

# Content Based Video Retrieval using Web Image of Specific Action

**Rajnor D.S**

Assistant Professor

Department of Computer Engineering  
SNJB's KBJ College of Engineering, Chandwad

## Abstract

Now a day it's very difficult to extract the specific videos with specific action from web. In my proposed system I am going to use the web action image for the Retrieval of corresponding videos shots. Here I am using Textual input at the search tool, first depending on the text enter in search tool we can extract the tag relevant videos of the specific action from the web, Simultaneously we can extract the web images of the specific action entered in the search tool. Then segment each video into the videos shots. After segmentation we will convert each video into separate bags of spatio temporal features [2] and apply visual ranking method to rank each video. On the web image I apply the pose let's to identify human and Non-human action [1], then matching each images with the video shots. Now my expectations will be the most relevant videos as compare to previous work done

**Keywords**—(SVM) support vector machine, (API) application Programming Interface, (EMD) Earth Mover's Distance, (A-MKL) Adaptive Multiple Kernel Learning

## I. INTRODUCTION

All Today's the world is running on Audio visual track .Due to the video database on the web increases then it is difficult to manage all the data in case of you tube and other like Kodak consumer video database social networking website. Tagging each video is very time consuming thing that's why we required an automatic tool that can extract the videos just as providing a keyword. Another thing is that sequencing of these videos required another man power for tag based video categorization. In this proposed system I take Web still images corresponding to given actions into account, with an intuition that the shots which are more similar to Web action images may be more likely relevant shots so that they should be ranked higher . I use here Bing API database for image collection and You tube database for video collection. In the previous work they are used only still images and work only on Human action, in case of non-human action the existing system is not capable. In this System we consider a video shot as a set of consecutive frames which represent a scene. If video shots corresponding to any action can be obtained automatically from Web source, so that building action database will become easier than before.

## II. RELATED WORK

**2.1L.** Bourdev and J. Malik have proposed the method "Body part detectors trained using 3D human pose annotations"[1] In this paper they address the classic problems of detection, segmentation and pose estimation of people in images with a novel definition of a part, a poselet. They postulate two criteria (i) It should be easy to find a poselet given an input image (ii) it should be easy to localize the 3D configuration of the person conditioned on the detection of a poselet. To permit this they have built a new dataset, H3D, of annotations of humans in 2D photographs with 3D joint information, inferred using anthropometric constraints. This enables us to implement a data-driven search procedure for finding poselets that are tightly clustered in both 3D joint configuration space as well as 2D image appearance. The algorithm discovers poselets that correspond to frontal and profile faces, pedestrians, head and shoulder views, among others. Each poselet provides examples for training a linear SVM classifier which can then be run over the image in a multiscale scanning mode. The outputs of these poselet detectors can be thought of as an intermediate layer of nodes, on top of which one can run a second layer of classification or regression. We show how this permits detection and localization of torsos or keypoints such as left shoulder, nose, etc. Experimental results show that we obtain state of the art performance on people detection in the PASCAL VOC 2007 challenge, among other datasets. We are making publicly available both the H3D dataset as well as the poselet parameters for use by other researchers.

**2.2M.** Blank, L. Gorelick have proposed the method "Actions as space-time shapes" [2] in this paper we found that Human action in video sequences can be seen as silhouettes of a moving torso and protruding limbs undergoing articulated motion. They regard human actions as three-dimensional shapes induced by the silhouettes in the space-time volume. They adopt a recent approach by Gorelick et al. (2004) for analyzing 2D shapes and generalize it to deal with volumetric space-time action shapes. Their method utilizes properties of the solution to the Poisson equation to extract space-time features such as local space-time saliency, action dynamics, shape structure and orientation. They show that these features are useful for action recognition, detection and

clustering. The method is fast, does not require video alignment and is applicable in (but not limited to) many scenarios where the background is known. Moreover, they demonstrate the robustness of their method to partial occlusions, non-rigid deformations, significant changes in scale and viewpoint, high irregularities in the performance of an action and low quality video

**2.3** N.I. Cinbins, R.G. Cinbins, and S. Sclaroff have proposed the method on “Learning actions from the Web” [3] this paper proposes a generic method for action recognition in uncontrolled videos. The idea is to use images collected from the Web to learn representations of actions and use this knowledge to automatically annotate actions in videos. Their approach is unsupervised in the sense that it requires no human intervention other than the text querying. Its benefits are two-fold: (1) we can improve retrieval of action images, and (2) we can collect a large generic database of action poses, which can then be used in tagging videos. They present experimental evidence that using action images collected from the Web, annotating actions is possible.

