based on graph	Semi-supervised method, improve concept annotation results.	Take more time for no of iteration.
Ontology	Explicit specification of a conceptualization, Hierarchical categories, classification of keywords.	Low level feature extraction. Complexity is directly proportional to cost.

III. Machine Learning

Low-level features can be extracted from the video or image. Various machine learning techniques such as support vector machine (SVM), Bayesian networks, Clustering, similarity and metric learning can be used. A framework for semantic video event annotation is presented, which exploits global feature, local feature and motion feature. Using these features, video clip can be encoded as a set of feature vectors. Then according to different features, SVM classifiers are trained, and a bicoded chromosome based genetic algorithm is performed to obtain optimal classifiers and relevant optimal weights based on training stage. With the optimal classifiers set and optimal weights, the maximum similarity between video clip in original database and unlabeled video clip is considered to be the final label result.

annotation is a supervised learning problem under Multiple-Instance Learning (MIL) framework. A novel Asymmetrical Support Vector Machine-based MIL algorithm is proposed, which extends the conventional Support Vector Machine. By increasing the pattern margins corresponding to the MIL constraints, ASVM-MIL converts the MIL problem to a traditional supervised learning problem. classifying weakly-annotated images, where just a small subset of the database is annotated with keywords. In this paper a new method by integrating semantic concepts extracted from text and by automatically extending annotations to the images with missing keywords is proposed. The model is inspired from the probabilistic graphical model theory. Bayesian networks are used.

IV Rule Learning

Visual features can be directly extracted from the video or images. These low-level features can be used for annotation but gap exists between the information that can be extracted automatically from visual data and the



ISSN: 2350-0328

International Journal of Advanced Research in Science, Engineering and Technology

Vol. 5, Issue 1 , January 2018

interpretation that the same data has for a user in a given situation. Rules are built to infer a set of high-level concepts from low-level descriptors. However; the used knowledge representations are predefined and static, limiting the adaptability to different contexts. A rule based video annotation system is proposed. The proposed system annotates video sequences automatically using knowledge from a pre-annotated dataset.

It creates representations from a set of low-level video features and infers the association rules between them and high-level concepts from a predefined lexicon. learning by means of Fuzzy Decision Trees (FDT), automatic rules based on a limited set of examples is proposed. Rules intended, in an exploitation step, to reduce the need of human usage in the process of indexation. Occurrence of some audiovisual features demonstrates remarkable patterns for detection of semantic events. present an approach for event detection and annotation of broadcast soccer video. A fuzzy rule-based reasoning system is designed as a classifier which adopts statistical information from a set of audiovisual features as its crisp input values and produces semantic concepts corresponding to the occurred events. A set of tuples is created by discretization and fuzzification of continuous feature vectors derived from the training data. We extract the hidden knowledge among the tuples and correlation between the features and related events by constructing a decision tree (DT).

V .Based On Graph

Based learning is a semi-supervised method. Graph with labeled and unlabeled vertices are used. These vertices are samples; the edges reflect the similarities between sample pairs. A function is estimated on the graph based on a label smoothness assumption.

Video annotation used for assigning the single or multiple constraint labels to a target data set, where the assignment is often done separately without considering the inter-concept relationship. Due to the fact that concepts do not occur in isolation (e.g., smoke and explosion); context-based video annotation with graph diffusion process. . Filtered tags are used; they are superior to a state-of-the-art semi-supervised technique for graph reinforcement learning on the initial user-supplied annotations. A multi-graph based semi-supervised learning method is proposed. This framework amounts to fusing graphs and then conducting semi-supervised learning on fused graph.

VI. Ontology

Ontology is defined as an explicit specification of a conceptualization. It is a large classification system that classifies different aspects of life into hierarchical categories. This is similar to classification by keywords, but the fact that the keywords belong to a hierarchy enriches the annotations . For example, it can be found that a “room” is a subclass of the class “house”. Ontology consists of entities and their relationships, which may be organized as classes and subclasses, each class may also consist of one or more instances. a framework for ontology enriched semantic annotation of CCTV video is proposed. Visual and text semantics are linked with appropriate keywords provided by domain experts. Video segmentation is done to find moving objects, which are classified as agent, action and recipient. These visual semantics are annotated by keywords of CCTV ontology. Video annotation based on ontology can also certain rules and/or machine learning. Semantic concept detectors can be linked to corresponding concepts in the ontology. A rule-based method for automatic semantic annotation is used; rule learning is built in SWRL. Concepts' relationship of co-occurrence and temporal consistency of video are used to improve performance of individual concept detectors.

VII. MODEL

Module 1

In the first phase any video, obtained from various sources like YouTube or flicker are transformed, where every video shot is converted into a set of images by extracting one frame every 1sec.



Let V denote the video shot that has been downloaded from YouTube. If the video shot lasts for T seconds, then extraction of a frame is done each second to get a cluster of images, these images are the various segments of the same video. These frames collectively represent the video to be processed.

Let $I = \{I_1, I_2, I_3, \dots, I_T\}$, where $I_i = V |_{\text{time}=i}$, $i = \{1, 2, 3, \dots, T\}$

For each image $I_i = \{I_1, I_2, I_3, \dots, I_T\}$, we extract their low-level features. Finally, we get a set of low-level features for the image cluster $F = \{f_1, f_2, f_3, \dots, f_T\}$. The specific type of visual features.

