

Evaluating multimedia features and fusion for example-based event detection

Gregory K. Myers · Ramesh Nallapati · Julien van Hout · Stephanie Pancoast · Ramakant Nevatia · Chen Sun · Amirhossein Habibian · Dennis C. Koelma · Koen E. A. van de Sande · Arnold W. M. Smeulders · Cees G. M. Snoek

Received: 23 January 2013 / Revised: 28 May 2013 / Accepted: 4 June 2013
© The Author(s) 2013. This article is published with open access at Springerlink.com

Abstract Multimedia event detection (MED) is a challenging problem because of the heterogeneous content and variable quality found in large collections of Internet videos. To study the value of multimedia features and fusion for representing and learning events from a set of example video clips, we created SESAME, a system for video SEArch with Speed and Accuracy for Multimedia Events. SESAME includes multiple bag-of-words event classifiers based on single data types: low-level visual, motion, and audio features; high-level semantic visual concepts; and automatic speech recognition. Event detection performance was evaluated for each event classifier. The performance of low-level visual and motion features was improved by the use of difference coding. The accuracy of the visual concepts was nearly as strong as that of the low-level visual features. Experiments with a number of fusion methods for combining the event detection scores from these classifiers revealed that simple fusion methods, such as arithmetic mean, perform as well as or better than other, more complex fusion methods. SESAME's

performance in the 2012 TRECVID MED evaluation was one of the best reported.

Keywords Multimedia event detection · Video retrieval · Content extraction · Difference coding · Late fusion

1 Introduction

The goal of multimedia event detection (MED) is to detect user-defined events of interest in massive, continuously growing video collections, such as those found on the Internet. This is an extremely challenging problem because the contents of the videos in these collections are completely unconstrained, and the collections include varying qualities of user-generated videos, often made with handheld cameras, and may have jerky motions, wildly varying fields of view, and poor lighting. The audio in these videos is recorded in a variety of acoustic environments, often with a single camera-mounted microphone, with no attempt to prevent background sounds from masking speech.

For purposes of this research, an event, as defined in the TRECVID MED evaluation task sponsored by the National Institute of Standards and Technology (NIST) [1], has the following characteristics:

- It includes a complex activity occurring at a specific place and time.
- It involves people interacting with other people and/or objects.
- It consists of a number of human actions, processes, and activities that are loosely or tightly organized and have significant temporal and semantic relationships to the overarching activity.
- It is directly observable.

G. K. Myers (✉) · R. Nallapati · J. van Hout · S. Pancoast
SRI International (SRI), 333 Ravenswood Avenue,
Menlo Park, CA 94025, USA
e-mail: gregory.myers@sri.com

R. Nevatia · C. Sun
Institute for Robotics and Intelligent Systems, University of
Southern California (USC), Los Angeles, CA 90089-0273, USA

A. Habibian · D. C. Koelma · K. E. A. van de Sande ·
A. W. M. Smeulders · C. G. M. Snoek
University of Amsterdam (UvA), Science Park 904, P.O. Box
94323, Amsterdam 1098 GH, The Netherlands

R. Nallapati
IBM Thomas J Watson Research Center, 1101 Kitchawan Rd,
Yorktown Heights, NY 10598, USA

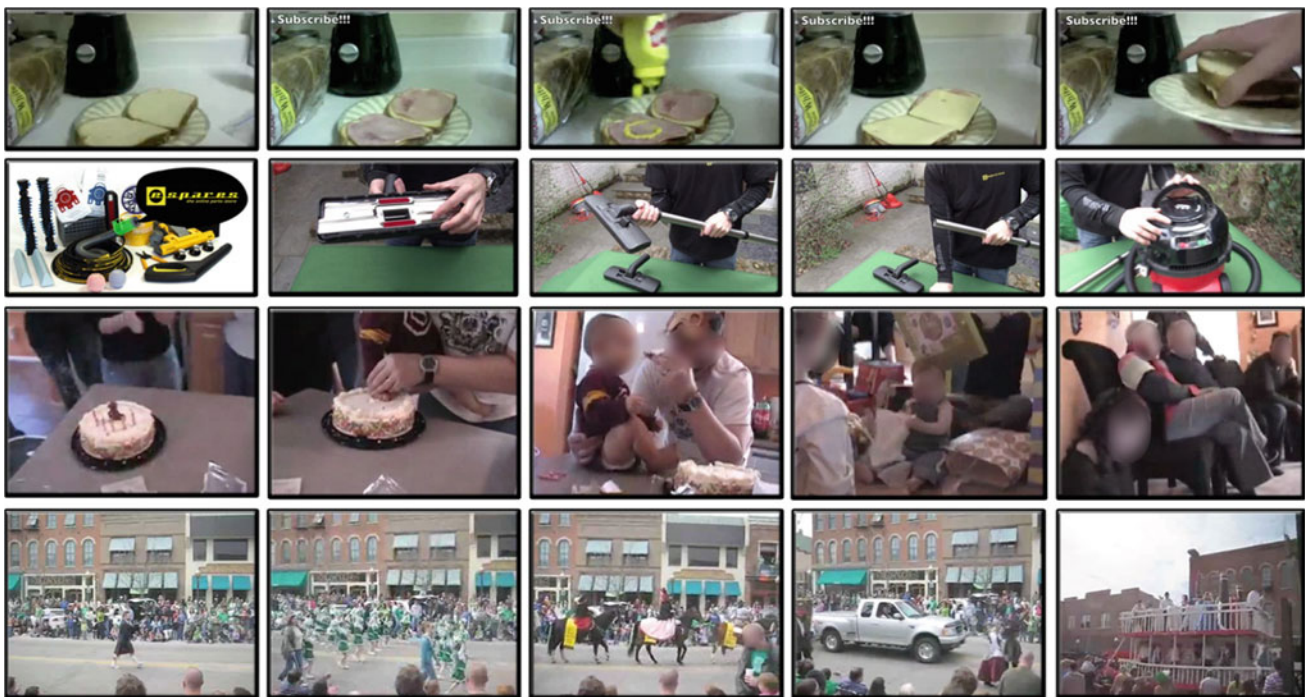


Fig. 1 Key-frame series from example videos for the events making a sandwich, repairing an appliance, birthday party, and parade

Figure 1 shows some sample video imagery from events in the TRECVID MED evaluation task. Events are more complex and may include actions (hammering, pouring liquid) and activities (dancing) occurring in different scenes (street, kitchen). Some events may be process-oriented, with an expected sequence of stages, actions, or activities (making a sandwich or repairing an appliance); other events may be a set of ongoing activities with no particular beginning or end (birthday party or parade). An event may be observed in only a portion of the video clip, and relevant clips may contain extraneous content.

Multimedia event detection can be considered as a search problem with a query-retrieval paradigm. Currently, videos in online collections, such as YouTube, are retrieved based on text-based search. Text labels are either manually assigned when the video is added to the collection or derived from text already associated with the video, such as text content that occurs near the video in a multimedia blog or web page. Videos are searched and retrieved by matching a text-based user query to videos' text labels, but performance will depend on the quality and availability of such labels.

Highly accurate text-based video retrieval requires the text-based queries to be comprehensive and specific. In the TRECVID MED evaluation, each event is defined by an “event kit,” which includes a 150–400 word text description consisting of an event name, definition, explication (textual exposition of the terms and concepts), and lists of scenes, objects, people, activities, and sounds that would indicate the presence of the event. Figure 2 shows an exam-

ple for the event working on a woodworking project. The user might also have to specify how similar events are distinguished from the event of interest (e.g., not construction in Fig. 2), and may have to estimate the frequency with which various entities occur in the event (e.g., often indoors). Subcategories and variations of the event may also have to be considered (e.g., operating a lathe in a factory).

Another approach to detect events is to define the event in terms of a set of example videos, which we call an example-based approach. Example videos are matched to videos in the collection using the same internal representation for each. In this approach, the system automatically learns a model of the event based on a set of positive and negative examples, taking advantage of well-established capabilities in machine learning and computer vision. This paper considers an example-based approach with both non-semantic and semantic representations.

