

# Detection Bank: An Object Detection Based Video Representation for Multimedia Event Recognition

Tim Althoff  
UC Berkeley EECS/ICSI  
althoff@icsi.berkeley.edu

Hyun Oh Song  
UC Berkeley EECS/ICSI  
song@eecs.berkeley.edu

Trevor Darrell  
UC Berkeley EECS/ICSI  
trevor@eecs.berkeley.edu

## ABSTRACT

While low-level image features have proven to be effective representations for visual recognition tasks such as object recognition and scene classification, they are inadequate to capture complex semantic meaning required to solve high-level visual tasks such as multimedia event detection and recognition. Recognition or retrieval of events and activities can be improved if specific discriminative objects are detected in a video sequence. In this paper, we propose an image representation, called *Detection Bank*, based on the detection images from a large number of windowed object detectors where an image is represented by different statistics derived from these detections. This representation is extended to video by aggregating the key frame level image representations through mean and max pooling. We empirically show that it captures complementary information to state-of-the-art representations such as Spatial Pyramid Matching and Object Bank. These descriptors combined with our Detection Bank representation significantly outperforms any of the representations alone on TRECVID MED 2011 data.

## Categories and Subject Descriptors

H.3.1 [Information Storage and Retrieval]: Content Analysis and Indexing

## General Terms

Experimentation, Performance, Algorithms

## Keywords

Object Detection, Representation, TRECVID, Multimedia Event Recognition

## 1. INTRODUCTION

There has been considerable progress in classifying videos using concept models formed from global feature representations. This approach learns a generic mapping from images

to a set of predefined scene labels, which are then used as features to classify and/or retrieve a specific target query. Currently this class of method achieves state-of-the-art performance on scene-level image classification and many video retrieval tasks [5, 12].

However, recognition and retrieval of events and activities may require discrimination based on specific objects in a scene; for example, the difference between a birthday party and a wedding ceremony may lie in particular clothing, objects such as candles and balloons in the scene, or even the type of cake.

Scene-level visual descriptors are generally inadequate to capture these fine-grained phenomena. Objects of interest may not cover the entire scene, and thus low-level visual descriptors that are pooled over the entire scene may fail to detect them. To be sensitive to the presence of individual objects in a scene requires visual analysis that considers windows or segments in a scene.

Recent advances in deformable part object detection have demonstrated that objects can be reliably found in many natural scenes, e.g., as demonstrated in the PASCAL VOC challenge [2]. Progress is ongoing, but some categories are reasonably well detected. Most localized detection models employ either a top-down window scanning methodology, or a bottom-up segment generation strategy. Recently, the Object Bank model was proposed as a representation for indexing video based on the max pooled score map output of a bank of individual object detectors, each of which operated over windows of the scene [8]. The detectors are based on the deformable part model [4] and were trained to recognize about 200 different object categories.

We extend this idea of using the output of object detectors to guide event recognition. While the Object Bank model provides a dense map of max pooled detection scores, it lacks a more immediate sense of whether or not there are objects present in the image and if so how many (e.g. the number of person detections presumably helps in differentiating “Attempting a board trick” from “Flash mob gathering”). We propose to compute different statistics from the detection images from object detectors to capture this information more directly. The detection images are obtained by thresholding the score maps from the object detectors, applying non-maximum suppression, and pooling across all scales such that the detection images contains all the final detections. Object Bank model omits all these steps that are standard in a detection pipeline which would result in a more sparse representation. We demonstrate that these steps along with the proposed detection count statistics lead

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

MM'12, October 29–November 2, 2012, Nara, Japan.

Copyright 2012 ACM 978-1-4503-1089-5/12/10 ...\$15.00.

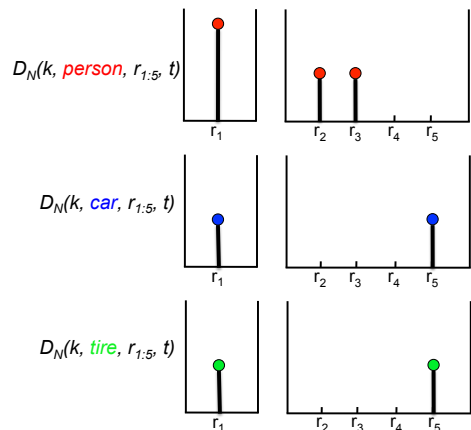
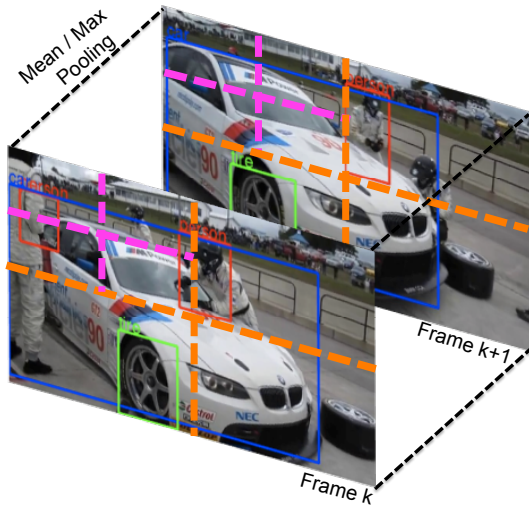