MODULE 2

A query image cluster $I_q = \{I_{q1}, I_{q2}, \dots, I_{qT}\}$ is extracted from the query video shot V , it includes n number of images. The features of images in this cluster are $F_q = \{f_{q1}, f_{q2}, \dots, f_{qT}\}$. Candidate sentence elements that may describe video contents are found first. These candidate sentence elements include four areas :
Object, Event, Scene, Adjective.

For the element object, event, scene and adjective, series of images are downloaded to explain its content, respectively:

These data sets are also set of images, from which comparison and selection of best element is done.

Object image data set (OI), event image data set (EI), scene

image data set (SI), and adjective image data set(AI)

Each image in OI, EI, SI, and AI has only one tag.

MODULE 3

Element selection with Weighted scoring Algorithm

We have weighted scoring algorithm to select best element in given frames. In terms of image cluster

$$I^q = \{I_1^q, I_2^q, \dots, I_T^q\}$$

We have obtained four sets of candidate sentence element

$$CE^c = \{CE_i^c\}_{i=1}^T, c = \{o, e, s, a\}.$$

From this we try to select tag that can best describe the cluster's content with calculating the relevance score.

Let the key frame in this video shot be I_k^q . Clearly, the tag I_k^q of should be given with the highest weight. I_{k+1}^q , I_k^q , and I_{k-1}^q are adjacent to each other in time sequence. Taking the temporal consistency of the video contents into consideration, the weights of images near I_k^q should be given with higher weights, and the image far from I_k^q .

MODULE 4

Selected elements are filtered and refined in this module. Refining of selected elements by a correlation graph algorithm is done. Relationships among these three elements object, event, and scene into consideration. Correlation between various elements is done in over here.

A graph consists of a set of nodes and a set of edges that connect the nodes. We model these three elements by a full connected undirected graph with only three nodes. The edge between two nodes measures their semantic correlation modeled by the normalized Google distance (NGD). The NGD is a semantic similarity measure derived from the number of hits returned by the Google search engine for a given set of keywords. Keywords with the same or meanings in a natural language sense tend to be close in units of NGD, while words with dissimilar meanings tend to be farther apart. Specifically, the NGD between two search terms x and y is –

$$NGD(x,y) = \frac{\max[\log f(x), \log f(y)] - \log f(x,y)}{\{\log N - \min[\log f(x), \log f(y)]\}}.$$

where N is the total number of web pages searched by Google; $f(x)$ and $f(y)$ is the number of Web pages containing search terms x and y , respectively; and $f(x, y)$ is the number of web pages on which both x and y occur.



ISSN: 2350-0328

International Journal of Advanced Research in Science, Engineering and Technology

Vol. 5, Issue 1 , January 2018

VIII. CONCLUSION

In this paper propose main approach for sentence generation from video. In this paper make a good use of well labeled image data set to find sentence related query image element.

It include four main part, object (the subject of the sentence), event (the action), scene (the place where the action happens), and adjective (modifier of the scene). Perform this experiment on more images with accurate tags in data set. And performance can be more satisfactory. Main two algorithm are effective based on the our experiment and how the weighted scoring algorithm and correlation graph algorithm. Need to collect no of image element in order to explain more no of videos.

REFERENCES

- [1] A. Altadmri and A. Ahmed, "A framework for automatic semantic video annotation: Utilizing similarity and commonsense knowledge bases," *Multimedia Tools Appl.*, vol. 72, no. 2, pp. 1167–1191, Mar. 2013.
- [2] A. Barbu et al., "Video in sentences out," in *Proc. UAI*, 2012, pp. 102–112.
- [3] Y. Yang, Y. Yang, Z. Huang, and H. T. Shen, "Transfer tagging from image to video," in *Proc. 19th ACM Int. Conf. MM*, 2011, pp. 1137–1140.
- [4] K. Yang, X.-S. Hua, M. Wang, and H.-J. Zhang, "Tag tagging: Towards more descriptive keywords of image content," *IEEE Trans. Multimedia*, vol. 13, no. 4, pp. 662–673, Aug. 2011.
- [5] TRECVID: TREC Video Retrieval Evaluation, accessed on 2005. [Online]. Available: <http://www.nlpir.nist.gov/projects/trecvid>
- [6] M. Wang, X.-S. Hua, R. Hong, J. Tang, G.-J. Qi, and Y. Song, "Unified video annotation via multigraph learning," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 19, no. 5, pp. 733–746, May 2009.
- [7] M. Wang, X.-S. Hua, Y. Song, X. Yuan, S. Li, and H.-J. Zhang, "Automatic video annotation by semi-supervised learning with kernel density estimation," in *Proc. 14th Annu. ACM Int. Conf. MM*, 2006, pp. 967–976.
- [8] E. Moxley, T. Mei, X.-S. Hua, W.-Y. Ma, and B. S. Manjunath, "Automatic video annotation through search and mining," in *Proc. IEEE ICME*, Apr./Jun. 2008, pp. 685–688.
- [9] J. Tang, X.-S. Hua, M. Wang, Z. Gu, G.-J. Qi, and X. Wu, "Correlative linear neighborhood propagation for video annotation," *IEEE Trans. Syst., Man, Cybern. B, Cybern.*, vol. 39, no. 2, pp. 409–416, Apr. 2009.
- [10] A. Ulges, C. Schulze, D. Keysers, and T. Beuel, "Content-based video tagging for online video portals," in *Proc. 3rd MUSCLE ImageCLEF Workshop Image Video Retr. Eval.*, 2007, pp. 1–4.