Current approaches for MED [2–7] rely heavily on kernel-based classifier methods that use low-level features computed directly from the multimedia data. These classifiers learn a mapping between the computed features and the category of event that occurs in the video. Videos and events are typically represented as “bag-of-words” (BOW) models composed of histograms of descriptors for each feature type, including visual, motion, and audio features. Although the performance of these models is quite effective, individual low-level features do not correspond directly to terms with semantic meaning, and therefore cannot provide

<p>Event name: <i>Working on a woodworking project</i></p> <p>Definition: <i>One or more people fashion an object out of wood</i></p> <p>Explication: <i>Woodworking is a popular hobby that involves crafting an object out of wood. Typical woodworking projects may range from creating large pieces of furniture to small decorative items or toys. The process for making objects out of wood can include cutting wood into smaller pieces with hand or machine tools, carving wood to shape it, sanding wood to smooth it, gluing wood pieces together, drilling holes into wood, and applying decorative finishes to the completed object. Woodworking is distinguished from construction, which typically involves creation of large, permanent structures such as houses, sheds, or buildings, which may or may not be made of wood.</i></p> <p>Evidential description:</p> <p>scene: <i>often indoors in a workshop, garage, artificial lighting, occasionally outdoors</i></p> <p>objects/people: <i>woodworking tools (automatic or non-automatic saws, sander, knife), paint, stains, sawhorses, toolbox, safety goggles</i></p> <p>activities: <i>cutting and shaping wood, attaching pieces of wood together, smoothing/sanding wood</i></p> <p>audio: <i>sounds from power tools, hand tools being used (hammer, saw, etc.); narration of the process</i></p>

Fig. 2 Event Kit for working on a woodworking project

human-understandable evidence of why a video was selected by the MED system as a positive instance of a specific event.

A second representation is in terms of higher-level semantic concepts, which are automatically detected in the video content [8–11]. The detectors are related to objects, like a flag; scenes, like a beach; people, like female; and actions, like dancing. The presence of concepts such as these creates an understanding of the content. However, except for a few entities such as faces, most individual concept detectors are not yet reliable [12]. In addition, training detectors for each concept requires annotated data, which usually involves significant manual effort to generate. In the future, it is expected that more annotated datasets will be available, and weakly supervised learning methods will help improve the efficiency of generating them. Event representations based on high-level concepts have started to appear in the literature [13–16].

For an example-based approach, the central research issue is to find an event representation in terms of the elements of the video that permits the accurate detection of the events. In our approach, an event is modeled as a set of multiple bags-of-words, each based on a single data type. Partitioning the representation by data type permits the descriptors for each data type to be optimized independently (specific multimodal combinations of features, such as bimodal audiovisual features [3], can be considered a single data type within this architecture). The data types we used included both low-level features (visual appearance, motion, and audio) and higher-level semantic concepts (visual concepts). We also used auto-

matic speech recognition (ASR) to generate a BOW model in which semantic concepts were expressed directly by words in the recognized speech. The resulting event model combined multiple sources of information from multiple data types and multiple levels of information.

As part of the optimization process for the low-level features, we investigated the use of difference coding techniques in addition to conventional coding methods. Because the information captured by difference coding is somewhat complementary to the information produced by the traditional BOW, we anticipated an improvement in performance. We conducted experiments to compare the performance of difference coding techniques with conventional feature coding techniques.

The remaining challenge is finding the best method for combining the multiple bags-of-words in the event-detection decision process. The most common approach is to apply late fusion methods [3, 5, 17] in which the results for each data type are combined by fusing the decision scores from multiple event classifiers. This is a straightforward way of using the information from all data types in proportion to their relative contribution to event detection on videos with widely diverse content. We evaluated the performance of several fusion methods.

The work described in this paper focused on evaluating the various data types and fusion methods for MED. Our approach for example-based MED, including methods for content extraction and fusion, is described in Sect. 2.

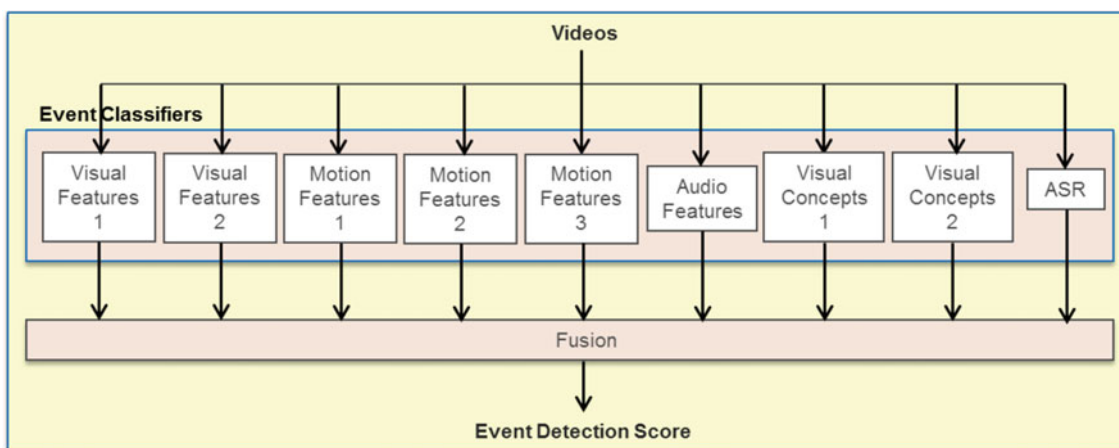


Fig. 3 Major components of the SESAME system

Experimental results are described in Sect. 3, and Sect. 4 contains a summary and discussion.

All the experiments for evaluating the performance of the MED capability were performed using the data provided in the TRECVID MED evaluation task. The MED evaluation uses the Heterogeneous Audio Visual Internet Collection (HAVIC) video data collection [18], which is a large corpus of Internet multimedia files collected by the Linguistic Data Consortium.

2 Approach for example-based MED

The work in this paper focuses on SEArch with Speed and Accuracy for Multimedia Events (SESAME), an MED system in which an event is specified as a set of video clip examples. A supervised learning process trains an event model from positive and negative examples, and an event classifier uses the event model to detect the targeted event. An event classifier was built for each data type. The results of all the event classifiers were then combined by fusing their decision scores. An overview of the SESAME system and methods for event classification and fusion are described in the following sections.

2.1 SESAME system overview

The major components of the SESAME system are shown in Fig. 3. A total of nine event classifiers generate event detection decision scores: two based on low-level visual features, three based on low-level motion features, one based on low-level audio features, two based on visual concepts, and one based on ASR. The outputs of the event classifiers are combined by the fusion process.

Figure 4 shows the processing blocks within each event classifier. Each event classifier operates on a single type of data and includes both training and event classification. Con-

tent is extracted from positive and negative video examples, and the event classifier is trained, resulting in an event model. The event model produces event detection scores when it is applied to a test set of videos. Figure 4 does not show off-line training and testing to optimize the parameter settings for the content extraction processes.

2.2 Content extraction methods

This section describes the feature coding and aggregation methods that were common to the low-level features and the content extraction methods for the different data types: low-level visual features, low-level motion features, low-level audio features, high-level visual features, and ASR.

2.2.1 Feature coding and aggregation

The coding and aggregation of low-level features share common elements that we describe here. We extracted local features and aggregated them by using three approaches: conventional BOW, vector of locally aggregated descriptors (VLAD), and Fisher vectors (FV).

The conventional BOW approach partitions low-level features into k -means clusters to generate a codebook. Given a set of features from a video, a histogram is generated by assigning each feature from the set to one or several nearest code words. Several modifications to this approach are possible. One variation uses soft coding, where instead of assigning each feature to a single code word, distances from the code words are used to weigh the histogram terms for the code words. Another variation describes code words by a Gaussian mixture model (GMM), rather than just by the center of a cluster.