Figure 1: Left: Example detections on two successive keyframes with a subset of object detectors. Red, green and blue bounding boxes show person, tire, and car detections respectively. Orange and magenta dotted lines show a subset of the spatial pyramid. Right: Illustration of the detection count feature for the three categories on the first 5 grid cells of frame  $k$ . Grid cell  $r_1$  has the extent of the whole image,  $r_2, r_3, r_4, r_5$  are the cells in the orange quadrant clockwise from the top left corner. Best viewed in color.

to significantly better classification performance. This image or key frame level representation is extended to a video level representation by mean and max pooling across keyframes which allows for an intuitive interpretation (see Section 3 for more detail). In light of the Object Bank model we call this new representation *Detection Bank*.

We empirically show that combining features of this proposed video representation with existing global features significantly outperforms any of the methods alone on TRECVID MED 2011 data. Our results suggest that high-level visual features from a bank of object detections can be used to complement present multi modal concept detection systems (such as [9]).

The rest of the paper is organized as follows. Section 2 gives an overview of related work in image representations, object detection, and multimedia event detection. Section 3 describes the proposed Detection Bank representation in more detail. An evaluation of the approach is presented in Section 4 before Section 5 concludes the paper with a short discussion of results as well as future work ideas.

## 2. RELATED WORK

The Spatial pyramid match (SPM) [7] computes histograms of visual words from gradient orientation statistics at multiple scales. A generative topic model on bag of visual words to model the image scene generation process has been explored in [3]. The Deformable part model (DPM) [4] learns an object detector model from a weakly supervised setting where only the root location is supervised and the locations of parts are not known. It is the current state-of-the-art object detector on PASCAL VOC. Object Bank (OB) [8] uses the max pooled intermediate score maps from deformable part model object detectors as a global feature representation of images. In order to capture the presence of certain objects more directly in our representation we propose to explicitly build detection images and compute detection count statistics instead of only using the intermediate score maps. Using DPM object detectors as high-level visual features for

the task of multimedia event detection has been explored in [9]. The authors report that their feature representation did not lead to a significant performance improvement (with the small exception of car detections and the “Getting vehicle unstuck” event). In contrast, we present an approach that leads to significant performance increase across all event categories.

## 3. DETECTION BANK REPRESENTATION

As motivated earlier, object detections can be used to discriminate between certain events and activities, e.g., a large number of occurrences of flags and cars provides a strong cue to distinguish “Parade” from “Flash mob gathering” while they might both contain multiple persons. Similarly, detecting a large number of candles and balloons but no wedding dresses and bridegrooms speaks in favor of “Birthday party” instead of “Wedding ceremony”.

As noted in the introduction, the Object Bank representation omits thresholding, non-maximum suppression, and pooling across different scales that are well-known mechanisms in the object detection community. Our Detection Bank representation explicitly builds detection images and computes the following detection count statistics (per object category) for each grid cell in a spatial pyramid (entire image,  $2 \times 2$ , and  $4 \times 4$  grid): the sum of scores of detections within that cell (above a certain threshold), the number of detections, and a single bit that indicates whether or not there is a detection within in that cell. By mean and max pooling these statistics across keyframes of a video we obtain a meaningful video-level representation capturing, e.g., the maximum number of detections, the average number of detections, and an empirical estimate of the detection probability for each grid cell and object category. It will be demonstrated in Section 4 that these statistics contain discriminative information that is complementary to both scene-level features and max pooled detection score maps (as used in Object Bank).

