

Unified Embedding and Metric Learning for Zero-Exemplar Event Detection

Noureddien Hussein, Efstratios Gavves, Arnold W.M. Smeulders
 QUVA Lab, University of Amsterdam

{nhussein, egavves, a.w.m.smeulders}@uva.nl

Abstract

Event detection in unconstrained videos is conceived as a content-based video retrieval with two modalities: textual and visual. Given a text describing a novel event, the goal is to rank related videos accordingly. This task is zero-exemplar, no video examples are given to the novel event.

Related works train a bank of concept detectors on external data sources. These detectors predict confidence scores for test videos, which are ranked and retrieved accordingly. In contrast, we learn a joint space in which the visual and textual representations are embedded. The space casts a novel event as a probability of pre-defined events. Also, it learns to measure the distance between an event and its related videos.

Our model is trained end-to-end on publicly available EventNet. When applied to TRECVID Multimedia Event Detection dataset, it outperforms the state-of-the-art by a considerable margin.

1. Introduction

TRECVID Multimedia Event Detection (MED) [1, 2] is a retrieval task for event videos, with the reputation of being realistic. It comes in two flavors: few-exemplar and zero-exemplar, where the latter means that no video example is known to the model. Although expecting a few examples seems reasonable, in practice this implies that the user must already have an index of any possible query, making it very limited. In this paper, we focus on event video search with zero exemplars.

Retrieving videos of never-seen events, such as “renovating home”, without any video exemplar poses several challenges. One challenge is how to bridge the gap between the visual and the textual semantics [3, 4, 5]. One approach [3, 6, 7, 8, 9, 10] is to learn a dictionary of concept detectors on external data source. Then, scores for test videos are predicted using these detectors. Test videos are then ranked and retrieved accordingly. The inherent weak-

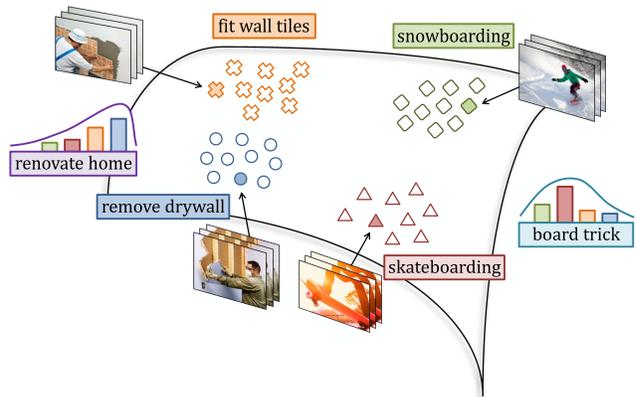


Figure 1. We pose the problem of zero-exemplar event detection as learning from a repository of pre-defined events. Given video exemplars of events “removing drywall” or “fit wall times”, one may detect a novel event “renovate home” as a probability distribution over the predefined events.

ness of this approach is that the presentation of a test video is reduced to a limited vocabulary from the concept dictionary. Another challenge is how to overcome the domain difference between training and test events. While Semantic Query Generation (SQG) [3, 8, 9, 11] mitigates this challenge by extracting keywords from the event query, it does not address how relevant these keywords to the event itself. For example, keyword “person” is not relevant to event “car repair” as it is to “flash mob gathering”.

Our entry to zero-exemplar events is that they generally have strong semantic correlations [12, 13] with other possibly seen events. For instance, the novel event “renovating home” is related to “fit wall tiles”, “remove drywall”, or even to “paint door”. Novel events can, therefore, be casted on a repository of prior events, for which knowledge sources in various forms are available beforehand, such as the videos, as in EventNet [14], or articles, as in WikiHow [15]. Not only do these sources provide video examples of a large –but still limited– set of events, but also they provide an association of text description of events with their corresponding videos. A text article can describe the event in words: what is it about, what are the

details and what are the semantics. We note that such a visual-textual repository of events may serve as a knowledge source, by which we can interpret novel event queries.

For Zero-exemplar Event Detection (ZED), we propose a neural model with the following novelties:

1. We formulate a unified embedding for multiple modalities (e.g. visual and textual) that enables a contrastive metric for maximum discrimination between events.
2. A textual embedding poses the representation of a novel event as a probability of predefined events, such that it spans a much larger space of admissible expressions.
3. We exploit a single data source, comprising pairs of event articles and related videos. A single source rather enables end-to-end learning from multi-modal individual pairs.

We empirically shows that our novelties result in performance improvement. We evaluate the model on TRECVID Multimedia Event Detection (MED) 2013 [1] and 2014 [2]. Our results show significant improvement over the state-of-the-art.

