# Analytical Review on Textual Queries Semantic Search based Video Retrieval

**Prof. Suvarna L. Kattimani [1], Miss. Saba Parveen Bougdadi [2]**

Assistant Professor, Department of Computer Science and Engineering, B.L.D.E.A's Dr. P.G. Halakatti College of

Engineering and Technology, Vijayapur, Karnataka, India[1]

PG Scholar, Department of Computer Science and Engineering, B.L.D.E.A's Dr. P.G. Halakatti College of Engineering

and Technology, Vijayapur, Karnataka, India[2]

**Abstract:** Semantic search based video retrieval is hard problem due to limited set of vocabulary. Textual queries semantic search based video retrieval is able to detect object from video frames, first the sentence is converted into parse tree than noun, verb, adverb present in that sentence is converted into semantic graph. And construct the semantic meaningful graph gives the semantic structure and matches the nouns, verb and adverb detected in the video frame and also detect the action and position of the object by using semantic meaningful graph. The developed approach is to first sentence is converted into parse tree and object detection takes place. First textual query is matched to concept. The advantage of textual queries approach is object appearance and motion by using structure prediction. Textual queries semantic search can contain temporal and spatial information about multiple objects like trees and building present in the scene.

**Keywords**: Parse tree, Textual queries, Semantic graph, Semantic structure, Structure prediction.

## 1. INTRODUCTION

The challenge is to perform video search retrieval by giving a semantic query. The main goal is to understand image semantically to retrieve or relevant candidates for these queries. And parsing of images is an extremely difficult.The main task is to perform textual queries semantic search based retrieval of videos. The benefit of retrieving video in large available data and has application in improving driver safety. This textual queries semantic search gives video frame and the associated description. And construct the semantic meaningful graph gives the semantic structure and matches the nouns, verb and adverb detected in the video frame and also detect the action and position of the object by using semantic meaningful graph. The developed approach is to first sentence is converted into parse tree and object detection takes place. First textual query is matched to concepts, which in turn uses linear program. The advantage of textual queries approach is object appearance and motion by using structure prediction.

Textual queries semantic search can contain temporal and spatial information about multiple objects like trees and building present in the scene.Videos are retrieves from database via bag of words. Content based image retrieval is a technique for image retrieval which is depends on feature of content of the image. And detect the object from video frame by giving semantic graph.Concept based is another method to retrieve videos and this method is based on a set of concepts detector which bridges the semantic query gap. Concept bank is defined as each concept in the concept bank is treated as node or leaf in the collection of the images known as imagenet. Corresponding particular image name is generated from the collection of the words knows as wordnet. Continuous word space bridges the "semantic query gap".

## 2. LITRARTURE SURVEY

In [1] Dong Wang et al., introduces the mapping concept for automatic generation of video, here first video query is mapped to concepts. Using this mapping technique it is possible to determine the text and visual information are solved by vector model. Also determine the relevant concept of a given query.

In [2] Li Zhang et al., presents a technique for global data for multi-object tracking. This technique proposed a network flow model which is based on optimization theory in which data is associated with multiple object tracking. This optimal method is based on cost of that particular node. A solution is based on iterative approach. This technique is efficient, and finally performance is comparison done between two results. The name of algorithm is polynomial global data links are helpful in reducing fragments of videos.

In [3] G. M. Snoek et al., review on concept based video retrieval, create problem for video retrieval by using old method because it create semantic gap. To overcome the semantic gap a technique is used which is based on concept, every concept is present in conceptnet. It uses machine learning and human computer interaction. Each having different characteristics of data, tasks, and creating baseline experiments.

In [4] Li-Jia Li et al., review on scene understanding and segment classification in automatic framework propose a model that classify scene and identify each segments of the object. Model performs three tasks in one framework. For example game consist of human, ball, grass etc. this model explain images through optical textual model.

The learning process is fully automatic can able to learn scene from data such as images and tags from Flicker. And demonstrate the framework by segmenting image, in which simultaneously objects recognition and segmentation takes place. This model provides a probabilistic tags and automatic training of images. Include classifying, labelling, segmenting complex scene in future they try to capture geometry and appearance of objects in contextual relationship.