While conventional BOW aggregation has been successfully used for many applications, it does not maintain any information about the distribution of features in the fea-

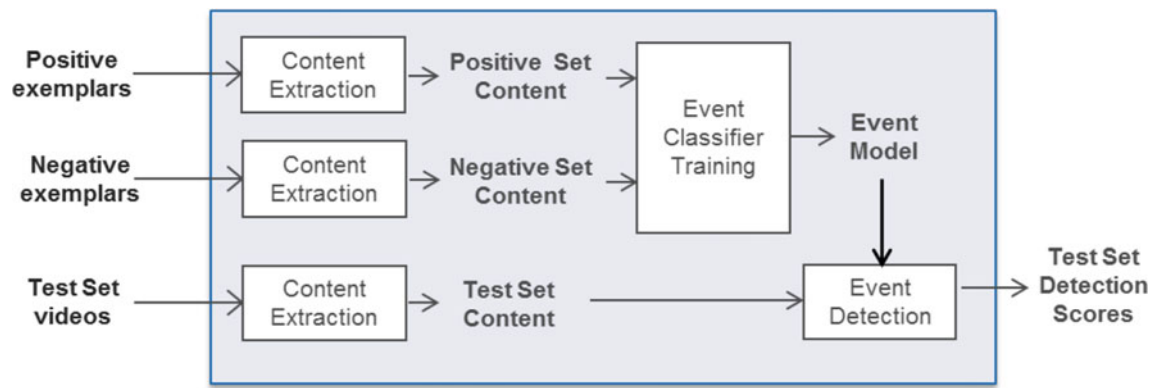


Fig. 4 Example-based event classifier for MED

ture space. FV has been introduced in previous work [19] to capture more detailed statistics, and has been applied to image classification and retrieval [20,21]. The basic idea is to represent a set of data by a gradient of its log-likelihood to model parameters and to measure the distance between instances with the Fisher kernel. For local features extracted from videos, it becomes natural to model their distribution as GMMs, forming a soft codebook. With GMM, the dimension of FV is linear in the number of mixtures and local feature dimensions.

Finally, VLAD [20] is proposed as a non-probabilistic version of FV. It uses k -means instead of GMM, and accumulates the relative positions of feature points to their single nearest neighbors in the codebook.

Compared with conventional BOW, FV and VLAD have the following benefits:

- FV takes GMM as the underlying generative model.
- Both FV and VLAD are derivatives, so feature points with the same distribution as the general model have no overall impact on the video-level descriptors; as a result, FV and VLAD can suppress noisy and redundant signals.

None of the above aggregation methods consider feature localization in space or in time. We introduced a limited amount of this information by dividing the video into temporal segments (for time localization) and spatial pyramids (for spatial localization). We then compute the features in each segment or block separately and concatenate the resulting features. The spatial pooling and temporal segmentation parameters that yielded the best performance were determined through experimentation.

2.2.2 Visual features

Two event classifiers were developed based on low-level visual features [22]. They both follow a pipeline consisting

of four stages: spatiotemporal sampling of points of interest, visual description of those points, encoding the descriptors into visual words, and supervised learning with kernel machines.

Spatiotemporal sampling The visual appearance of an event in video may have a dependency on the spatiotemporal viewpoint under which it is recorded. Salient point methods [23] introduce robustness against viewpoint changes by selecting points, which can be recovered under different perspectives. To determine salient points, Harris–Laplace relies on a Harris corner detector; applying it on multiple scales makes it possible to select the characteristic scale of a local corner using the Laplacian operator. For each corner, the Harris–Laplace detector selects a scale invariant point if the local image structure under a Laplacian operator has a stable maximum.

Another solution is to use many points by dense sampling. For imagery with many homogenous areas, such as outdoor snow scenes, corners may be rare, therefore relying on a Harris–Laplace detector can be suboptimal. To counter the shortcomings of Harris–Laplace, we used dense sampling, which samples an image grid in a uniform fashion, using a fixed pixel interval between regions.

In our experiments, we used an interval distance of six pixels and sampled at multiple scales. Appearance variations caused by temporal effects were addressed by analyzing video beyond the key-frame level [24]. Taking more frames into account during analysis allowed us to recognize events that were visible during the video, but not necessarily in a single key frame. We sampled one frame every 2 s. Both Harris–Laplace and dense sampling give an equal weight to all key-points, regardless of their spatial location in the image frame. To overcome this limitation, Lazebnik et al. [25] suggested repeated sampling of fixed subregions of an image, e.g., 1×1 , 2×2 , 4×4 , etc., and then aggregating the different resolutions into a spatial pyramid, which allows for region-specific weighting. Since every region is an image in itself, the spatial pyramid can be combined with both the Harris–

Laplace point detector and dense point sampling. We used a spatial pyramid of 1×1 and 1×3 regions in our experiments.

Visual descriptors In addition to the visual appearance of events in the spatiotemporal viewpoint under which they are recorded, the lighting conditions during recording also play an important role in MED. Properties of color features under classes of illumination and viewing features, such as viewpoint, light intensity, light direction, and light color, can change, specifically for real-world datasets as considered within TRECVID [26]. We followed [22] and used a mixture of SIFT, OpponentSIFT, and C-SIFT descriptors. The SIFT feature proposed by Lowe [27] describes the local shape of a region using edge-orientation histograms. Because the SIFT feature is normalized, the gradient magnitude changes have no effect on the final feature. OpponentSIFT describes all the channels in the opponent color space using SIFT features. The information in the O3 channel is equal to the intensity information, while the other channels describe the color information in the image. The feature normalization, as effective in SIFT, cancels out any local changes in light intensity. In the opponent color space, the O1 and O2 channels still contain some intensity information. To add invariance to shadow and shading effects, the C-invariant [28] eliminates the remaining intensity information from these channels. The C-SIFT feature uses the C-invariant, which can be seen as the gradient (or derivative) for the normalized opponent color space O1/I and O2/I. The I intensity channel remains unchanged. C-SIFT is known to be scale-invariant with respect to light intensity. We computed the SIFT and C-SIFT descriptors around salient points obtained from the Harris–Laplace detector and dense sampling. We then reduced all descriptors to 80 dimensions with principal component analysis (PCA).

Word encoding To avoid using all low-level visual features from a video, we followed the well-known codebook approach. We first assigned the features to discrete codewords from a predefined codebook. Then, we used the frequency distribution of the codewords as a compact feature vector representing an image frame. Based on [22], we employed codebook construction using k -means clustering in combination with average codeword assignment and a maximum of 4,096 codewords. The traditional hard assignment can be improved using soft assignment through kernel codebooks [29]. A kernel codebook uses a kernel function to smooth the hard assignment of (image) features to codewords by assigning descriptors to multiple clusters weighted by their distance to the center. We also used difference coding, with VLAD performing k -means clustering of the PCA-reduced descriptor space with 1,024 components. The output of the word encoding is a BOW vector using either hard average coding or soft VLAD coding. The BOW vector forms the foundation for event detection.

Kernel learning Kernel-based learning methods are typically used to develop robust event detectors from audiovisual features. As described in [22], we relied predominantly on the support vector machine (SVM) framework for supervised learning of events: specifically, the LIBSVM¹ implementation with probabilistic output. To handle imbalance in the number of positive versus negative training examples, we fixed the weights of the positive and negative classes by estimating the prior probabilities of the classes on training data. We used the histogram intersection kernel and its efficient approximation as suggested by Maji et al. [30]. For difference coded BOWs, we used a linear kernel [19].

Experiments We evaluated the performance of these two event classifiers on a set of 12,862 drawn from the training and development data from the TRECVID MED evaluation. This SESAME Evaluation dataset consisted of a training set of 8,428 videos and a test set of 4,434 videos sampled from 20 event classes and other classes that did not belong to any of the 20 events. To make good use of the limited number of available positive instances of events, the positives were distributed so that, for each event, there were approximately twice as many positives in the training set as there were in the test set. Separate classifiers were trained for each event based on a one-versus-all paradigm. Table 1 shows the performance of the two event classifiers measured by mean average precision (MAP). Color-average coding with a histogram intersection kernel (HIK) SVM slightly outperformed color-difference soft coding with a linear SVM. For events, such as changing a vehicle tire and town hall meeting, the average HIK was the best event representation. However, for some events, such as flash mob gathering and dog show, the difference coding was more effective. To study whether the representations complement each other, we also performed a simple average fusion; the results indicate a further increase in event detection performance, improving mean average precision from 0.342 to 0.358 and giving the best overall performance for the majority of events.

2.2.3 Motion features

Many motion features for activity recognition have been suggested in previous work; [4] provides a nice evaluation of motion features for classifying web videos on the NIST MED 2011 dataset. Based on our analysis of previous work and some small-scale experiments, we decided to use three features: spatiotemporal interest points (STIPs) and dense trajectories (DTs) [31], and MoSIFT [32]. STIP features are computed at corner-like locations in the 3D spatiotemporal volume. Descriptors consist of histograms of gradient and optical flow at these points. This is a very commonly used

¹ <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>.

