

NTT-Fudan Team @ TRECVID 2015: Multimedia Event Detection

Yongqing Sun¹, Zuxuan Wu², Xi Wang², Kyoko Sudo¹, Yukinobu Taniguchi^{1,3},
Tetsuya Kinebuchi¹, and Yu-Gang Jiang²

¹ NTT Media Intelligence Laboratories, Japan

² Fudan University, China

³ Faculty of Engineering, Tokyo University of Science, Japan

1 Summary

In this notebook paper, we present an overview and results analysis of our system designed for TRECVID 2015 [1] Multimedia Event Detection (MED) task. Motivated by the great success of deep learning, we focus on exploiting various deep features to capture visual appearance and temporal dynamics in video clips. In order to fully utilize knowledge from existing large-scale image and video benchmarks, our system also incorporates high-level semantic features generated by pre-trained Convolutional Neural Networks. Then we performed classification with SVMs using different features and average the results carefully to obtain the final prediction scores. We submitted results for full evaluation of both Pre-Specified (PS) and Ad-Hoc (AH) sub-tasks in the 010Ex training condition. Our runs are submitted below.

Table 1. Summary of submitted runs for TRECVID 2015 MED

| | | |
|----|------------|---|
| AH | baseline-1 | IDT + MFCC + VGG19- f_{c_6} + VGG19- f_{c_7} + C-20K + C-233 |
| | contrast-1 | VGG19- f_{c_6} + VGG19- f_{c_7} + C-20K + C-233 |
| PS | baseline-1 | IDT + MFCC + VGG19- f_{c_6} + VGG19- f_{c_7} + C-20K + C-233 |
| | contrast-1 | IDT + MFCC + VGG19- f_{c_6} + VGG19- f_{c_7} + LSTM + C-20K + C-233 |
| | contrast-2 | IDT + MFCC + VGG19- f_{c_6} + VGG19- f_{c_7} |
| | contrast-3 | VGG19- f_{c_6} + VGG19- f_{c_7} + C-20K + C-233 |
| | contrast-4 | IDT + MFCC |

2 System Overview

In this section, we elaborate the technical components of our system. First, we describe the adopted features as well as their corresponding encoding strategies. Then we introduce the classifiers for model training and different fusion approaches.

2.1 Feature Representation

In TRECVID 2015, we adopt three sets of features in our system to capture the rich multimodal information in videos, including traditional features, deep features and concept representations. All the features used in our MED system are summarized in Table 2 and the detailed descriptions are given below.

Table 2. Features adopted in our MED system

| | Features |
|----------------------|--|
| Traditional Features | IDT (MBH, HOG, HOF) MFCC |
| Deep Features | VGG19- f_{c_6} , VGG19- f_{c_7} , LSTM |
| Concept Feature | C-20K, C-233 |

Traditional Features

– Improved Dense Trajectory (IDT): We extract the state-of-the-art improved dense trajectory features [2], which exhibit top-notch performance on action recognition tasks. Along with the densely extracted trajectories, three features are computed: HOG, HOF, and MBH. We first reduce the dimension of HOG, HOF and MBH descriptors by a factor of two using Principal Component Analysis (PCA). Then these features are further quantized respectively using the FV representation with the vocabulary size being 256.

– MFCC: In addition to the above visual features, audio features can provide complementary clues. For this, we adopt the well-known Mel-Frequency Cepstral Coefficients (MFCC). It is first computed over each 32ms time-window (with 16ms overlap) of the soundtrack and then all the descriptors are quantized into a single BoW feature representation.

Deep Features

– VGG19- f_{c_6} and VGG19- f_{c_7} : Inspired by the great success of CNN, we adopt VGG19 model proposed by Simonyan [3] in our system. Compared to AlexNet, VGG19 not only further reduces the size of convolutional filters and the stride, but more importantly, it also extends the depth of the network. With this much deeper architecture, VGG19 possess strong capabilities of learning more discriminative features and the high-level final predictions. It can produce a 7.1% top-5 error rate on the ILSVRC-2012 validation set. In order to increase the generalization ability of the VGG19 model, we finetune the model using the full ImageNet dataset, which consists of 14 million images annotated into 20K classes. Given a video clip, we extract the outputs from the two fully-connected layers (i.e., VGG19- f_{c_6} and VGG19- f_{c_7}) of each frame and then average them frame-level features into video-level representations.

