# Review Paper on Video Content Analysis into Text Description

Vandana D. Edke
M.E 2nd Year student Dept. Computer Engg, Dr. DY Patil School of Engg and Technology, Pune.

Ramesh M. Kagalkar
Research Scholar and   Asst.Professor, Computer  Engg. Department, Dr. D Y Patil School of Engg and Technology, Pune.

## ABSTRACT
This paper reviews video content analysis from the various situations into matter version. The totally different researchers are applied different technique to unravel the approaches. It is a tendency to tend to jointly obtaining down addressing the required down siting extracting the frames from video, comparison the frames; pattern matching and generating the corresponding text description is address here. Hence additionally created a discussion, observation and comparison of quick work applied during this work. It is a tendency to mix the output of progressive object and activity detectors with "real-world" data to pick the foremost probable subject-verb-object triplet for describing a video. It is a tendency to show that this data, mechanically well-mined from web-scale text corpora, hence projected choice rule by providing it discourse information and results in a four-fold increase in activity identification. In contrast to previous ways mentioned in literature survey, therefore in this approach will annotate absolute videos while not requiring the high-priced assortment and annotation of an analogous coaching video corpus.

## Keywords
Natural language generation, concept hierarchy, semantic primitive, position/posture and estimation of human Case frame.

## 1. INTRODUCTION
Recognizing activities in real-world videos may be a difficult AI downside with several sensible applications [1, 16]. Hence a tendency to gift a unique approaches to with efficiency constructing activity recognizers by effectively combining three various techniques. First, it is a tendency to use natural-language descriptions of videos as "weak" supervising for coaching activity recognizers [15]. To mechanically develop a group of activities at the side of a tagged coaching corpus by cluster the verbs utilized in sentences describing videos. Second use antecedently trained object recognizers to mechanically observe objects in video and use this data to assist establish connected activities [14]. As an example, sleuthing "horse" within the image helps classify the activity as "riding". Third, we have a tendency to mine an

Oversized corpus of generic, raw natural-language text to find out the correlations between activities (verbs) and their connected objects (nouns). By mining an oversized corpus and collection statistics on however seemingly completely different verbs attach to specific nouns and estimate the chance of specific activities given specific objects. Integration these three ways permits for the fast development of fairly correct activity recognizers while not ever expressly providing coaching labels for videos. By combining text mining to each mechanically infer tagged activities and extract relevant world-knowledge connecting activities and objects, at the side of computer-vision techniques for each object and activity

recognition, work demonstrates the utility of integration ways in tongue process and pc vision to develop effective AI systems. Experiments on a sizeable corpus of You Tube videos annotated with natural-language descriptions [3] verify that our approach improves the accuracy of a typical activity recognizer for real-world videos. Figure 1 shows sample frames from some of videos with their linguistic descriptions.

The remainder of the paper is organized as follows. In section 2 Literature survey of various work carried out, observations of its techniques and results, comparison of information sets, features; technique utilized by completely different analysis teams acting on video content description in text is additionally given. In Section 3 describes the proposed system design well and steps needed to unravel the matter. In final Section 4 presents conclusions and discussion.

## 2.  LITERATURE SURVEY
Video activity recognition has become a full of life space of analysis in recent years [8, 32, 35]. However, the set of activity categories square measure continuously expressly provided a tendency to mechanically discover the set of activities from matter descriptions. Scene context [26] and object context [14, 27, 31, 36, 37] has antecedently been accustomed aid activity recognition. However most of this previous work uses a awfully unnatural set of activities, whereas we have a tendency to address a various set of activities in real-world You Tube videos. Also, in contrast to previous work, and tendency to mechanically extract correlations between activities and objects from an oversized text corpus. There has been work victimization text related to videos within the sort of scripts or closed captions to assist activity recognition [10, 20, 19, 4, 15, 36]. However, these strategies do not use deeper tongue process. By contrast; we have a tendency to demonstrate the advantage of full parsing of associate unrelated corpus to mine noises connecting objects and activities.

A particular connected project that uses natural-language descriptions to mechanically annotate videos with activity labels. However, in [19], the set of activity categories square measure pre-specified, whereas a tendency to mechanically generate activity categories from matter descriptions by cluster verbs victimization WorldNet because the solely supply of previous data or direction. Also, the approach in [19, 37] needs coaching set within which linguistic descriptions square measure annotated with the pre-specified activities, whereas our approach doesn't need any specially annotated text. whereas current approach uses coaching videos every represented by many completely different natural-language sentences provided by multiple human annotators recruited on Amazon Mechanical Turk uses freely accessible picture show scripts downloaded from the web. In table 1 show offers the detail review of observation and Comparison of quick work meted out thus far.