**2.4** Duan, D. Xu, I. W. Tsang, and J. Lu have proposed method on “Visual event recognition in videos by learning from web data” [4] in this paper we found that they proposed a visual event recognition framework for consumer domain videos by leveraging a large amount of loosely labeled web videos (e.g., from YouTube). First, they propose a new aligned space-time pyramid matching method to measure the distances between two video clips, where each video clip is divided into space-time volumes over multiple levels. They calculate the pair-wise distances between any two volumes and further integrate the information from different volumes with Integer-flow Earth Mover’s Distance (EMD) to explicitly align the volumes. Second, they propose a new cross-domain learning method in order to 1) fuse the information from multiple pyramid levels and features (i.e., space-time feature and static SIFT feature) and 2) cope with the considerable variation in feature distributions between videos from two domains (i.e., web domain and consumer domain). For each pyramid level and each type of local features, they train a set of SVM classifiers based on the combined training set from two domains using multiple base kernels of different kernel types and parameters, which are fused with equal weights to obtain an average classifier. Finally, they propose a cross-domain learning method, referred to as Adaptive Multiple Kernel Learning (A-MKL), to learn an adapted classifier based on multiple base kernels and the relearned average classifiers by minimizing both the structural risk functional and the mismatch between data distributions from two domains. Extensive experiments demonstrate the effectiveness of our proposed framework that requires only a small number of labeled consumer videos by leveraging web data.

**2.5** D. H. Nga and K. Yanai have proposed other method i.e. “Automatic construction of an action video shot database using web videos” [5] there are a huge number of videos with text tags on the Web nowadays. In this paper, they propose a method of automatically extracting from Web videos video shots corresponding to specific actions with just only providing action keywords such as “walking” and “eating”. Their proposed method consists of three steps: (1) tag-based video selection, (2) segmenting videos into shots and extracting features from the shots, and (3) visual-feature-based video shot selection with tag-based scores taken into account. Firstly, they gather video IDs and tag lists for 1000 Web videos corresponding to given keywords via Web API, and then calculate tag relevance scores for each video using a tag-co-occurrence dictionary which is constructed in advance. Secondly, they fetch the top 200 videos from the Web in the descending order of the tag relevance scores, and segment each downloaded video into several shots. From each shot they extract spatio-temporal features, global motion features and appearance features, and convert them into the bag-of-features representation. Finally, we apply the Visual Rank method to select the video shots which describe the actions corresponding to the given keywords best after calculating a similarity matrix between video shots. In the experiments, they achieved the 49.5% precision at 100 shots over six kinds of human actions by just providing keywords without any supervision. In addition, they made large-scale experiments on 100 kinds of action keywords.

**2.6** Action snippets: How many frames does human action recognition require? In this paper we found that Visual recognition of human actions in video clips has been an active field of research in recent years. However, most published methods both analyze an entire video and assign it a single action label, or use relatively large look-ahead to classify each frame. Contrary to these strategies, human vision proves that simple actions can be recognized almost instantaneously. In this paper they present a system for action recognition from very short sequences (snippets) of 1-10 frames, and systematically evaluate it on standard data sets. It turns out that even local shape and optic flow for a single frame are enough to achieve a 90% correct recognitions, and snippets of 5-7 frames (0.3-0.5 seconds of video) are enough to achieve a performance similar to the one obtainable with the entire video sequence.

### III. FINDINGS

From the survey I found that if we design the system that take the textual input and extracting web videos as well as web images, and match these videos with the images which is also simultaneously extracted from the web we will be got the accurate videos as compared to the existing system.