Formally, an object detector for category  $c$  searches over

every location  $(x, y)$  of an image  $\mathcal{I}$  and outputs  $P$  predicted locations of the object in terms of bounding boxes  $\mathcal{B}_c$  on  $\mathcal{I}$ ,  $\mathcal{B}_c = [\mathbf{b}_{c,1}, \dots, \mathbf{b}_{c,P}]$  where  $\mathbf{b}_{c,i} = [x_{i1}, y_{i1}, x_{i2}, y_{i2}, \text{score}_i]$ . We consider the following three statistics per key frame  $k$  from windowed object detectors at detection threshold  $t$ :

$$\begin{aligned} D_S(k, c, r, t) &= \sum_{i=1}^P \mathbb{I}[\overline{\mathbf{b}_{c,i}} \in \mathcal{I}(r)] \mathbb{I}[s(\mathbf{b}_{c,i}) \geq t] s(\mathbf{b}_{c,i}) \\ D_N(k, c, r, t) &= \sum_{i=1}^P \mathbb{I}[\overline{\mathbf{b}_{c,i}} \in \mathcal{I}(r)] \mathbb{I}[s(\mathbf{b}_{c,i}) \geq t] \\ D_0(k, c, r, t) &= \mathbb{I}\left[\sum_{i=1}^P (\mathbb{I}[\overline{\mathbf{b}_{c,i}} \in \mathcal{I}(r)] \mathbb{I}[s(\mathbf{b}_{c,i}) \geq t]) > 0\right] \end{aligned} \quad (1)$$

where  $\mathbb{I}[\cdot]$  is the indicator function,  $t$  is a specific detection threshold, and  $\overline{\mathbf{b}_{c,i}}$  denotes the center of the bounding box.  $\mathcal{I}(r)$  denotes a spatial grid cell indexed by  $r$ .  $s(\mathbf{b}_{c,i})$  denotes the score of the bounding box.  $D_S(k, c, r, t)$  is a detection statistic summing over the scores of detections within spatial grid region  $r$  for a specific category  $c$  and detection threshold  $t$ .  $D_N(k, c, r, t)$  computes the number of detections within region  $r$ .  $D_0(k, c, r, t)$  is a binary feature whether there are any detections within the region. As more conservative threshold values are applied the true positive rate increases at the expense of less detections. Overall, thresholding scores and applying non-maximum suppression lead to a representation that is easier for a linear classifier to learn as demonstrated in Section 4.

Finally, we aggregate the key frame level statistics to the feature vector  $\mathcal{F}$  for each video  $\mathcal{V} = \{\mathcal{I}_1, \dots, \mathcal{I}_K\}$  by mean and max pooling  $\mathcal{P}$  across keyframes. Denote by  $K, C, R, T$  the total number of frames, object categories, spatial regions, and threshold levels respectively.

$$\begin{aligned} \mathcal{P}(\mathcal{V}, c, r, t) &: \mathbb{R}^{3K} \mapsto \mathbb{R}^3 \\ \mathcal{P}_{max}(\mathcal{V}, c, r, t) &= \begin{bmatrix} \max_{k \in K} D_S(k, c, r, t) \\ \max_{k \in K} D_N(k, c, r, t) \\ \max_{k \in K} D_0(k, c, r, t) \end{bmatrix} \\ \mathcal{P}_{mean}(\mathcal{V}, c, r, t) &= \begin{bmatrix} \text{mean}_{k \in K} D_S(k, c, r, t) \\ \text{mean}_{k \in K} D_N(k, c, r, t) \\ \text{mean}_{k \in K} D_0(k, c, r, t) \end{bmatrix} \end{aligned} \quad (2)$$

$$\mathcal{F}_{max}(\mathcal{V}) = (\mathcal{P}_{max}(\mathcal{V}, c_1, r_1, t_1), \dots, \mathcal{P}_{max}(\mathcal{V}, c_C, r_R, t_T))$$

$$\mathcal{F}_{mean}(\mathcal{V}) = (\mathcal{P}_{mean}(\mathcal{V}, c_1, r_1, t_1), \dots, \mathcal{P}_{mean}(\mathcal{V}, c_C, r_R, t_T))$$

Figure 1 illustrates our Detection Bank feature representation on two successive keyframes on a video from ‘‘Changing a vehicle tire’’ event.

## 4. EXPERIMENTS

We evaluated the proposed detection-based video representation in the realm of multimedia event classification using a forced-choice classification paradigm on the TRECVID MED 2011 Event Kit that contains 2025 videos from 15 different events [10] (see Table 1). Note that this study differs from the detection paradigm on the larger DEV-T and DEV-O collections of the TRECVID MED 2011 data set in order

EventID	Event Name
1	Attempting a board trick
2	Feeding an animal
3	Landing a fish
4	Wedding ceremony
5	Working on a woodworking project
6	Birthday party
7	Changing a vehicle tire
8	Flash mob gathering
9	Getting a vehicle unstuck
10	Grooming an animal
11	Making a sandwich
12	Parade
13	Parkour
14	Repairing an appliance
15	Working on a sewing project