# 3. PROPOSED SYSTEM ARCHITECTURE

In figure 1 show planned system design. Wherever it takes a video clip as input and generates result as language description. At present, it is not very easy to estimate correct posture and motion of extremely articulated object like human in real time. For this reason, to assume the subsequent 3 clues which may be obtained by comparatively light-weight processes are enough for normal cases, e.g. shown within the figure one video clip of elephant with folk's activity scenes, to notice a posture of a human [34]

- Position of head implies not solely an edge wherever the human is however additionally a posture whether or not he/she is standing or sitting.

- Direction of head implies what he/she is viewing.

- Positions of hands imply a form of gesture and interaction with objects.

- Movement of elephant implies what the activity goes on.

For additional correct understanding of interaction with an object it is necessary to estimate the relative position and motion of the things and in addition on establish the object.

For each frame of input video pictures, the body, skin regions of a person's and animal structure is extracted by calculative distinction of colours between input and background pictures component by component. Positions of the top and also the hands are found by perspective transformation. Orientation of the top is additionally calculable by calculative correlation between input head region and also the head models with multiple aspects that are ready beforehand. On the opposite hand, action of transferring an object above all is detected in a very separate approach. Examining shapes of the regions of a human and an object appeared on difference images, it can be verified whether the human and elephant action and take it out. Additionally, the things are often known by scrutiny edges and color histograms of extracted object region with those of object models [35].

Next, abstract descriptions of actions are generated for every part by applying domain information, like allocation of apparatus in a very area, to the position/posture of the human obtained from the video pictures. In language, ideas of motion verbs embody variety of linguistics primitives. Which means of a verb tends to be additional concrete and specific because the range of linguistics primitives will increase. Therefore we tend to construct idea hierarchies of actions for everybody components classified by combination of linguistics primitives. By creating correspondence between linguistics primitive of action and a feature extracted from the video pictures, the foremost acceptable predicate, object, etc. are selected. Hence to fill these syntactical elements into a case frame that is commonly used as a linguistics illustration of a sentence within the space of language process.

Finally, analysis of frames activity in body components is integrated into a frame expressing total body action. Applying syntactical rules and natural word lexicon, the frame activity is translated into a language sentence.
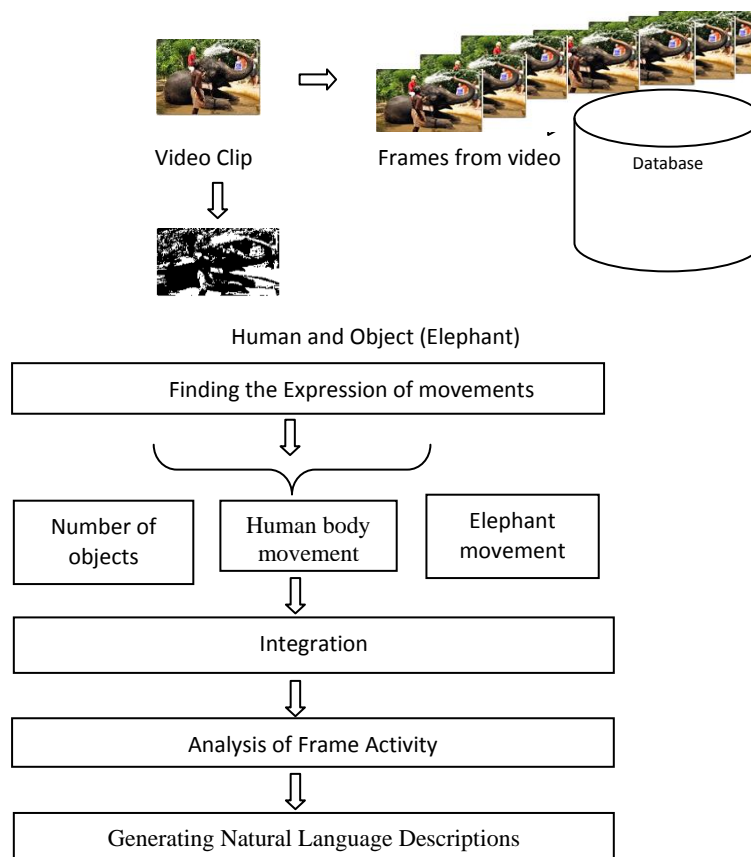