Table 1 Mean average precision (MAP) of event classifiers with low-level visual features and their fusion for 20 TRECVID MED evaluation event classes

Event	Average coding with HIK SVM	Difference coding with linear SVM	Fusion
Birthday_party	0.275	0.229	0.261
Changing_a_vehicle_tire	0.305	0.269	0.302
Flash_mob_gathering	0.602	0.644	0.636
Getting_a_vehicle_unstuck	0.457	0.496	0.494
Grooming_an_animal	0.280	0.222	0.275
Making_a_sandwich	0.268	0.278	0.314
Parade	0.416	0.415	0.427
Parkour	0.464	0.413	0.450
Repairing_an_appliance	0.486	0.469	0.498
Working_on_a_sewing_project	0.378	0.388	0.400
Attempting_a_bike_trick	0.398	0.350	0.408
Cleaning_an_appliance	0.138	0.077	0.135
Dog_show	0.595	0.651	0.636
Giving_directions_to_a_location	0.123	0.130	0.134
Marriage_proposal	0.058	0.093	0.071
Renovating_a_home	0.229	0.273	0.285
Rock_climbing	0.488	0.466	0.507
Town_hall_meeting	0.531	0.463	0.502
Winning_a_race_without_a_vehicle	0.237	0.284	0.263
Working_on_a_metal_crafts_project	0.109	0.134	0.153
Mean for all events	0.342	0.337	0.358

Best result per event is denoted in bold

descriptor; more details may be found in [33]. Dense trajectory features are computed on a dense set of local trajectories (typically computed over 15 frames). Each trajectory is described by its shape and by histograms of intensity gradient, optical flow, and motion boundaries around it. Motion boundary features are somewhat invariant to camera motion. MoSIFT, as its name suggests, uses SIFT feature descriptors; its feature detector is built on motion saliency. STIP and DT were extracted using the default parameters as provided²; the MoSIFT features were obtained in the form of coded BOW features.³

After the extraction of low-level motion features, we generated a fixed-length video-level descriptor for each video. We experimented with the coding schemes described in Sect. 2.2.1 for the STIP and DT features; for MoSIFT, we were able to use BOW features only. We used the training and test sets described above.

We trained separate SVM classifiers for each event and each feature type. Training was based on a one-versus-all paradigm. For conventional BOW features, we used the χ^2 kernel. We used the Gaussian kernel for VLAD and FV. To

select classifier-independent parameters (such as the codebook size), we conducted fivefold cross validation of 2,062 videos from 15 event classes. We conducted fivefold cross validation on the training set to select classifier-dependent parameters. For BOW features, we used 1,000 codewords; for FV and VLAD, we used 64 cluster centers. More details of the procedure are found in [34].

We compared the performance of conventional BOW, FV, and VLAD for STIP features; BOW and FV for DT features; and BOW for MoSIFT, using the SESAME Evaluation dataset. Table 2 shows the results.

We can see that FV gave the best MAP for both STIP and DT. VLAD also improved MAP for STIP, but was not as effective as the FV features. We were not able to perform VLAD and FV experiments for MoSIFT features, but would expect to have seen similar improvements there.

2.2.4 Audio features

The audio is modeled as a bag-of-audio-words (BOAW). The BOAW has recently been used for audio document retrieval [35] and copy detection [36], as well as MED tasks [37]. Our recent work [38] describes the basic BOAW approach. We extracted the audio data from the video files and converted them to a 16 kHz sampling rate. We extracted Mel frequency cepstral coefficients (MFCCs) for every 10 ms interval using

² We obtained the STIP code from <http://www.di.ens.fr/~laptev/download/stip-1.1-winlinux.zip>, and DT code from http://lear.inrialpes.fr/people/wang/dense_trajectories.

³ MoSIFT features were provided by Dr. Alex Hauptmann of Carnegie-Mellon University.

Table 2 Mean average precision of event classifiers with motion features for 20 TRECVID MED evaluation event classes

Event	BOW + MoSIFT	BOW + STIP	VLAD + STIP	FV + STIP	BOW + DT	FV + DT
Birthday_party	0.191	0.217	0.217	0.189	0.225	0.293
Changing_a_vehicle_tire	0.126	0.064	0.165	0.136	0.190	0.217
Flash_mob_gathering	0.463	0.535	0.579	0.569	0.564	0.567
Getting_a_vehicle_unstuck	0.337	0.284	0.316	0.365	0.403	0.439
Grooming_an_animal	0.290	0.093	0.116	0.147	0.216	0.247
Making_a_sandwich	0.164	0.154	0.193	0.225	0.198	0.234
Parade	0.326	0.260	0.364	0.457	0.446	0.419
Parkour	0.295	0.366	0.404	0.369	0.413	0.459
Repairing_an_appliance	0.368	0.357	0.370	0.385	0.417	0.443
Working_on_a_sewing_project	0.270	0.292	0.346	0.386	0.352	0.433
Attempting_a_bike_trick	0.640	0.104	0.234	0.235	0.245	0.438
Cleaning_an_appliance	0.090	0.058	0.088	0.074	0.066	0.089
Dog_show	0.488	0.361	0.489	0.557	0.600	0.632
Giving_directions_to_a_location	0.085	0.194	0.148	0.191	0.069	0.052
Marriage_proposal	0.027	0.040	0.107	0.173	0.059	0.118
Renovating_a_home	0.157	0.182	0.201	0.255	0.277	0.361
Rock_climbing	0.465	0.156	0.326	0.352	0.470	0.425
Town_hall_meeting	0.519	0.285	0.286	0.462	0.317	0.370
Winning_a_race_without_a_vehicle	0.273	0.187	0.174	0.260	0.179	0.216
Working_on_a_metal_crafts_project	0.116	0.148	0.064	0.032	0.072	0.128
MAP	0.285	0.217	0.259	0.291	0.289	0.329

Best result per event is denoted in bold

a hamming window with 50 % overlap. The features consist of 13 values (12 coefficients and the log-energy), along with their delta and delta-delta values. We used a randomized sample of the videos from the TRECVID 2011 MED evaluation development set to generate the codebook. We performed k -means clustering on the MFCC features to generate 1,000 clusters. The centroid for each cluster is taken as a code word. The soft quantization process used the codebook to map the MFCCs to code words. We trained an SVM classifier with a histogram intersection kernel on the soft quantization histogram vectors of the video examples, and used the classifier to detect the events. Evaluation with the SESAME Evaluation dataset showed that the audio features achieved a MAP of 0.112.

2.2.5 Visual concepts

Two event classifiers were based on concept detectors. We followed the pipeline proposed in [39]. We decoded the videos by uniformly extracting one frame every 2 s. We then applied all available concept detectors to the extracted frames. After we concatenated the detector outputs, each frame was represented by a concept vector. Finally, we aggregated the frame representations into a video-level representation by averaging and normalization. On top of this concept

representation per video, we used either a HIK SVM or a random forest as an event classifier.

To create the concept representation, we needed a comprehensive pool of concept detectors. We built this pool of detectors using the human-annotated training data from two publicly available resources: the TRECVID 2012 Semantic Indexing task [40] and the ImageNet Large-Scale Visual Recognition Challenge 2011 [41]. The former has annotations for 346 semantic concepts on 400,000 key frames from web videos. The latter has annotations for 1,000 semantic concepts on 1,300,000 photos. The categories are quite diverse and include concepts from various types; i.e., objects like helicopter and harmonica, scenes like kitchen and hospital, and actions like greeting and swimming. Leveraging the annotated data available in these datasets, we trained 1,346 concept detectors in total.

We followed the state-of-the-art for our implementation of the concept detectors. We used densely sampled SIFT, OpponentSIFT, and C-SIFT descriptors, as we had for our event detector using visual features, but this time, we used difference coding with FV [19]. We used a visual vocabulary of 256 words. We again used the full image and three horizontal bars as a spatial pyramid. The feature vectors representing the training images formed the input for a linear SVM.

Experiments with the SESAME Evaluation dataset, summarized in Table 3, show that the random forest classifier

Table 3 Mean average precision of event classifiers with visual concept features for 20 TRECVID MED evaluation event classes

Event	RF	SVM
Birthday_party	0.339	0.324
Changing_a_vehicle_tire	0.251	0.241
Flash_mob_gathering	0.542	0.542
Getting_a_vehicle_unstuck	0.454	0.426
Grooming_an_animal	0.254	0.231
Making_a_sandwich	0.283	0.257
Parade	0.373	0.306
Parkour	0.550	0.479
Repairing_an_appliance	0.422	0.404
Working_on_a_sewing_project	0.390	0.394
Attempting_a_bike_trick	0.475	0.472
Cleaning_an_appliance	0.097	0.149
Dog_show	0.595	0.529
Giving_directions_to_a_location	0.058	0.097
Marriage_proposal	0.077	0.066
Renovating_a_home	0.295	0.325
Rock_climbing	0.412	0.401
Town_hall_meeting	0.411	0.417
Winning_a_race_without_a_vehicle	0.198	0.167
Working_on_a_metal_crafts_project	0.099	0.162
Mean for all events	0.341	0.330

Best result per event is denoted in bold

is more successful than the non-linear HIK SVM for event detection using visual concepts, although the two approaches are quite close on average. Note that the event detection results using visual concepts are close to our low-level representation using visual or motion features.