In [5] Simon Ting et al., presents a pattern recognition technique for decision making system to make intelligent cricket decision. It uses snikometer used to detect and analyse the signals coming from different video frequency. This snikometer is used to make better decision in matches. The different objects used in sports are detected with slow motion video to make correct decision.it improve accuracy in pattern recognition.

In [6] Ali Farhadi et al., review on image tells a story, in which sentences are generating from images. Humans can able to describe the image and focusing on what is main in the picture, and demonstrate automatic method by score linking an image to a sentence. Which uses model based on AND-OR method. And gives the description of scene generated from visual data. It provides intermediate space of meaning and then generated the sentences for images. The mapping of images into meanings takes place.

It represents triplets in the sense object, action and scene, Uses linear combination. Object stand for 'O' and action stand for 'A' and scene stand for 'S'. Involve the learning of nodes and edges.

In [7] Pablo F. Alcantarilla et al., review onlarge-scale construction of image from stereo imagery proposed the method stereo imagery in which camera motion and action measurements are known. This technique is able to reconstruct image accurately in urban environment. The advantage of stereo matching method is able to build dense and accurate maps. It uses efficient data technique to perform geometric and photometric in validation.

Finally grid filtering technique is used for storage requirement. In this method automatically discards the obstacles from image. And compare the approach with respect to PMVSand Stereo Scan. It uses the efficient large scale stereo matching method provides high quality inequality maps. And require low storage requirements. In future incremental localization and mapping approaches will be used.

In [8] Hamed Pirsiavash et al., review on optimal algorithm for tracking number of object. They analyse the problem of multi-object tracking in video frames in sequence, which contain birth and death states. The problem is solved by using greedy algorithm involve shortest path computation. It also uses tracking algorithm which gives linear time in number of objects. This algorithm is fast, simple, and scalable.

In [9] Dr. Tariq Mahmood et al., presents a technique A-Eye for the game, that is able to decide who is out in the run-out situation. Which require the correct decision of the third umpire, which is more accurate and has the potential to minimize the human error to make correct decisions,Also increases the rating performance. But have some limitations is that it cannot operate at certain height from the ground

In [10] Andreas Geiger and Philip Lenz et al., Presents a technique for autonomous driving using benchmark suite and visual recognition system comes under robotics applications. This method takes the advantage of autonomous driving platform to develop benchmark for stereo flow.

In [11] Nathan silberman et al., presents the indoor segmentation from RGBD image. Physical support to interact with object, and technique used is object recognition. Working with RGB plus depth restricted to indoor scenes. Also provide larger variation in scene. RGB image is converted into segmentation which support inference and scene parsing consist of input goes to major surface uses integer program formulations.

In [12] Marcus Rohrbach et al., introduces the concept of Translating Video Content to Natural Language in which humans uses natural language to communicate visual perception which provide description of natural language for

visual content. And generate semantic representation of visual content include object. And the concept of translating video content uses CRF models. CRF prediction is done by predict the hamming distance. Here video is translated into text and also image into text.

In [13] J. Dalton et al., review on zero-shot video defined as where no training data is given and queries consist of text extracted from frame. And also able to recognized in speech of audio. Source extracted to build textual representation, and uses MRF which means markov random field. Zero-shot video holds both text and semantic video concepts. It uses zero-shot video retrieval technique which requires no training data. Zero-shot retrieval uses queries for video retrieval holds semantic concepts.

In [14] Sanja Fidler et al., Presents a sentence is equivalent to a maximum number of pixels which uses holistic scene for understanding where images are present. Proposed a holistic conditional random field model about which object are present in the particular scene and gives image information as input.

And this technique is use to detect object from the sentences by using ranking candidate detection. Holistic model gives information about text in the form of sentence.

Future work plan is to utilize the textual information in webpages e.g., Wikipedia.

In [15] Cynthia Matuszek et al., Introduced a technique named joint model for grounded attribute learning in which robots become more popular. The goal is to bring out the meaning of language which is linked with physical world which is useful in learning of languages for grounded attribute. Uses latent variable concepts where learning is performed by using training algorithm, which is able to learn accurate languages to identify the shape and colour of the objects.

In [16] Richard Socher et al., review on parsing with compositional grammar uses natural languages with small set of discrete vectors. And improve the lexical phrasesAnd splitting of address, also introduce the compositional vector grammar. And it is fast to train and improve the performance of semantic information.