**Figure1: Shows the proposed system architecture.**

**Table 1.Shows observations and comparison of data sets, features; technique used by different research groups working on video content description in text.**

| Sr. No. | Citation | Work carried out | Conclusion |
|---|---|---|---|
| 1 | [20] Dong X u ,Shih-Fu Chang | Single level EMD, Multi frame based detection. | Event recognition of unconstrained broadcast news videos Has Discussed. |
| 2 | [21]Nivedakrishnamoorthy,Girish Malkarnenkar,RaymondMooney, KateSaenko,SergioGuadarrama | Holistic data driven technique is used. | This paper generates a natural language descriptions of short videos by using holistic data driven approach |
| 3 | [22]Motwani,T and Mooney,R | Text mining is used to learn the correlations between these verbs and related objects. Cluster verb is used to describe videos to discover classes of activities. | It shows that automatically discovering activity classes from natural language description of videos. |
| 4 | [23] Marcus Rohrbach Wei Qiu Ivan Titov Stefan Thater Manfred Pinkal BerntSchiele | Conditional Random Field is used to predict the semantic representation. | It uses semantic representation of visual content including object and activity labels and generates sentences as target language. |
| 5 | [24] Muhammad Usman Ghani Khan YoshihikoGotoh | ROUGE score, conventional image processing features extraction, is used to describe the videos. | The natural language description is generated from video stream by using CFG and by using rough score evaluation is done between human annotated and machine generated sentences. |
| 6 | [25] Andrei Barbu, Alezander Bridge Zachary Burchill, Dan Coroian Sven Dickinson, Sanja Fidler Aaron Michaux, | Clustering, Smoothing, hidden markov models is used for sentence generation of video. | It uses coordination, verbs,, nouns, adjectives,  prepositions, lexical pps, determiners, particles, pronouns, adverbs ,auxiliary and  appropriate sentences are generated. |
| 7 | [26]Ding,D.;Metze,F.;Rawat,S.;S chulam,P.;Burger,S.;Younesian,E. ;Bao,L.;Christel,M.;and Hauptmann, | Topic Oriented Multimedia Summarization (TOMS). | Automatically generates a paragraph of natural language by using topic oriented multimedia summarization system. |
| 8 | [27]Lee,M.;Hakeem,A.;Haering, N.;andZhu,S | These methods such as SAVE, VEML HPSG are used for annotates a large collection of videos. | SAVE framework that provides an end-to-end automatic system for parsing video, extracting visual event content, and providing semantic and text annotation. |
| 9 | [28]J.Sivic,M. veringham, and A.Zisserman | Seamless integration, detection, recognition is used for detecting person specific classifier from video. | Learning person specific classifiers Using kernel combination can improve the accuracy of automatic naming of characters in TV Video. |
| 10 | [29]S.Guptaand R.J.Mooney | For video activity recognition leave-one-game-out is used. | Automatically train a video activity recognizer without requiring any manual labeling of video clips by using closed caption. |
| 11 | [30] D.Xu  and S-F.Chang | Kernel based method is used for visual event recognition of news videos. | Here the problem of visual event recognition of unconstrained broadcast news videos has described. |
| 12 | [31] A.Hakeem and M.Shah | For recognizing multiple events of videos event browsing annotation indexing is used. | Detecting events In video involving multiple agents and their interaction was identified. |
| 13 | [32]Kojima,T.Tamura,and K.Kukunaga | Concept hierarchy of action, semantic representation is used for natural language description of human activities | In this work appropriate verb and objects can be selected in a proper manner. |
| 14 | [33] Douglas Ayers and Mubarakak Shah | Scene detection, Scene change detection, tracking is used for monitoring human behavior. | Human actions can recognize such as entering, picking up a phone. etc |

## 4. CONCLUSION

In this way proposed work has been compared and contrasted utterly totally different approaches for video content analysis to come up with text description. Thus this paper tends to finally created three necessary contributions to video activity recognition. First, it's introduced a unique technique for mechanically discovering activity categories from natural-language descriptions of videos. Second, it is incontestable however an existing activity recognition system is improved victimization object context along with correlations between objects and activities. Third, shows however language process is accustomed mechanically extract the requisite information concerning the correlation between objects and activities from a corpus of general text. In this way paper tend to process the means of reworking video pictures drawn as geometrical numerical information into matter descriptions as abstract information. We tend to propose a brand new construct of actions to extraction of meanings of pictures by substantiative correspondence between linguistics options of human actions and therefore the linguistic communication ideas. Consequently, applicable verbs and objects are designated in a very subtle means.