– LSTM Feature: In order to further model the long-term dynamic information that is mostly discarded in the spatial CNNs, we utilize our recently developed LSTM model, as shown in Figure 1. Different from a traditional Recurrent Neural Network (RNN) unit, the LSTM unit has a built-in memory cell.

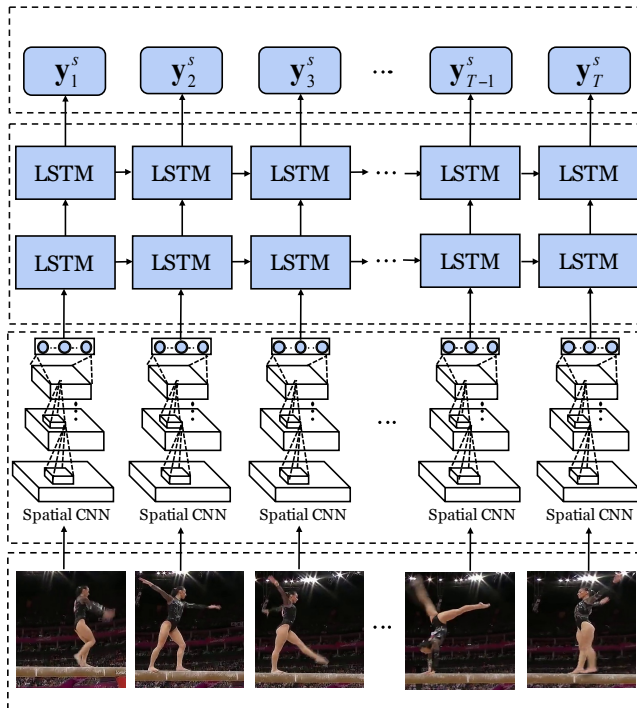


Fig. 1. The structure of the LSTM network.

Several non-linear gates are used to govern the information flow into and out of the cell, which enables the model to explore long-range dynamics by taking the feature representations from CNN at each time step. Due to time constraint of the evaluation, we directly adopt LSTM model trained with another video dataset (the UCF-101 dataset) and use the average output from all the time-steps of the last LSTM layers as the feature (512-d).

Concept Features

– C-20K: Since the softmax output of our finetuned VGG19 model demonstrates the probability of the 20K objects in a frame, we adopt this as our high-level semantic concept detector in our system. For each key frame in a given video, we obtain a 20,574-d concept score with the trained model and frame-level scores are then averaged to generate a video-level concept feature vector for further classification.

– C-233: We trained 233 concept detectors on the newly released Fudan-Columbia Video Dataset [4] using a VGG19 CNN model. Given a video clip, we obtain the 233 concept detector responses using the softmax layer of the CNN model. Then, a video level concept representation is obtained by average pooling the scores of all frames.

2.2 Classification and Fusion

To train event detection models, we employ two different types of classifiers in our system:

- Linear SVMs: To enhance classification performance, we first perform early fusion with the appearance feature and motion feature by concatenating them into a long vector. Since the concatenated vector is discriminative enough in the high-dimensional space, we adopt linear SVMs with C fixed to 100 to train the model.
- χ^2 SVMs: For MFCC audio feature, deep features and concept scores, we first map them into χ^2 -kernel separately. Then, we train independent classifiers for each of these features.

With multiple classifiers, each video clip is accordingly associated with multiple output scores, which are then fused to compute the final prediction.

3 Results and Analysis

Our MED system is designed to combine multiple feature representations to fully model multiple clues in videos. We submitted 2 runs for AH task and 5 runs for the PS task in order to investigate the effectiveness of different features.

Figure 2 shows the results of all the submissions. The official performance measure is infoAP200 for both AH and PS tasks. For AH task, we can see that traditional features are highly complementary with features extracted from deep models (i.e., deep features and concept features). For PS task, as a first trial, we incorporate LSTM features trained on UCF-101 in order to capture the long-term temporal dynamics, which promote the performance by 0.6% (PS baseline-1 vs PS-contrast-1). We claim that if the LSTM models are trained on large video corpus, the features can be more discriminative and will offer better performance. Comparing PS baseline-1 and PS contrast-2, we found that concept features can slightly improve the performance. In addition, we can see that the deep learning based features (PS contrast-3) are significantly better than the conventional features (PS contrast-4), which corroborates the fact that deep learning features trained on ImageNet usually possess high generalization ability [5].