2. Related Work

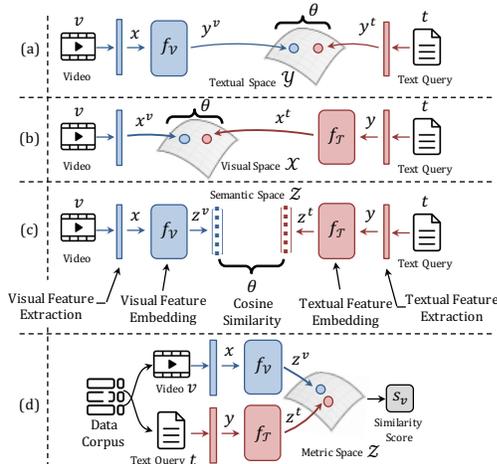


Figure 2. Three families of methods for zero-exemplar event detection: (a), (b) and (c). They build on top of feature representations learned a priori (i.e. initial representations), such as CNN features x for a video v or word2vec features y for event text query t . In a post-processing step, the distance θ is measured between the embedded features. In contrast, our model rather falls in a new family, depicted in (d), for it learns unified embedding with metric loss using single data source.

We identify three families of methods for ZED, as in figure 2 (a), (b) and (c).

Visual Embedding and Textual Retrieval. As in figure 2(a), given a video v_i represented as $x \in \mathcal{X}$ and a related text t represented as $y \in \mathcal{Y}$. Then, a visual model f_v is trained to project x as $y^v \in \mathcal{Y}$ such that the distance is minimized between (y^v, y) . In test time, video ranking and

retrieval is done using distance metric between the projected test video y^t and test query representation y .

[16, 17] project the visual feature x of a web video v into term-vector representation y of the video’s textual title t . However, during training, the model makes use of the text query of the test events to learn better term-vector representation. Consequently, this limits the generalization for novel event queries.

Textual Embedding and Visual Retrieval. As in figure 2(b), a given text query t is projected into $x^t \in \mathcal{X}$ using pre-trained or learned language model f_T .

[18] makes use of freely-available weekly-tagged web videos. Then it propagates tags to test videos from its nearest neighbors. Methods [7, 8, 9, 10, 3] have similar approach. Given a text query t , Semantic Query Generation (SQG) extracts N most related concepts $\{c_i, i \in N\}$ to the test query. Then, pre-trained concept detectors predict probability scores $\{s_i, i \in N\}$ for a test video v . Aggregating these probabilities results in the final video score s_v , upon which videos are ranked and retrieved. [9] learns weighted averaging.

The shortcoming of this family is that expressing a video as probability scores of few concepts is under-representation. Any concept that exists in the video but is missing in the concept dictionary is thus unrepresented.

Visual-Textual Embedding and Semantic Retrieval. As in figure 2(c), visual f_v and textual f_T models are trained to project both of the visual x and textual y features into a semantic space \mathcal{Z} . During test, ranking score is the distance between the projections z^v, z^t in the semantic space \mathcal{Z} .

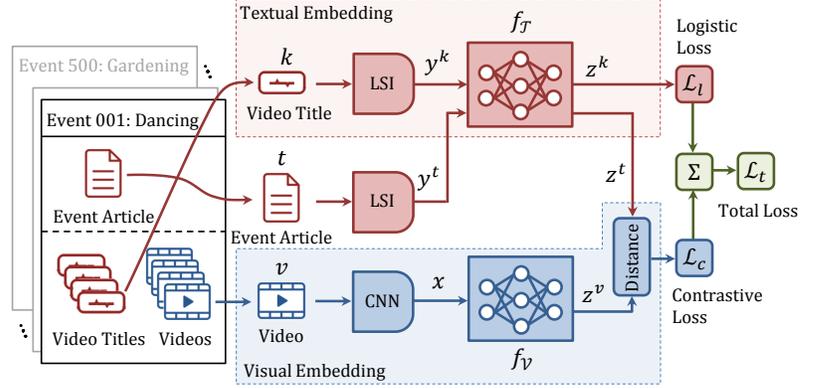
[19] projects video concepts into a high-dimensional lexicon space. Separately, it projects concept-based features to the space, which overcomes the lexicon mismatch between the query and the video concepts. [20] embeds a fusion of low and mid-level visual features into distributional semantic manifold [21, 22]. In a separate step, it embeds text-based concepts into the manifold.

The third family, see figure 2(c), is superior to the others, see figure 2(a), (b). However, one drawback of [19, 20] is separately embedding both the visual and textual features z^v, z^t . This leads to another drawback, having to measure the distance between (z^v, z^t) in a post-processing step (e.g. cosine similarity).

Unified Embedding and Metric Learning Retrieval Our method rather falls into a new family, see figure 2(d), and it overcomes the shortcomings of [19, 20] by the following. It is trained on a single data source, enabling a unified embedding for features of multiple modalities into a metric space. Consequently, the distance between the embedded features is measured by the model using the learned metric space.