In [17] Ashok Kumar et al., review on cricket activity detection is used to identify aspects of such contents like batting. This technique splits into shots which combine with segments to identify the batsman stroke, which require optical flow analysis. Stroke is classified into four directions.If classify the view and detect the shots, that identify the directions and achieved the good accuracy. It detects the highlighted detection based on audio. This approach uses large datasets for event detection.

In [18] Chen Kong et al., Presents a technique of co-referencing text to image.co-reference of text to image means two or more expression in a text refer to the same image. It uses natural language for RGB scene in order to improve parsing speed. Particularly each object each is refers as noun in the image, which allows statistical machine translation for visual information. It also proposed a structure prediction model.

Uses lingual descriptions to improve visual scene, in which visual object the text is reference. Also improve classification accuracy and reliable. Proposed model is holistic of object detection and scene categorised between text and visual object. In future work plan to employ information over a large set of objects.

In [19] Rahul Anand Sharma et al., introduces the fine annotation of cricket videos, recognition of human activities is problem in video understanding. It uses action recognition technique, which contain small set of actions. Settings are weakly supervised, contain actions and semantic descriptions. First step is segmented into scenes to extract information and then classify the video shots, means phrase is mapped to video shots. Involve machine learning algorithm for complex actions.

In [20] D. Karmaker et al., review on cricket shot classifications, different sports objects shots cannot be detected from single video and without multiple view camera and extracting features from videos is difficult task. Human body parts created several movement is different directions of optical flow which is related to motion using 3d action, and angle ranges to detect important shots on motion vector to measure the angle and also provide accuracy.

In [21] Ali javed et al., introduces a technique on hybrid approach on summarization of videos propose automatic method for detection of video in matches in which large set of multimedia content is available. The framework is designed based on video summarization in large dataset which is reliable with low bandwidth.

ANAYSIS ON TEXTUAL QUERIES SEMANTIC SEARCH BASED VIDEO RETRIEVAL

TABLE1. TEXTUAL QUERIES SEMANTIC SEARCH BASED VIDEO RETRIEVAL

| Sl. No | Area of objective | Author | Year | Major contribution | Method Used |
|---|---|---|---|---|---|
| 1 | Query concept mapping for video retrieval | D. Wang et al. | 2007 | Automatic generation of concept of a given query | Vector computational model |
| 2 | Data association for multi-object tracking | L. Zhang etal. | 2008 | Global data links are helpful in reducing fragments of video | Network flow model |
| 3 | Concept based video retrieval | G.M. Snoek etal. | 2008 | Improve machine learning and human computer interaction | Concept retrieval technique |
| 4 | Scene understanding, classification and segmentation | L. Li etal. | 2009 | Simultaneously object recognition and segmentation takes places | Automatic classification framework |
| 5 | Pattern recognition technique for cricket decision making | Simon Ting et al. | 2009 | Used to make better decision in matches | Sniko-pattern recognition technique |
| 6 | Generating sentences from pictures | A. Farhadi etal. | 2010 | Represents triplets action, object and scene | Automatic linking method |
| 7 | Large scale construction of image | F. Alcantarillaet al. | 2010 | Perform photometric and geometric validation | Stereo imagery technique |
| 8 | Optimal algorithm for tracking number of object | Hamed pirsiavash etal. | 2010 | Analyse the problem of multi-object tracking | Globally-optimal greedy algorithm |
| 9 | A-Eye for the game | Dr. Tariq Mahmood et al. | 2011 | Increase rating performance provide more accurate result | A-Eye third umpire decision |
| 10 | Ready for autonomous driving? Kitti benchmark | A. Geigeret al. | 2012 | Used in robotics applications | Stereo flow model |
| 11 | Indoor segmentation from rgbd image | N. Silberman etal. | 2013 | Provide physical support to interact with objection | Integer program formulation |
| 12 | Translating video content to natural language | M. Rohrbach etal. | 2013 | Used to communicate visual perception | Conditional random field model |
| 13 | Zero shot video retrieval concepts | J.Daltonet al. | 2013 | Automatic detect relevant concepts given in text query | Zero shot video retrieval technique |
| 14 | Single sentence contain large number of pixels | S. Fidleretal. | 2013 | Gives image information | Holistic CRF model |
| 15 | Joint model for grounded attribute | C. Matuszeket al. | 2013 | Able to learn accurate language to identify the shape and colour of the object | Grounded attribute learning |
| 16 | Parsing use compositional vector grammar | R. Socher et al. | 2013 | improve the performance of semantic information | Discrete vector lexical phrases |
| 17 | Sports activity detection | Ashok Kumar et al. | 2014 | Classify the view and detect the shots | Optical flow analysis |
| 18 | Co-referencing text to image | Chen Kong et al. | 2014 | Extremely useful for robotics application | Object retrieval from image to text |
| 19 | Fine grain annotation of sports video | R. Sharma et al. | 2015 | Human activity action recognition | Machine learning algorithm |
| 20 | Cricket shot classification using motion vector | D. Karmakeret al. | 2015 | Detect important shots in different direction | Motion vector classification method |
| 21 | A hybrid approach for summarization | Ali Javed et al. | 2016 | Detect videos from large form of multimedia | Automatic framework technique |