### IV. PROPOSED METHOD

The Method presented in this paper is built on the framework of automatic construction of an action video shot database using Web videos , which aims to extract most relevant video shots to given keywords from a large number of tagged YouTube videos in an unsupervised manner. The introduction of Web images into video shot ranking process make it more possible to obtain relevant shots in the case that tag noisy causes the failure on collecting relevant videos. In this paper, I enhance previous work by

introducing Web images into video shot

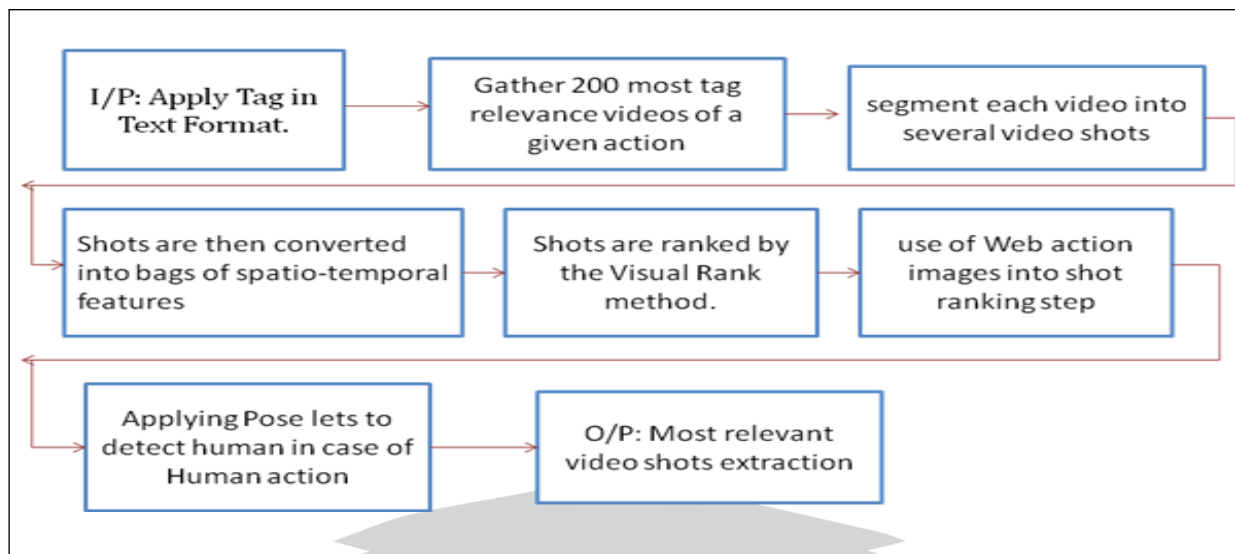


Figure 4.1: Overview of the proposed Method

I added one more steps in previous proposed method:

- 4.1: First perform tag-based videoselection; then video is converted to video shot by segmentation; do spatio-temporal feature extraction from shots; bag of features representation; similarity matrix between shots calculation,
- 4.2: Download Web images form API BING; for only human actions: selection using Poselets-based human detection [1]
- 4.3: Performing visual feature extraction from shots and images; feature matching based similarity calculation Between image and videos
- 4.4: Performing spatio-temporal-feature-based videoshot selection with shots-image similarity taken into account.

## V. CONCLUSION

In this paper we will apply a web image for automatically extraction of web video form web here will uses two database, you tube and Bing for extraction of web videos and web Images respectively, after segmentation of videos into video shots we extract the spatio-temporal feature from that videos in to separate sections. Then extracting visual feature from video shot as well as web action images, According to the matching the similarities between video shots and image we will rank the videos by visual ranking method. From this paper we conclude that will get most relevant videos of specific action image.

## REFERENCES

- [1] L. Bourdev and J. Malik. Poselets: Body part detector trained using 3d human pose annotations. In *ICCV*, 2009.
- [2] M. Blank, L. Gorelick, E. Shechtman, M. Irani, and R. Basri. Actions in space-time shapes. In *ICCV*, 2005.
- [3] N. I. Cinbins, R. G. Cinbins, and S. Sclaroff. Learning actions from the web. In *ICCV*, pp. 995–1002, 2009. [4] L. Duan, D. Xu, I. W. Tsang, and J. Luo. Visual event recognition in videos by learning from web data. In *CVPR*, 2010.
- [5] D. H. Nga and K. Yanai. Automatic construction of an action video shot database using web videos. In *ICCV*, 2011. [6] K. Schindler and L. van Gool. Action snippets: How many frames does human action recognition require? In *CVPR*, 2008.