2.2.6 Automatic speech recognition

Spoken language content is often present in user-generated videos and can potentially contribute useful information for detecting events. The recognized speech has direct semantic information that typically complements the information contributed by low-level visual features. We used DECIPHER, SRI's ASR software, to recognize spoken English. We used acoustic and language models obtained from an ASR system [42] trained on speech data recorded in meetings with a far-field microphone. Initial tests on the audio in user-generated videos revealed that the segmentation process, which distinguishes speech from other audio, often misclassified music as speech. Therefore, before running the speech recognizer on these videos, we constructed a new segmenter, which is described below.

The existing segmenter was GMM based and had two classes (speech and non-speech). For this effort, we leveraged the availability of annotated TRECVID video data and

built a segmenter better tuned to audio conditions in user-generated videos. We built a segmenter with four classes: speech, music, noise, and pause. We measured the effectiveness of the new segmentation by the word-error rates (WERs) obtained by feeding the speech-segmented audio to our ASR system. We found that the new segmentation helped reduce the WER from 105 to 83 %. This confirmed that the new segmentation models were a better match to the TRECVID data than models trained on meeting data. For reference, when all the speech segments were processed by the ASR, the WER obtained by our system was 78 % (this oracle segmentation provided the lowest WER that could be achieved by improving the segmentation).

To create features for the event classifiers, we used ASR recognition lattices to compute the expected word counts for each word and each video. This approach provided significantly better results compared to using the 1-best ASR output, because it compensated for ASR errors by including words with lower posteriors that were not necessarily present in the 1-best. We computed the logarithm of the counts for each word, appended them to form a feature vector of dimension 34,457, and used a linear SVM for the event classifiers. More details may be found in [43]. Evaluation with the SESAME Evaluation dataset showed that the ASR event classifiers achieved a MAP of 0.114.

2.3 Fusion

We implemented a number of fusion methods, all of which involved a weighted average of the detection scores from the individual event classifiers. The methods for determining the weights considered several factors:

- *Event dependence and learned weights* Because the set of most reliable data types for different events might vary, we considered the importance of learning the fusion weights for each event using a training set. However, when there is limited data available for training, aggregating the data for all events and computing a fixed set of weights for all events may yield more reliable results. Another strategy is to set the weights without training with any data at all. For example, in the method of fusing with the arithmetic mean of the scores, all the weights are equal.
- *Score dependence* For weights learned via cross-validation on a training set, a single set of fixed weights might be learned for the entire range of detection scores. Alternatively, the multidimensional space of detection scores might be partitioned into a set of regions, with a set of weights assigned to each region. In general, more data is needed for score-dependent weights to avoid overfitting.
- *Adjustment for missing scores* When the scores for some types of data (particularly for ASR and MFCC) are miss-

ing, a default value, such as an average for the missing score, might be used, but this could provide a misleading indication of contribution. Other ways of dealing with missing scores include renormalizing the weights of the non-missing scores, or learning multiple sets of weights, each set for a particular combination of non-missing scores.

We evaluated the fusion models described below. All the models operated on detection scores were normalized using a Gaussian function (i.e., computing the z score by removing the mean and scaling by the standard deviation)

Arithmetic mean (AM) In this method, we compute the AM of the scores of the observed data types for a given clip. Missing data types for a given clip are ignored, and the averaging is performed over the scores of observed data types.

Geometric mean (GM) In this method, we compute the uniform GM of the scores of the observed data types for a given clip. As we do for AM, we ignore missing data types and compute the geometric mean of the scores from observed data types.

Mean average precision-weighted fusion (MAP) This fusion method weighs scores from the observed data types for a clip by their normalized average precision scores, as computed on the training fold. Again, the normalization is performed only over the observed data types for a given clip.

Weighted mean root (WMR) This fusion method is a variant of the MAP-weighted method. In this method, we compute the fusion score as we do for MAP-weighted fusion, except the final fused score x' is determined by performing a power normalization of the MAP-based fused score x :

$$x' = x^{\frac{1}{\alpha}} \quad (1)$$

where α is the number of non-missing data types for that video. In other words, the higher the number of data types from which the fusion score is computed, the more trustworthy the output.

Conditional mixture model This model combines the detection scores from various data types using mixture weights that were trained by the expectation maximization (EM) algorithm on the labeled training folds. For clips that are missing scores from one or more data types, we provide the expected score for that data type based on the training data.

Sparse mixture model (SMM) This extension of the conditional mixture model addresses the problem of missing scores for a clip by computing a mixture for only the observed data types [44]. This is done by renormalizing the mixture

weights over the observed data types for each clip. The training was done with the EM algorithm, but the maximization step no longer had a closed-form solution, therefore we used gradient-descent techniques to learn the optimal weights.

SVMLight This fusion model consists of training an SVM using the scores from various data types as the features for each clip. Missing data types for a given clip are assigned zero scores. We used the SVMLight⁴ implementation with linear kernels.

Distance from threshold This is a weighted averaging method [3] that dynamically adjusts the weights of each data type for each video clip based on how far the score is from its decision threshold. If the detection score is near the threshold, the correct decision is presumed to be somewhat uncertain, and a lower weight is assigned. A detection score that is much greater or much lower than the threshold indicates that more confidence should be placed in the decision, and a higher weight is assigned.

Bin accuracy weighting This method tries to address the problem of uneven distribution of detection scores in the training set. For each data type, the range of scores in the training fold is divided into bins with approximately equal counts per bin. During training, the accuracy of each bin is measured by computing the proportion of correctly classified videos whose scores fall within the bin. During testing, for each data type, the specific bin that the scores fall into is determined, and the corresponding bin accuracy scores for each data type are used as fusion weights.

Table 4 summarizes the fusion methods and their characteristics.

3 Experimental results

We evaluated the performance of our SESAME system using the data provided in the TRECVID MED evaluation task. Although the MED event kit contained both a text description and video examples for each event, the SESAME system implemented the example-based approach in which only the video examples were used for event detection training.

3.1 Evaluation by data type

Table 5 lists results on the SESAME Evaluation dataset. In terms of the performance of the various data types, the visual features were the strongest performers across all events. The accuracy of the visual concepts was nearly as strong as that

⁴ <http://svmlight.joachims.org/>.

Table 4 Fusion methods and their characteristics

Fusion method	Event-independent?	Learned on a training set?	Score-dependent?	Adjustment for missing scores?
Arithmetic mean	Yes	No	No	Yes
Geometric mean	Yes	No	No	Yes
MAP weighted	No	Yes	No	Yes
Weighted mean root	No	Yes	No	No
Conditional mixture model	No	Yes	No	No
Sparse mixture model	No	Yes	No	Yes
SVMLight	No	Yes	Yes	No
Distance from threshold	No	Yes	Yes	No
Bin accuracy weighting	No	Yes	Yes	No

Table 5 Experiment results in terms of mean average precision for individual event classifiers