## 3. CONCLUSION

Textual queries semantic search able to solve the problem of semantic using natural language queries and can able to detect objects in the video frame with high accuracy by parsing into semantic graph then the graph is matched to the optical concepts, thus concluded that matching between the nouns and objects detected in video. In future, plan to improve the output of tracked object in selected video frames.

## REFERENCES

[1] D. Wang, X. Li, J. Li, and B. Zhang. The importance of query-concept-mapping for automatic video retrieval. In Proc. of ACM Multimedia, 2007

[2] L. Zhang, Y. Li, and R. Nevatia. Global data association for multi-object tracking using network flows. In CVPR'08

[3] G. M. Snoek and M. Worring. Concept-based video retrieval.Foundations and Trends in information Retrieval, 2008.

[4] L. Li, R. Socher, and L. Fei-Fei. Towards total scene understanding: classification, annotation and segmentation in an automatic framework. In CVPR, 2009.

[5] Simon Ting and M. V. Chilukuri. Novel Pattern Recognition Technique For An Intelligent Cricket Decision Making System I2MTC 2009.

[6] A. Farhadi, M. Hejrati, M. Sadeghi, P. Young, C. Rashtchian, J. Hockenmaier, and D. Forsyth. Every picture tells a story: Generating sentences for images. In ECCV, 2010.

[7] Pablo F. Alcantarilla et al., review onlarge-scale construction of  image, In 2010.

[8] Hamed Pirsiavash et al., Optimal algorithm for tracking number of object.ICS 2010.

[9] Dr. Tariq Mahmood, Syed Obaid Ahmed. A-Eye: Automating The Role Of The Third Umpire In The Game Of cricket, 2011.

[10] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In CVPR, 2012.

[11] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus. Indoor segmentation and support inference from rgbd images. In ECCV, 2012.

[12] Marcus Rohrbach,Wei Qiu,Ivan Titov, Stefan Thater,Manfred Pinkal: Translating Video Content to Natural Language Descriptions In ICCV ,2013.

[13] J. Dalton, J. Allan, and P. Mirajkar. Zero-shot video retrieval using content and concepts. In CIKM, 2013.

[14] S. Fidler, A. Sharma, and R. Urtasun. A sentence is worth a thousand pixels. In CVPR, 2013.

[15] C. Matuszek, N. FitzGerald, L. Zettlemoyer, L. Bo, and D. Fox. A joint model of language and perception for grounded attribute learning. In ICML, 2013

[16] R. Socher, J. Bauer, C. D. Manning, and A. Y. Ng. Parsing with compositional vector grammars. In ACL, 2013.

[17] Ashok Kumar, Javesh Garg and Amitabha Mukerjee. Cricket activity detection, IPAS 2014.

[18] C. Kong, D. Lin, M. Bansal, R. Urtasun, and S. Fidler. What are you talking about? text-to-image coreference. In CVPR, 2014.

[19] Rahul Anand Sharma, C. V. Jawahar. Fine grain annotation of cricket videos, IAPR 2015.

[20] D Karmaker, Chowdhury†, M S U Miah. Cricket Shot Classification Using Motion Vector, 2015.

[21] Ali Javed, Khalid bashir bajwa, hafiz malik.A hybrid Approach for summarization of circket videos. ICCE 2016.