

Audio-Visual Classification Video Browser

David Scott¹, Zhenxing Zhang¹, Rami Albatat¹, Kevin McGuinness¹,
Esra Acar², Frank Hopfgartner², Cathal Gurrin¹, Noel E. O'Connor¹,
and Alan F. Smeaton¹

¹ Dublin City University, Glasnevin, Dublin 9, Ireland

² Technische Universität Berlin, Ernst-Reuter-Platz 7, 10587 Berlin, Germany

Abstract. This paper presents our third participation in the Video Browser Showdown. Building on the experience that we gained while participating in this event, we compete in the 2014 showdown with a more advanced browsing system based on incorporating several audio-visual retrieval techniques. This paper provides a short overview of the features and functionality of our new system.

1 Introduction

In recent years, interactive audio-visual search techniques are playing a paramount role within video retrieval systems. In order to enhance the user search experience and outperform the current state-of-the-art, we embrace five multi-modal retrieval methodologies in an advanced system which does not rely on text retrieval approaches. In this paper we present our VBS Multimedia Retrieval system incorporating HTML5 interface, that communicates with a middle-ware layer utilizing the audio-visual search elements which are described in the following sections.

2 System Features

2.1 Video Segmentation

We segment the video into shots based on techniques outlined by Pickering et al.[5]. Further, we extract multiple key frames for each shot to represent the shot in the graphical user interface. A segmentation into shots allows us to provide a quick overview over different scenes within a video, a strategy which proved successful in the TRECVID known item search task. It is likely that some key frames will be very similar to other key frames of the same shot. To avoid including these similar frames, we remove duplicates by comparing the global color layout of all key frames within each shot.

2.2 Visual Concept Classification

We use trained models to classify the visual content based on concepts such as person, landscape or buildings. We use the judgements from the classification to

act as ranked list when used alone. We also use the concept lists to either filter or boost content when used in tandem with other searches. The initial classifiers were trained on TrecVid 2013 Semantic Indexing task (SIN) training set [7]. In this task 60 concepts were requested to be identified in the shots of the SIN video dataset. Three visual content descriptors were used: two types of Opponent SIFT Bag of Visual words, the first is based on dense sampling, the second on sparse sampling; the third descriptor is a concatenation of a normalized RGB Histogram and a normalized Gabor transform. A Support Vector Machine (RBF-euclidean distance kernel)[1] was trained for each of the three descriptor and each of the 60 concepts. To provide a judgement about the existence of a particular concept in a shot, the correspondent three classifiers of that concept were used to evaluate the visual features extracted from the shot, then a weighted sum of the three judgement scores generated by the classifiers is performed to provide a final score. An initial framework for feature extraction and classification parameter evaluation and optimisation is developed and tested in this task; the framework is designed to be extendible to work on a large scale data, it is installed on the machines of the Irish Centre for High-End Computing (ICHEC)¹.

2.3 Audio Concept Classification

We generate models to detect audio concepts such as explosions, gunshots and screams. In a manner similar to the visual concept classification, we use audio concept lists to either filter or boost content when used in tandem with other searches. The classifiers are trained on the MediaEval 2013 Violent Scenes Detection task (VSD) training set using high-level audio concept annotations provided in the VSD dataset. We employ Mel-Frequency Cepstral Coefficients (MFCC) features as low-level audio features. For the representation of video shots, we use a Bag-of-Audio Words (BoAW) approach based on MFCC with a sparse coding scheme. We adopt the dictionary learning technique presented in [4]. In the coding phase, we construct the sparse representations of audio signals by using the LARS algorithm [2]. In order to generate the final sparse representation of video shots which are a set of MFCC feature vectors, we apply the *max-pooling* technique. A Support Vector Machine (SVM) with an RBF kernel is trained for each audio concept using sparse audio representations. In order to provide a judgment about the existence of a particular concept in a shot, the probability estimates of SVM models are used. Normally, in a basic SVM, only class labels or scores are output. The class label results from thresholding the score, which is not a probability measure. The scores output by the SVM are converted into probability estimates using the method explained in [9].