References

1. Over, P., Awad, G., Michel, M., Fiscus, J., Kraaij, W., Smeaton, A.F., QuHenot, G., Ordelman, R.: Trecvid 2015 – an overview of the goals, tasks, data, evaluation mechanisms and metrics. In: TRECVID 2015. (2015)
2. Wang, H., Schmid, C.: Action recognition with improved trajectories. In: ICCV. (2013)
3. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. CoRR (2014)

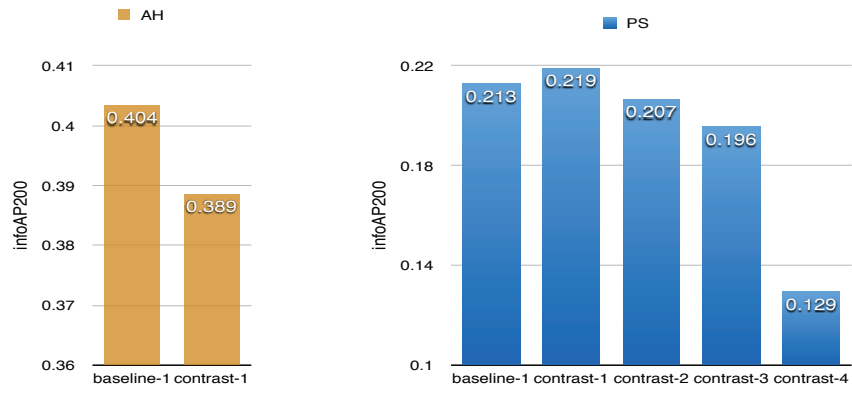


Fig. 2. The results of our submissions.

4. Jiang, Y.G., Wu, Z., Wang, J., Xue, X., Chang, S.F.: Exploiting feature and class relationships in video categorization with regularized deep neural networks. arXiv preprint arXiv:1502.07209 (2015)
5. Razavian, A.S., Azizpour, H., Sullivan, J., Carlsson, S.: CNN features off-the-shelf: an astounding baseline for recognition. CoRR (2014)

NTT-Fudan Team@TRECVID2015: Surveillance Event Detection

Haruka YONEMOTO[†] Kazuhiko MURASAKI[†] Kyoko SUDO[†]
Yukinobu TANIGUCHI^{† ‡} Tetsuya KINEBUCHI[†]

[†]NTT Media Intelligence Laboratories, NTT Corporation,
1-1 Hikari-no-oka, Yokosuka, Kanagawa, 239-0847 Japan

1 Summary

We present an event detection system for the SED task of TRECVID 2015[1]. Its simple framework is challenged with three events, CellToEar, Pointing and ObjectPut, which are the most difficult events to detect in SED. The proposed system uses STIP-HOG/HOF as the low level feature, and Fisher Vector encoding to represent each spatio-temporal sub-volume extracted by the sliding window approach. Linear SVM is used for event detection. Our results (minimum DCR metric) are compared to the results of other teams.

2 Introduction

The TRECVID SED task requires detecting observable events in real surveillance video sequences taken in London Gatwick International Airport[1]. It consists of two parts, retrospective SED and interactive SED. The former performs off-line event detection and the latter allows manually filtering of the system's result in 25 minutes. We built the system for retrospective SED. The target events in SED are PersonRuns, Pointing, CellToEar, ObjectPut, Embrace, PeopleMeet and PeopleSplitUp; Pointing and CellToEar are the most difficult events to detect because they demonstrate less change in appearance and motion. We focus on these events and, to overcome these difficulties, design a framework on the sliding window approach, where window size is adaptively changed for each event. We also employ the commonly used feature, STIP-HOG/HOF, as a low level feature. Its superiority has been demonstrated in many action recognition studies. We tested our system in three events, CellToEar, Pointing, and ObjectPut. Our system did not provide a significant improvement over past results but our evaluation results are comparable to the results of other teams in terms of minimum DCR metric.