Event	Low-level visual features		Visual concept features		Motion features			Audio		Fusion
	SIFT-AVG	SIFT-DC	RF	SVM	STIP	DT	MOSIFT	MFCC	ASR	AM
Birthday_party	0.275	0.229	0.339	0.324	0.189	0.293	0.191	0.146	0.062	0.372
Changing_a_vehicle_tire	0.305	0.270	0.251	0.241	0.136	0.217	0.126	0.024	0.209	0.343
Flash_mob_gathering	0.603	0.644	0.542	0.542	0.569	0.567	0.463	0.139	0.017	0.644
Getting_a_vehicle_unstuck	0.457	0.496	0.454	0.426	0.365	0.439	0.337	0.040	0.011	0.586
Grooming_an_animal	0.280	0.222	0.254	0.231	0.147	0.247	0.290	0.038	0.024	0.352
Making_a_sandwich	0.267	0.278	0.283	0.257	0.225	0.234	0.164	0.038	0.378	0.392
Parade	0.416	0.414	0.373	0.306	0.457	0.419	0.326	0.119	0.013	0.578
Parkour	0.464	0.414	0.550	0.479	0.369	0.459	0.295	0.029	0.009	0.564
Repairing_an_appliance	0.486	0.469	0.422	0.404	0.385	0.443	0.368	0.449	0.517	0.591
Working_on_a_sewing_project	0.378	0.388	0.390	0.394	0.386	0.433	0.270	0.192	0.276	0.551
Attempting_a_bike_trick	0.398	0.350	0.475	0.472	0.235	0.438	0.640	0.019	0.003	0.703
Cleaning_an_appliance	0.138	0.077	0.097	0.149	0.074	0.089	0.090	0.050	0.144	0.174
Dog_show	0.591	0.650	0.595	0.529	0.557	0.632	0.488	0.183	0.002	0.672
Giving_directions_to_a_location	0.123	0.130	0.058	0.097	0.191	0.052	0.085	0.075	0.066	0.193
Marriage_proposal	0.057	0.093	0.077	0.066	0.173	0.118	0.027	0.044	0.010	0.179
Renovating_a_home	0.229	0.273	0.295	0.325	0.255	0.361	0.157	0.099	0.145	0.461
Rock_climbing	0.488	0.466	0.412	0.401	0.352	0.425	0.465	0.020	0.005	0.615
Town_hall_meeting	0.531	0.463	0.411	0.417	0.462	0.370	0.519	0.433	0.341	0.649
Winning_a_race_without_a_vehicle	0.237	0.284	0.198	0.167	0.260	0.216	0.273	0.074	0.005	0.295
Working_on_a_metal_crafts_project	0.109	0.133	0.099	0.162	0.032	0.128	0.116	0.024	0.044	0.209
Mean for all events	0.342	0.337	0.341	0.330	0.291	0.329	0.285	0.112	0.114	0.456

The data type with the highest MAP score for each event is in bold. The AM fusion of the individual event classifiers is listed in the last column

of the low-level visual features. The motion features also showed strong performance. Although the performance of low-level audio features and ASR was significantly less, ASR had the highest performance for events containing a relatively large amount of speech content, including a number of instructional videos. The best scores for each event are distributed among all of the data types, indicating that fusion of these data should yield improved performance. Indeed, the

AM fusion of the individual event classifiers, which is listed in the last column of Table 5, shows a significant boost in performance: a 33 % improvement over the best single data type.

3.2 Evaluation of fusion methods

We tested the late fusion methods described in Sect. 2.3 using the SESAME Evaluation dataset. For all our fusion

Table 6 MED performance of fusion methods with all event classifiers

Fusion method	Macro MAP	Standard deviation
Arithmetic mean	0.456	0.0000
Geometric mean	0.456	0.0000
MAP-weighted	0.437	0.0006
Weighted Mean Root	0.451	0.0005
Conditional mixture model	0.403	0.0054
Sparse mixture model	0.443	0.0007
SVMLight	0.451	0.0036
Distance from threshold	0.407	0.0005
Bin accuracy weighting	0.401	0.0031

experiments, we trained each event classifier on the training set, and executed the classifier on the test set to produce detection scores for each event. To produce legitimate fusion scores over the test set, we used tenfold cross validation, with random fold selection, to generate the detections, and then obtained a micro-averaged average precision over the resulting detections. The micro-averaged MAP was computed by averaging the average precision for each event. To gauge the stability of the fusion methods, we repeated this process 30 times and computed the macro average and standard deviation of the micro-averaged MAPs. Because the AM and GM methods are untrained, their micro-averaged MAPs will be the same regardless of fold selection; thus, the standard deviations for their micro-averaged MAPs are zero.

Table 6 shows the MED performance of various fusion methods. The comparison indicates that the simplest fusion methods, such as AM and GM, performed as well as or better than other, more complex fusion methods. Also note that most of the top-performing fusion methods (AM, GM, MAP, WMR, and SMM) adjusted their weights to accommodate missing scores.

3.3 Evaluation of MED performance in TRECVID

As the SESAME team, we participated in the 2012 TRECVID MED evaluation and submitted the detection results for a system configured nearly the same as that described in this paper⁵. The event classifiers were trained with all the positives from the event kit and negatives from the TRECVID MED training and development material. The test set consisted of the 99,000 videos used in the formal evaluation.

Figure 5 shows the performance of the primary runs of 17 MED systems in this evaluation in terms of miss and false alarm rates [45]. The performance of the SESAME run was one of the best among the evaluation participants.

⁵ It included a poorer-performing ASR capability instead of the one described in Sect. 2.2.6, and a video OCR capability that contributed minimally to overall performance.

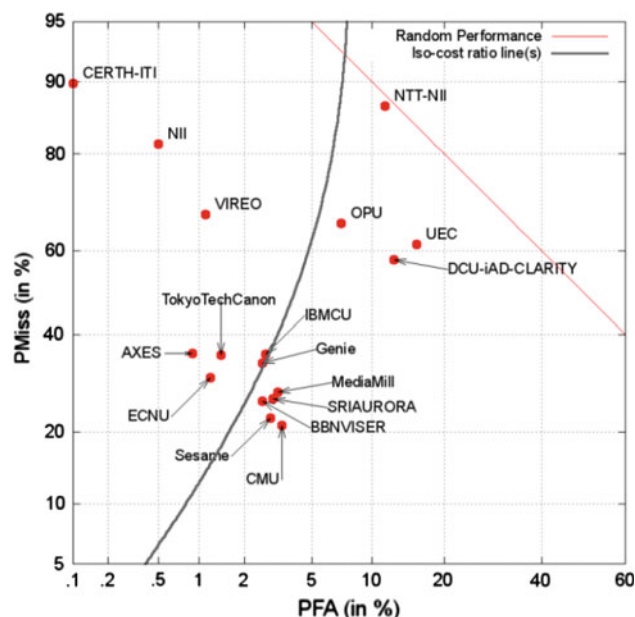


Fig. 5 Performance of the primary runs of 17 MED systems in the 2012 TRECVID MED evaluation

4 Summary and discussion

Search with Speed and Accuracy for Multimedia Events, a MED capability that learns event models from a set of example video clips, includes a number of BOW event classifiers based on single data types: low-level visual, motion, and audio features; high-level semantic visual concepts; and ASR. Partitioning the representation by data type permits the descriptors for each data type to be optimized independently. We evaluated the detection performance for each event classifier and experimented with a number of fusion methods for combining the event detection scores from these classifiers. Our experiments using multiple data types and late fusion of their scores demonstrated strongly reliable MED performance.

Major conclusions from this effort include:

- The relative contribution of visual, motion, and audio features varies according to the specific event. This is due to differences in the relative distinctiveness and consistency of certain features for each event category. Across all events, score-level fusion resulted in a 33 % improvement over the best single data type, indicating that different types of features contribute to the representation of heterogeneous video data.
- The use of difference coding in low-level visual and motion features significantly improved performance. We surmise that difference coding works better than the traditional BOW because it measures differences from the general model, which is likely to be dominated by the background features. We expect additional gains in per-

formance if difference coding was applied to low-level audio features.

- The set of 1,346 high-level visual features was nearly as effective as the set of low-level visual features. It appears that, in comparison to the 5,000 or so concepts predicted to be needed for sufficient performance in event detection [46], this number of high-level features begins to span the space of concepts reasonably well. Therefore, analogous sets of motion and audio concepts should further improve the overall performance.
- Although the performance of ASR was lower than that of the visual and motion features, its performance was highly event dependent, and it performed reasonably well for events containing a relatively large amount of speech content, such as instructional videos.
- The simplest fusion methods for computing event detection scores were very effective compared to more complex fusion methods. One possible explanation for this is that the reliability of the scores is roughly equal across all data types. Another possible reason is that the limited number of positive training examples (an average of about 70 per event) is not enough to achieve the full benefit of the more complex fusion models.

While our relatively straightforward BOW approach was quite effective, we view it as a baseline capability that could be improved in several ways:

- Since the current approach aggregates low-level visual and motion features within fixed spatial partitions, the usage of local information is limited. Features of an object divided by our predefined partition, for example, will not be aggregated as a whole. We expect that the use of dynamic spatial pooling, which is better aligned to the structure and content of the video imagery, will improve performance. Segmenting the image into meaningful homogeneous regions would be even better, as it allows more salient characteristics to be extracted, and would eventually lead to better classification.
- The current approach ignores the temporal information within each video clip; all the visual, motion, and audio features are aggregated. However, events consist of multiple components that appear at different times, therefore using time-based information for event modeling and detection should improve performance. In addition, aggregating low-level features according to the temporal structure of the video may yield feature sets that better represent the video contents.
- All the classifiers in our approach operate on a histogram of features and do not leverage any relationships between the features. Features occurring in video data are not generally independent. In particular, the combination of particular high-level semantic concepts could become strong discriminatory evidence, since their co-occurrence might

be associated with a subset of relevant video content. For example, although the concepts balloons and singing occur in many contexts, the occurrence of both might be more common to birthday party than to other video content. Exploiting the spatiotemporal dependencies among the features would better characterize the video contents and offer a richer set of data with which to build event models.