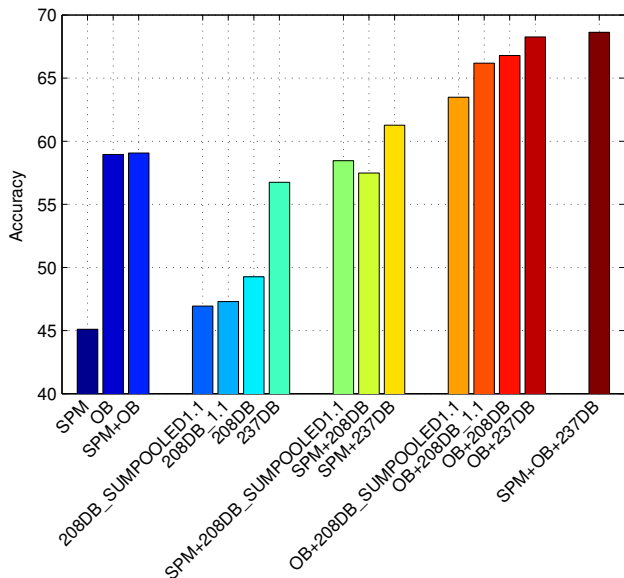
**Table 1: Short description for all the 15 events in the TRECVID MED 2011 dataset.**

to allow for an extensive comparison of the expressiveness and suitability of the proposed representation and study of the model parameters (e.g. number of models, number of thresholds, influence of different detection statistics etc.).

We randomly split the videos into a training (40%), validation (20%), and test (40%) set. Support Vector Machines (SVM) are used to classify concatenated features of different representations as one of the 15 events. We also experimented with Multiple Kernel Learning [11] and Double Fusion [6] which did not lead to a significant increase of performance in our case. In our experience, linear kernel (one-vs-rest) SVM yielded the best performance in almost all cases. We also investigated both mean and max pooling techniques to pool the computed representation across all keyframes of a given video. Generally, max pooling worked best for Object Bank (OB) features while mean pooling worked best for both Spatial Pyramid Matching (SPM) features and Detection Bank (DB) features. Due to space constraints we only report the best classification accuracies per feature representation in Figure 2.

The first three bars represent the baseline accuracies using SPM (45.10%) and OB (58.95%) features (and their combination; 59.07%). Object Bank performs already quite well due to the fact that it captures the likelihood of presence of certain objects that are relevant to the events (e.g. several animals, wheels, different tools, balloons, wedding gowns, bridegrooms etc.). The next four bars show the performance of our proposed representation alone where the number before DB refers to the number of DPM models used in the representation. We used 208 models from Object Bank, 20 models from PASCAL, and 9 additional models specifically trained for the defined task (boat, bread, cake, candle, fish, goat, jeep, scissors, and tire). These event-specific models were learned automatically by choosing categories based on available textual descriptions of the MED Events and acquiring available training data from ImageNet [1]. Using only the  $D_S$  statistic and only one threshold ( $t = -1.1$ ) on all 208 models from Object Bank (208DB\_SUMPOOLED1.1) we obtain 46.94%, using all statistics ( $D_S, D_N, D_0$ ) but only one threshold (208DB\_1.1) 47.30%, using all statistics and four thresholds ( $t \in \{-1.1, -0.9, -0.7, -0.5\}$ ; 208DPM) 49.26%, and using all 237 models, all statistics and four thresholds (237DB) 56.74% accuracy. We have found these features to have only 30% non-zero values (sparsity 70%).

Our representation provides complementary discriminative information to both scene-level (SPM) and window-based object features (OB) by improved classification performance on the combined feature sets. Using only a relatively compact representation (SPM+208DB\_SUMPOOLED1.1)



**Figure 2: Classification accuracies on the TRECVID MED 2011 Event Kit for different feature combinations. The proposed Detection Bank (DB) representation provides complementary information to both scene-level (SPM) and window-based features (OB).**

we can obtain Object Bank-like performance of 58.46% (8568D compared to 44604D). Using all three statistics and four thresholds (SPM+208DB), or additional models (SPM+237DB) we obtain 57.48% and 61.27%, respectively.