3 Overview

We overview the proposed system in Figure 1. It consists of two main components, model learning and event detection. In model learning, the video sequence

[‡] The current affiliation is Faculty of Engineering, Tokyo University of Science, 1-3 Kagurazaka, Shinjuku-Ku, Tokyo, 162-8601, Japan

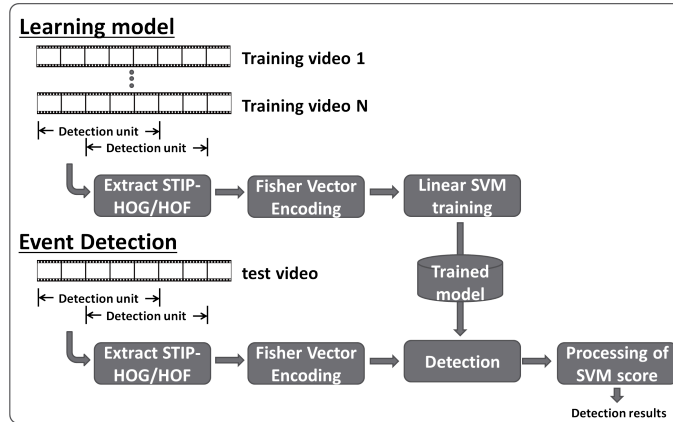


Fig. 1. Our system overflow.

is segmented into detection units by sliding windows and STIP-HOG/HOF is extracted as the low level feature. To represent the detection units, we employ the Fisher Vector encoding technique. Linear SVM is trained using these units as samples. In event detection, we read the video sequence from the first frame by sliding windows. STIP-HOG/HOF is extracted from each window and coded by Fisher Vector encoding. Linear SVM accepts the units as input and outputs a sequence of likelihood scores. The positions of events and decision scores are decided by processing this likelihood sequence.

4 Video representation

We segment the video sequence into detection units using fixed length, l , and fixed step size by using the sliding window approach. Window size l is set by calculating the average event duration and step size is set to $l/10$. The STIP detector[2] with HOG/HOF feature is commonly used to capture human motion and moving objects. In this paper, we employ it as the low level feature and employ PCA to reduce its dimension from 162 to 80. To represent each unit, the Fisher Vector encoding[3] technique is used.

5 Model learning

We use Linear SVM[4] for event detection due to its lower computation cost. To train a detector, such as Linear SVM which can predict labels with high accuracy, a large dataset with annotated labels is needed. However, such datasets are usually not available for the SED task and real-world applications. To augment what labeled samples are available, we consider a detection unit as a positive sample if more than 50% of its frames have a positive label. The negative dataset is made of samples randomly produced from the units without positive labels. We make as many negative samples as there are positive ones. The penalty score, C , of Linear SVM is decided by grid search using the training data. All elements of encoded features are normalized to the range of -1 to 1.

6 Event detection

In detecting events, we read the video sequence from the first frame by the sliding window approach. As in model learning, the STIP-HOG/HOF feature is extracted and encoded by Fisher Vector encoding. Then, a sequence of likelihood scores is generated by using encoded features as inputs to the trained Linear SVM. In post processing, we decide the positions of predicted events and calculate the decision scores from the likelihood scores. Consecutive units whose likelihoods exceed 0 are treated as one prediction event. For each prediction, the average of the likelihood scores is calculated and used as the decision score.

7 Results

Table 1. Results of event detection (mDCR-score). the left column is actual DCR. the right is actual DCR.

| | aDCR | mDCR |
|-----------|--------|--------|
| CellToEar | 1.5153 | 1.0006 |
| Pointing | 3.0253 | 1.0006 |
| ObjectPut | 4.7264 | 0.9965 |

In training the model, we used all data and annotations provided by NIST. Window size, l , was set at 16, 20, and 10 when making the positive dataset and the penalty score, C , in Linear SVM was set at $2^{12.2}$, 2^{-5} , 2^{15} , respectively, for CellToEar, Pointing, and ObjectPut. Table 1 shows our results for the EVAL15 data. We conjecture that the proposed system offered only low precision because it uses Fisher Vector encoding which eliminates the spatial information of the active subjects. Localizing the active subjects might improve the prediction score of these two events.

References

1. Paul Over, George Awad, Martial Michel, Jonathan Fiscus, Wessel Kraaij, Alan F. Smeaton, Georges Quénot, and Roeland Ordeman. Trecvid 2015 – an overview of the goals, tasks, data, evaluation mechanisms and metrics. In *Proceedings of TRECVID 2015*. NIST, USA, 2015.
2. Ivan Laptev. On space-time interest points. *International Journal of Computer Vision*, Vol. 64, No. 2-3, pp. 107–123, 2005.
3. Florent Perronnin, Jorge Sánchez, and Thomas Mensink. Improving the fisher kernel for large-scale image classification. In *Computer Vision–ECCV 2010*, pp. 143–156. Springer, 2010.
4. Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. Liblinear: A library for large linear classification. *The Journal of Machine Learning Research*, Vol. 9, pp. 1871–1874, 2008.