Acknowledgments This work was supported by the Intelligence Advanced Research Projects Activity (IARPA) via the Department of Interior National Business Center, Contract Number D11PC0067. The US Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of IARPA, DoI/NBC, or the US Government.

Open Access This article is distributed under the terms of the Creative Commons Attribution License which permits any use, distribution, and reproduction in any medium, provided the original author(s) and the source are credited.

References

1. Smeaton, A.F., Over, P., Kraaij, W.: Evaluation campaigns and TRECVID. In: Proceedings of the 8th ACM international workshop on multimedia information retrieval, Santa Barbara, 26–27 October 2006 (MIR '06). ACM Press, New York, pp. 321–330 (2006)
2. Jiang, Y.-G., Bhattacharya, S., Chang, S.-F., Shah, M.: High-level event recognition in unconstrained videos. *Int. J. Multimed. Inform. Retr.* 1–29 (2012)
3. Natarajan, P., Wu, S., Vitaladevuni, S., Zhuang, X., Tsakalidis, S., Paurk, U., Prasad, R.: Multimodal feature fusion for robust event detection in web videos. In: Proceedings of the IEEE computer society conference on computer vision and pattern recognition (CVPR), pp. 1298–1305 (2012)
4. Sawhney, H., Cheng, H., Divakaran, A., Javed, O., Liu, J., Yu, Q., Ali, S., Tamrakar, A.: Evaluation of low-level features and their combinations for complex event detection in open source videos. *CVPR*, 2496–2499 (2012)
5. Jiang, Y.: Super: towards real-time event recognition in internet videos. *ACM Int. Conf. Multimed. Retr. (ICMR)* (2012) (article no. 33)
6. Ballan, L., Bertini, M., Del Bimbo, A., Seidenari, L., Serra, G.: Event detection and recognition for semantic annotation of video. *Multimed. Tools Appl.* 51(1), 279–302 (2011)
7. Xu, D., Chang, S.-F.: Video event recognition using kernel methods with multilevel temporal alignment. *IEEE Trans. Pattern Anal. Mach. Intell. (IEEE TPAMI)* 30(11), 1985–1997 (2008)
8. Snoek, C.G.M., Worring, M.: Concept-based video retrieval. *Found. Trends Inf. Retr.* 2(4), 214–322 (2009)
9. Felzenszwalb, P., Girshick, R., McAllester, D., Ramanan, D.: Object detection with discriminatively trained part-based models. *IEEE TPAMI* 32(9), 1627–1645 (2010)
10. Li, L., SU, H., Xing, E., Fei-Fei, L.: Object bank: a high-level image representation for scene classification and semantic feature sparsification. *Adv. Neural Inf. Process. Syst.*, 24 (2010)
11. Sadanand, S., Corso, J.J.: Action bank: a high-level representation of activity in video. *CVPR* (2012)

12. Snoek, C.G.M., Smeulders, A.W.M.: Visual-concept search solved? *IEEE Comput.* **43**(6), 76–78 (2010)
13. Merler, M., Huang, B., Xie, L., Hua, G., Natsev, A.: Semantic model vectors for complex video event recognition. *IEEE Trans. Multimed. (TMM)* **14**(1), 88–101 (2012)
14. Althoff, T., Song, H., Darrell, T.: Detection bank: an object detection based video representation for multimedia event recognition. *ACM Multimed. (MM)* (2012)
15. Tsampoulatidis, I., Gkalelis, N., Dimou, A., Mezaris, V., Kompatsiaris, I.: High-level event detection in video exploiting discriminant concepts. In: *Proceedings of the 1st ACM international conference on multimedia retrieval*, pp. 85–90 (2011)
16. Habibian, A., van de Sande, K.E.A., Snoek, C.G.M.: Recommendations for video event recognition using concept vocabularies. In: *Proceedings of the ACM international conference on multimedia retrieval*, pp. 89–96 Dallas (2013)
17. Perera, A.G.A., Oh, S., Leotta, M., Kim, I., Byun, B., Lee, C.-H., McCloskey, S., Liu, J., Miller, B., Huang, Z.F., Vahdat, A., Yang, W., Mori, G., Tang, K., Koller, D., Fei-Fei, L., Li, K., Chen, G., Corso, J., Fu, Y., Srihari, R.: GENIE TRECVID 2011 multimedia event detection: late-fusion approaches to combine multiple audio-visual features. In: *NIST TRECVID, workshop* (2011)
18. Strassel, S., Morris, A., Fiscus, J., Caruso, C., Lee, H., Over, P., Fiumara, J., Shaw, B., Antonishek, B., Michel, M.: Creating HAVIC: heterogeneous audio visual internet collection. In: Calzolari N., Choukri K., Declerck T., Uğur Doğan M., Maegaard B., Mariani J., Odijk J., Piperidis S. (eds.) *Proceedings of the eighth international conference on language resources and evaluation, Istanbul* (2012)
19. Jaakkola, T., Haussler, D.: Exploiting generative models in discriminative classifiers. In: *Proceedings of the 1998 conference on advances in neural information processing systems II*, pp. 489–493 (1999)
20. Jégou, H., Perronnin, F., Douze, M., Sanchez, J., Pérez, P., Schmid, C.: Aggregating local image descriptors into compact codes. *IEEE TPAMI* **34**(9), 1704–1716 (2012)
21. Perronnin, F., Dance, C.: Fisher kernels on visual vocabularies for image categorization. *CVPR*, (2007)
22. Snoek, C.G.M., van de Sande, K.E.A., Habibian, A., Kordumova, S., Li, Z., Mazloom, M., Pinteá, S.L., Tao, R., Koelma, D.C., Smeulders, A.W.M.: The MediaMill TRECVID 2012 semantic video search engine. In: *Proceeding of the TRECVID workshop, Gaithersburg* (2012)
23. Tuytelaars, T., Mikolajczyk, K.: Local invariant feature detectors: a survey. *Found. Trends. Comput. Graph. Vis.* **3**(3), 177–280 (2008)
24. Snoek, C.G.M., Worring, M., Geusebroek, J.-M., Koelma, D.C., Seinstra, F.J.: On the surplus value of semantic video analysis beyond the key frame. In: *Proceedings of the IEEE international conference on multimedia and expo* (2005)
25. Lazebnik, S., Schmid, C., Ponce, J.: Beyond bags of features: spatial pyramid matching for recognizing natural scene categories. *CVPR 2*, 2169–2178 (2006) (New York)
26. van de Sande, K.E.A., Gevers, T., Snoek, C.G.M.: Evaluating color descriptors for object and scene recognition. *IEEE TPAMI* **32**(9), 1582–1596 (2010)
27. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vis.* **60**, 91–110 (2004)
28. Geusebroek, J.-M., Boomgaard, R., Smeulders, A.W.M., Geerts, H.: Color invariance. *IEEE TPAMI* **23**(12), 1338–1350 (2001)
29. van Gemert, J.C., Snoek, C.G.M., Veenman, C.J., Smeulders, A.W.M., Geusebroek, J.-M.: Comparing compact codebooks for visual categorization. *Comput. Vis. Image Underst.* **114**(4), 450–462 (2010)
30. Maji, S., Berg, A.C., Malik, J.: Classification using intersection kernel support vector machines is efficient. In: *Proceedings of the IEEE computer society conference on CVPR*, pp. 619–626, Anchorage (2008)
31. Wang, H., Kläser, A., Schmid, C., Cheng-Lin, L.: Action recognition by dense trajectories. *CVPR*, 3169–3176 (2011)
32. Chen, M.-Y., Hauptmann, A.: MoSIFT: recognizing human actions in surveillance videos. *CMU-CS-09-161*. Carnegie Mellon Univ. (2009)
33. Laptev, I.: On space-time interest points. *Int. J. Comput. Vis.* **64**(2/3), 107–123 (2005)
34. Sun, C., Nevatia, R.: Large scale web video classification by use of Fisher vectors. In: *Workshop on applications of computer vision, Clearwater* (2013) (January)
35. Chechik, G., Ie, E., Rehn, M., Bengio, S., Lyon, D.: Large-scale content-based audio retrieval from text queries. In: *Proceedings of 1st ACM international conference on multimedia information retrieval (MIR '08)*, pp. 105–112, New York (2008)
36. Uchida, Y., Sakazawa, S., Argawal, M., Akbacak, M.: KDDI labs and SRI international at TRECVID 2010: content-based copy detection. In: *NIST TRECVID 2010 evaluation, workshop* (2010)
37. Jiang, Y., Zeng, X., Ye, G., Ellis, D., Shah, M., Chang, S.: Columbia-UCF TRECVID 2010 multimedia event detection: combining multiple modalities, contextual concepts, and temporal matching. In: *NIST TRECVID, workshop* (2010)
38. Pancoast, S., Akbacak, M.: Bag-of-audio-words approach for multimedia event detection. In: *Proceedings of interspeech* (2012)
39. Merler, M., Huang, B., Xie, L., Hua, G., Natsev, A.: Semantic model vectors for complex video event recognition. *IEEE Trans. Multimed.* **14**(1), 88–101 (2012)
40. Over, P., Awad, G., Michel, M., Fiscus, J., Sanders, G., Shaw, B., Kraaij, W., Smeaton, A.F., Quénot, G.: TRECVID 2012—an overview of the goals, tasks, data, evaluation mechanisms, and metrics. In: *Proceedings of TRECVID* (2012) <http://www-nlpir.nist.gov/projects/tvpubs/tv12.papers/tv12overview.pdf>
41. Berg, A., Deng, J., Satheesh, S., Su, H., Li, F.-F.: Imagenet large scale visual recognition challenge (2011) <http://www.image-net.org/challenges/LSVRC/2011/>
42. Janin, A., Stolcke, A., Anguera, X., Boakye, K., Çetin, Ö., Frankel, J., Zheng, J.: The ICSI–SRI spring 2006 meeting recognition system, MLMI'06. In: *Proceedings of the third international conference on machine learning for multimodal, interaction*, pp. 444–456 (2006)
43. van Hout, J., Akbacak, M., Castaneda, D., Yeh, E., Sanchez, M.: Extracting audio and spoken concepts for multimedia event detection. In: *International conference on acoustics, speech, and signal processing (ICASSP)* (2013)
44. Nallapati, R., Yeh, E., Myers, G.: Sparse mixture model: late fusion with missing scores for multimedia event detection. *Algorithms and systems VII. SPIE Multimed. Content Access* (2012)
45. Fiscus, J., Michel, M.: TRECVID 2012 multimedia event detection task. In: *NIST TRECVID 2012 evaluation, workshop* (2012)
46. Hauptmann, A., Yan, R., Lin, W.-H., Christel, M., Wactlar, H.: Can high-level concepts fill the semantic gap in video retrieval? A case study with broadcast retrieval. *IEEE Trans Multimed.* **9**(5), 958–966 (2007)