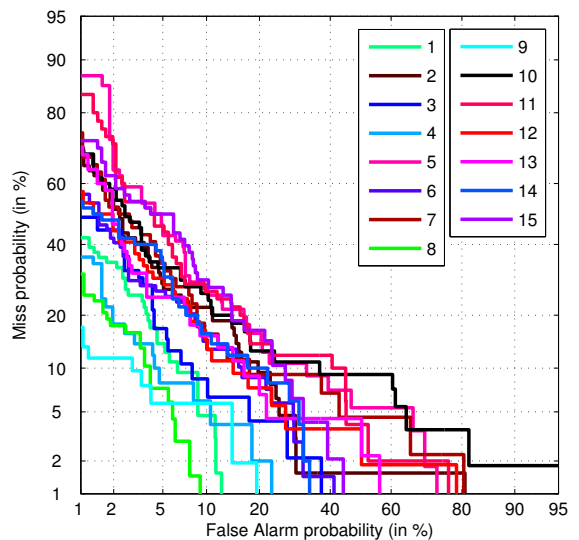
The main result of this paper is that we can, without using more object models, use the proposed representation on top of Object Bank to significantly improve performance. Adding the  $D_S$  statistic for 208 models with only one threshold we obtain 63.48%, using all three statistics 66.18%, using all statistics and four thresholds 66.79%, and using the 29 additional models 68.26% (almost 10% higher than OB alone). If we also add SPM features to the last one, we reach the highest classification accuracy of 68.63% (108528D feature). Figure 3 shows a DET curve for this feature combination for all 15 events.

In addition, we investigated the value of the additional categories. Removing the 9 event-specific categories from the full combination of features leads to a drop of 0.7% and removing all 29 additional categories leads to a drop of 1.2% in classification accuracy.

## 5. CONCLUSIONS

We proposed a feature representation on videos which we call Detection Bank based on the detection images from a large number of windowed object detectors. The Detection Bank representation provides complementary discriminative information to current state-of-the-art image representations such as Spatial Pyramid Matching and Object Bank. We demonstrated that these combined with our Detection Bank representation provide a significant improvement in multimedia event classification on TRECVID MED 2011 data.

As future work, we plan to evaluate the proposed representation in a detection scenario on the full TRECVID MED 2011 data set. Further, we want to investigate the influence of event-specific object categories that could be learned on the fly and to what degree detecting meaningful categories



**Figure 3: DET curves for all 15 events using the all features. Event IDs in order of intersection with horizontal axis are: 8, 1, 9, 4, 14, 3, 6, 15, 13, 5, 11, 12, 7, 2, and 10. The mapping from ID to Event is given in Table 1. Best viewed in color.**

differs from simply having a large number of independent measurements from random filter responses.

## 6. REFERENCES

- [1] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. *CVPR*, 2009.
- [2] M. Everingham, L. Van Gool, C. Williams, J. Winn, and A. Zisserman. The PASCAL visual object classes challenge, VOC 2009 results (2009).
- [3] L. Fei-Fei and P. Perona. A bayesian hierarchical model for learning natural scene categories. *CVPR*, 2005.
- [4] P. F. Felzenszwalb, R. B. Girshick, D. A. McAllester, and D. Ramanan. Object detection with discriminatively trained part-based models. *PAMI*, 2010.
- [5] A. Hauptmann, R. Yan, W. Lin, M. Christel, and H. Wactlar. Can high-level concepts fill the semantic gap in video retrieval? a case study with broadcast news. *IEEE TMM*, 2007.
- [6] Z. Lan, L. Bao, S. Yu, W. Liu, and A. Hauptmann. Double fusion for multimedia event detection. *Advances in Multimedia Modeling*, 2012.
- [7] S. Lazebnik, C. Schmid, and J. Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. *CVPR*, 2006.
- [8] L.-J. Li, H. Su, E. P. Xing, and L. Fei-Fei. Object bank: A high-level image representation for scene classification & semantic feature sparsification. *NIPS*, 2010.
- [9] P. Natarajan, P. Natarajan, V. Manohar, S. Wu, S. Tsakalidis, S. Vitaladevuni, X. Zhuang, R. Prasad, G. Ye, D. Liu, et al. BBN VISER TRECVID 2011 multimedia event detection system. 2011.
- [10] P. Over, G. Awad, J. Fiscus, B. Antonishek, M. Michel, A. Smeaton, W. Kraaij, and G. Quénot. TRECVID 2010—an overview of the goals, tasks, data, evaluation mechanisms, and metrics. 2011.
- [11] S. Sonnenburg, G. Rätsch, C. Schäfer, and B. Schölkopf. Large scale multiple kernel learning. *JMLR*, 2006.
- [12] L. Torresani, M. Szummer, and A. Fitzgibbon. Efficient object category recognition using classemes. *ECCV*, 2010.

ACKNOWLEDGMENT: Supported by the Intelligence Advanced Research Projects Activity (IARPA) via Department of Interior National Business Center contract number D11PC20066. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon. Disclaimer: The views and conclusion contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsement, either expressed or implied, of IARPA, DOI/NBC, or the U.S. Government.