## 5. REFRENCES

[1] J. K. Aggarwal and S. Park, "Human motion: Modeling and recognition of actions and interactions', in 3DPVT, (2004).

[2] Marie Catherine De Marneffe, Bill Maccartney, and Christopher D.Manning, 'Generating typed dependency parses from phrase structureparses', in LREC, (2006).

[3] David L. Chen and William B. Dolan, 'Collecting highly parallel data for paraphrase evaluation', in ACL, (2011).

[4] Timoth´eeCour, Chris Jordan, EleniMiltsakaki, and Ben Taskar, 'Movie/script: Alignment and parsing of video and text transcription', in ECCV, (2008).

[5] Ke Chen and Kazuaki Maeda David Graff, Junbo Kong,'English gigaword second edition', in LDC, (2005).

[6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-FeiLi, 'Imagenet: A large-scale hierarchical image database', in CVPR,(2009).

[7] Pedro Domingos and Michael Pazzani, 'On the optimality of the simple Bayesian classifier under zero-one loss, ML,(1997).

[8] Alexei A. Efros, Alexander C. Berg, Er C. Berg, Greg Mori, and JitendraMalik, 'Recognizing action at a distance', in ICCV, (2003).

[9] Mark Everingham, Luc J. Van Gool, Christopher K. I. Williams, John M. Winn, and Andrew Zisserman, 'The Pascal visual object classes (voc) challenge', IJCV, (2010).

[10] Mark Everingham, Josef Sivic, and Andrew Zisserman, 'Hello! Myname is... buffy" – automatic naming of characters in TV video', inBMVC, (2006).

[11] P. F. Felzenszwalb, R. B. Girshick, and D. McAllester.Discriminativelytrained deformable part models, release 4.

[12] http://people.cs.uchicago.edu/ pff/latent-release4/.Pedro F. Felzenszwalb, Ross B. Girshick, David A. McAllester, and Deva Ramanan, 'Object detection with discriminatively trained part based models', IEEE Trans. Pattern Anal. Mach. Intell., (2010).

[13] AdrienGaidon, MarcinMarszalek, and CordeliaSchmid, 'Mining visual actions from movies', in BMVC, (2009).

[14] Abhinav Gupta and Larry S. Davis, 'Objects in action: An approach for combining action understanding and object perception', in CVPR, (2007).

[15] Sonal Gupta and Raymond J. Mooney, 'Using closed captions as supervision for video activity recognition', in AAAI, (2010).

[16] AnthonyHoogs and A. G. AmithaPerera, "Video Activity Recognition in the Real World.' in AAAI, (2008).

[17] Jay J. Jiang and David W. Conrath, 'Semantic similarity based on corpus statistics and lexical taxonomy', CoRR, (1997).

[18] Ivan Laptev, "On space-time interest points", IJCV, (2005).

[19] Ivan Laptev, MarcinMarszalek, CordeliaSchmid, and Benjamin Rozenfeld, 'Learning realistic human actions from movies', in CVPR,(2008).

[20] Dong X u,Shih-Fu Chang, "Video Event Recognition Using Kernal Methods With multilevel Temporal Alignment".IEEE (Transaction on pattern analysis and Machine intelligence, Vol. 30 No.11, 2008.

[21] Niveda krishnamoorthy, Girish Malkarnenkar, RaymondMooney,KateSaenko,SergioGuadarrama,"Gene rating Natural –Language Video Descriptions Using Text-Mined Knowledge". Association for a advancement of Artificial Intelligence 2013.

[22] Motwani,T., and Mooney, R."Improving Video Activity Recognition Using Object Recognition and Text Mining".In European Conference On Artificial Intelligence (ECAI) ,2012.

[23] Marcus Rohrbach, Wei Qiu, Ivan TitovStefan, Thater Manfred , Pinkal BerntSchiele," Translating Video Content To Natural Language"proc.of IEEE International Conference On Computer Vision( ICCV)Dec.2013.

[24] Muhammad Usman Ghani Khan, YoshihikoGotoh "Describing Video Contents in Natural Language" Proceeding of the workshop on Innovative Hybrid Approaches to the processing Of Textual data (Hybrid2012) EACL 2012.

[25] Andrei Barbu, Alexander Bridge, Zachary Burchill,Dan Coroian, Sven Dickinson, Sanja Fidler, Aaron Michaux, Sam Mussman,"Video In Sentence Out"2012.