Author Biographies



Gregory K. Myers is an associate director in the Robotics Program at SRI International. He specializes in developing computer vision, video, speech, robotics, and biometrics technologies and systems for specialized applications. He is the Principal Investigator of the SESAME project. His recent research includes the development of a video OCR capability for overlay and in-scene text. He led the development and transitioning of mail address reading

technology for nationwide deployment for the US Postal Service. He holds a BSEE from the Massachusetts Institute of Technology and an MSEE from University of California at Berkeley.



Ramesh Nallapati received his Ph.D. in computer science from University of Massachusetts at Amherst in 2006. He did his postdoctoral studies at Carnegie Mellon University and has held research positions at Stanford University and SRI International. Currently, he is a research staff member at IBM T. J. Watson Research Center. His research interests are in machine learning and their applications to information retrieval, data mining and natural language processing.



Julien van Hout studied electrical engineering and computer science at Ecole Polytechnique, France where he received his B.Sc. in 2010. He then obtained his M.Sc. degree in electrical engineering (2012) from the University of California, Los Angeles. He is now a research engineer at the Speech Technology and Research laboratory of SRI International in Menlo Park, CA. His

research interests are in automatic speech recognition applied to keyword spotting and multimedia event detection.



Stephanie Pancoast is currently a Ph.D. candidate in electrical engineering at Stanford University. She also works as a student associate in the Speech Technology and Research Laboratory at SRI International in Menlo Park, CA. She received her M.Sc. in electrical engineering at Stanford University in 2012 and B.Sc. at Cornell University in 2010. Her research interests include audio and speech processing applied to multimedia.



Ramakant Nevatia received his Ph.D. from Stanford University and is currently a professor of computer science and electrical engineering at the University of Southern California, where he is the founding director of the Institute for Robotics and Intelligent Systems (IRIS). Prof. Nevatia has more than 30 years of experience in computer vision research and has published more than 100 refereed technical papers. He is a fellow of the Institute of Electrical and Electronic Engineers (IEEE) and American Association for Artificial Intelligence (AAAI) and a member of the ACM. In recent years, his research has been focused on detecting and tracking humans and vehicles in video and on techniques for representation and recognition of events.



Chen Sun is a Ph.D. student at Institute for Robotics and Intelligent Systems, University of Southern California. His main research interest is on video analysis.



Amirhossein Habibian received his B.Sc. in computer engineering (2008) and his M.Sc. in artificial intelligence (2011), both from Electrical and Computer Engineering Department of University of Tehran, Iran. Currently, he is a Ph.D. candidate at the University of Amsterdam. His research interest is image and video understanding, covering computer vision, statistical pattern recognition and information retrieval.



Dennis C. Koelma is senior scientific programmer at the UvA. He received his M.Sc. and Ph.D. degrees in computer science from the University of Amsterdam in 1989 and 1996, respectively. The subject of his thesis is "A software environment for image interpretation". His research interests include image and video processing, software architectures, parallel programming, databases, graphical user interfaces and visual information systems. He is the

lead designer and developer of Impala: a software architecture for accessing the content of digital images and video. The software serves as a platform for consolidating software resulting from ISIS research. It has been licensed by the UvA spin-off EUvision where he has a part-time affiliation.



Koen E. A. van de Sande received his B.Sc. in computer science (2004), a B.Sc. in artificial intelligence (2004), M.Sc. in computer science (2007) and Ph.D. in computer science (2011) from the University of Amsterdam, The Netherlands. Currently, he is a parttime researcher at the University of Amsterdam and works in R&D at Euvision Technologies. His research interests include computer vision, object recognition and localisation, machine learning, parallel

computing and large-scale benchmark evaluations. At Euvision Technologies, he is bridging the gap between the latest techniques developed in academia and real-world, large-scale computer vision applications.



Arnold W. M. Smeulders graduated from Delft Technical University in Physics on texture in medical images. He received his Ph.D. from Leyden Medical Faculty. Since 1994, he is professor at the University of Amsterdam, leading the ISIS group on visual search engines. The search engine has been a top three performer in the TRECvid competition for the last 8 years. ISIS came out best in the 6-yearly international review in 2003 and 2009 (shared the maximum with

few). In 2010 he co-founded Euvision, a university spin-off. In addition, recently, he is director of COMMIT, the large public private-ICT-research program of the Netherlands. He is past associate editor of the IEEE trans PAMI and currently of IJCV. He is IAPR Fellow and an honorary member of NVPBHV. He was visiting professor in Hong kong, Tsukuba, Modena, Cagliari and Orlando.



Cees G. M. Snoek received his M.Sc. degree in business information systems (2001) and Ph.D. degree in computer science (2005) from the University of Amsterdam, The Netherlands. He is currently an assistant professor in the Intelligent Systems Lab at the University of Amsterdam. He was a visiting scientist at Carnegie Mellon University, Pittsburgh, PA (2003) and at the University of California, Berkeley, CA (2010–2011). His research interest is video and

image search. Dr. Snoek is the lead researcher of the MediaMill Semantic Video Search Engine, which is a consistent top performer in the yearly NIST TRECVID evaluations. He is member of the editorial boards for IEEE Multimedia and IEEE Transactions on Multimedia. Cees is recipient of an NWO Veni award (2008), an NWO Vidi award (2012) and the Netherlands Prize for ICT Research (2012). Several of his Ph.D. students have won best paper awards, including the IEEE Transactions on Multimedia Prize Paper Award.