2.4 Visual Similarity Search

The aim of visual similarity search is to offer our system the ability to find the most visual relevant shots to a given shot query and then provide the most

¹ <http://www.ichec.ie/>

possible shots for users to identify. Supposing the searching topic is happened in front of a special scene, we could easy filter out the shots which did not happened in that scene.

Different from the approach of our previous participation using locally aggregation descriptors (VLAD) and nearest neighbours searching, we followed the route of [6] and built a linear discriminative object classifier for each query image. The main benefit of this approach is that a unique weighting score will be learnt to determine the most discriminative visual features for retrieval from training data which contains one positive data and many negative data.

For each keyframe, local feature descriptors are extracted and an aggregation descriptor is generated to represent it. A large size of negative training set is created and reused for every classifier training. In the online process, the same dimensional descriptor is extracted to the search query image and a linear classifier is trained by using the open source library [3]. Finally, each video shot from dataset can be sorted by calculating the inner product of weighting vector from classifier and their feature vectors.

2.5 Face Browsing and Search

We anticipate that providing functionality to allow users to get a high-level overview of all the human faces appearing in a video will be useful for queries involving people. To this end, we provide a face view that shows all the faces found in the video, and allows users to quickly navigate to the locations in the video in which selected faces appear. We use the Viola-Jones face detector [8] to first locate faces in the videos, and then to cluster these faces by using agglomerative techniques. This clustering also allows face-based search to be easily implemented: when the user chooses to search for similar faces on a given key frame, all images associated with the clusters containing any faces that appear in the key frame can be retrieved and displayed.

3 User Interface

Our user interface features a standard multi-modal platform powered by the python based Django framework. Users will be presented with metrics for accessing each of the system features explained above and will allow for multiple users to issue collaborative search for items of interest in the collection. As an optional module to the system we will use a mobile device to capture live screenshots of the example videos to use as input to a similarity based query. Our interface will attempt to maximise the usable canvas to enable users to find items as quickly as possible.

4 Conclusion

In this paper we present a technical overview of the system which will be presented at MMM Video browser showdown 2014. In this paper we outline the

technologies used and the interaction with the system. Having come in second place overall last year, we have made some modifications in order to have a more rounded and better system which will challenge for the top spot this year.

References

1. Cortes, C., Vapnik, V.: Support-vector networks. In: Machine Learning, pp. 273–297 (1995)
2. Efron, B., Hastie, T., Johnstone, I., Tibshirani, R.: Least angle regression. *The Annals of Statistics* 32(2), 407–499 (2004)
3. Fan, R.-E., Chang, K.-W., Hsieh, C.-J., Wang, X.-R., Lin, C.-J.: LIBLINEAR: A library for large linear classification. *Journal of Machine Learning Research* 9, 1871–1874 (2008)
4. Mairal, J., Bach, F., Ponce, J., Sapiro, G.: Online learning for matrix factorization and sparse coding. *The Journal of Machine Learning Research* 11, 19–60 (2010)
5. Pickering, M.J., Rüger, S.M.: Evaluation of key frame-based retrieval techniques for video. *Computer Vision and Image Understanding* 92(2-3), 217–235 (2003)
6. Shrivastava, A., Malisiewicz, T., Gupta, A., Efros, A.A.: Data-driven visual similarity for cross-domain image matching. *ACM Transaction of Graphics (TOG) (Proceedings of ACM SIGGRAPH ASIA)* 30(6) (2011)
7. Smeaton, A.F., Over, P., Kraaij, W.: Evaluation campaigns and trecvid. In: *MIR 2006: Proceedings of the 8th ACM International Workshop on Multimedia Information Retrieval*, pp. 321–330. ACM Press, New York (2006)
8. Viola, P., Jones, M.: Rapid object detection using a boosted cascade of simple features, pp. 511–518 (2001)
9. Wu, T.-F., Lin, C.-J., Weng, R.C.: Probability estimates for multi-class classification by pairwise coupling. *The Journal of Machine Learning Research* 5, 975–1005 (2004)