[26] Ding,D.;Metze,F.;Rawat,S.;Schulam,P.;Burger,S.;Youne sian,E.;Bao,L.;Christel,M.;and Hauptmann,"A.Beyond Audio and Video Retrival:Towards multimedia summarization".In Proceeding of the 2nd ACM International Conference on multimedia Retrival.2012.

[27] Lee,M.; Hakeem, A.; Haering, N.; and Zhu, S."Save: A Framework For Semantic Annotation of Visual Events".In IEEE Computer Vision and Pattern Recognition Workshops (CVPR-W), 2008.

[28] J.Sivic, M.Everingham, and A.Zisserman,"Who Are You?-Learning Persons Specific Classifiers from Video,"proc.IEEEConf. Computer Vision and pattern Recognition, 2009.

[29] S.Gupta and R.J.Mooney,"Using Closed Caption as Supervision for Video Activity Recognition"Proc.24th AAAI Conf.Artificial Intelligenc, pp1083-1088, july 2010.

[30] D.Xu and S-F.Chang."Visual Event recognition in News Video Using Kernel Methods With Multi-Level Temporal Alignment,"Proc.IEEE Conf.Computer Vision And Pattern Recognition, 2007.

[31] A.Hakeem and M.Shah,"Learning Detection and Representation of multiple Agent Events in Videos" Artificial Intelligence Journal, 2007.

[32] A.Kojima,T.Tamura,and.Kukunaga,"Natural Language Description of Human Activities of Video Image Based on Concept Hierarchy of action,"International Journal of computer vision, vol.50, pp.171-184, 2002.

[33] Douglas Ayers, Mubarakak Shah," Monitoring Human Behavior from Video Taken in an Office Environment", Image and vision computing 2001.

[34] Ramesh.M.Kagalkar,Mrityunjaya.V.Latteand Basavaraj.M.Kagalkar "Template Matching Method For Localization Of Suspicious Area And Classification Of Benign Or Malignant Tumors Area In Mammograms", International Journal on Computer Science and Information Technology (IJCECA), ISSN 0974-2034, Vol.25, Issue1, 2011.

[35] Ramesh.M.Kagalkar Mrityunjaya .V. Latte and Basavaraj.M.Kagalkar ""An Improvement In Stopping Force Level Set Based Image Segmentation", International Journal on Computer Science and Information Technology(IJCEIT), ISSN 0974-2034,Vol 25,Issue1,Page 11-18,2010.

[36] Mrunmayee Patil and Ramesh Kagalkar "An Automatic Approach for Translating Simple Images into Text Descriptions and Speech for Visually Impaired People", International Journal of Computer Applications (0975 – 8887) Volume 118 – No. 3, May 2015.

[37] Mrunmayee and Ramesh Kagalkar "A Review on Conversion of Image to Text As well As Speech Using

Edge Detection and Image Segmentation" International Journal of Advance Research in Computer Science Management Studies, Volume 2, and Issue 11 (November-2014) publish on 29th November to 30th November 2014.

# 6. AUTHORS PROFILE

**Vandana D. Edke** is M.E 2nd year student of Computer Engg.Department, Dr. DY Patil School of Engg and Technology, Lohegaon, Pune. My main research interest includes Image processing and Gesture recognition.

**Ramesh. M. Kagalkar** was born on Jun 1st, 1979 in Karnataka, India and presently working as an Assistant. Professor, Department of Computer Engineering, Dr.D.Y.Patil School Of Engineering and Technology, Charoli, B.K.Via – Lohegaon, Pune, Maharashtra, India. He has 13.5 years of teaching experience at various institutions. He is a Research Scholar in Visveswaraiah Technological University, Belgaum, He had obtained M.Tech (CSE) Degree in 2006 from VTU Belgaum and He received BE (CSE) Degree in 2001 from Gulbarga University, Gulbarga. He is the author of text book Advance Computer Architecture which cover the syllabus of final year computer science and engineering, Visveswaraiah Technological University, Belgaum. One of his research article "A Novel Approach for Privacy Preserving" has been consider as text in LAP LAMBERT Academic Publishing, Germany (Available in online). He is waiting for submission of two research articles for patent right. He has published more than 25 research papers in International Journals and presented few of there in international conferences. His main research interest includes Image processing, Gesture recognition, speech processing, voice to sign language and CBIR. Under his guidance four ME students awarded degree in SPPU, Pune, five students at the edge of completion their ME final dissertation reports and two students started are started new research work and they have publish their research papers on International Journals and International conference.