Auxiliary Methods Independent to the previous works, the following techniques have been used to improve the results:

Figure 3. Model overview. Using dataset \mathcal{D}^z of M event categories and N videos. Each event has a text article and a few videos. Given a video x with text title k , belonging to an event with article t , we **extract** features x, y^k, y^t respectively. At the top, network $f_{\mathcal{T}}$ **learns** to classify the title feature y^k into one of M event categories. In the middle, we borrow the network $f_{\mathcal{T}}$ to **embed** the event article’s feature y^t as $z^t \in \mathcal{Z}$. Then, at the bottom, the network $f_{\mathcal{V}}$ **learns** to embed the video feature x as $z^v \in \mathcal{Z}$ such that the distance between (z^v, z^t) is minimized, in the learned metric space \mathcal{Z} .



self-paced reranking [23], pseudo-relevance feedback [24], event query manual intervention [25], early fusion of features (action [26, 27, 28, 29, 30] or acoustic [31, 32, 33]) or late fusion of concept scores [17]. All these contributions may be applied to our method.

Visual Representation. ConvNets [34, 35, 36, 37] provide frame-level representation. To tame them into video-level counterpart, literature use: i- frame-level filtering [38] ii- vector encoding [39, 40] iii- learned pooling and re-counting [10, 41] iv- average pooling [16, 17]. Also, low-level action [28, 29], mid-level action [26, 27] or acoustic [31, 32, 33] features can be used. **Textual Representation.** To represent text, literature use: i- sequential models [42] ii- continuous word-space representations [22, 43] iii- topic models [44, 45] iv- dictionary-space representation [17].

3. Method

3.1. Overview

Our goal is zero-exemplar retrieval of event videos with respect to their relevance to a novel textual description of an event. More specifically, for the zero-exemplar video dataset $\mathcal{D}^z = \{v_i^z\}, i = 1, \dots, L$ and given any future, textual event description t^z , we want to learn a model $f(\cdot)$ that ranks the videos v_i^z according to the relevance to t^z , namely:

$$t^z : v_i^z \succ v_j^z \rightarrow f(v_i^z, t^z) > f(v_j^z, t^z). \quad (1)$$

3.2. Model

Since we focus on zero-exemplar setting, we cannot expect any training data directly relevant to the test queries. As such, we cannot directly optimize our model for the parameters $W_{\mathcal{T}}, W_{\mathcal{V}}$ in eq. (3). In the absence of any direct data, we resort to external knowledge databases. More specifically, we propose to cast future novel query descriptions as a convex combination of known query descriptions in external databases, where we can measure their relevance to the database videos.

We start from a dataset $\mathcal{D}^z = \{v_i, k_i, l_j, t_j\}, i = 1, \dots, N, j = 1, \dots, M$ organized by an event taxonomy, where we do not expect nor require the events to overlap with any future event queries. The dataset is composed of M events. Each event is associated with a textual, article description of the event, analyzing different aspects of it, such as: (i) the typical appearance of subjects and objects (ii) it’s procedures (iii) the steps towards completing task associated with it. The dataset contains in total N videos, with v_i denoting the i -th video in the dataset with metadata k_i , e.g. the title of the video. A video is associated with an event label l_i and the article description t_i of the event it belongs to. Since multiple videos belong to the same event, they share the article description of such event.

The ultimate goal of our model is zero-exemplar search for event videos. Namely, provided unknown text queries by the user, we want to retrieve those videos that are relevant. We illustrate our proposed model during training in figure 3. The model is composed of two components, a textual embedding $f_{\mathcal{T}}(\cdot)$, a visual embedding $f_{\mathcal{V}}(\cdot)$. Our ultimate goal is the ranking of videos, $v_i \succ v_j \succ v_k$ with respect to their relevance to a query description, or in pairwise terms $v_i \succ v_j, v_j \succ v_k$ and $v_i \succ v_k$.

Let us assume a pair of videos v_i, v_j and query description t , where video v_i is more relevant to the query t than v_j . Our goal is a model that learns to put videos in the correct relative order, namely $(v_i, t) \succ (v_j, t)$. This is equivalent to a model that learns visual-textual embeddings such that $d_i^{tv} < d_j^{tv}$, where d_i^{tv} is the distance between visual-textual embeddings of (v_i, t) , d_j^{tv} is the same for (v_j, t) . Since we want to compare distances between pairs $(v_i, t), (v_j, t)$, we pose the learning of our model as the minimization of a contrastive loss [46]:

$$\mathcal{L}_{con} = \frac{1}{2N} \sum_{i=1}^N h_i \cdot d_i^2 + (1 - h_i) \max(1 - d_i, 0)^2, \quad (2)$$

$$d_i = \|f_{\mathcal{T}}(t_i; W_{\mathcal{T}}) - f_{\mathcal{V}}(v_i; W_{\mathcal{V}})\|_2, \quad (3)$$

where $f_{\mathcal{T}}(t_i; W_{\mathcal{T}})$ is the projection of the query description t_i into the **unified** metric space \mathcal{Z} parameterized by

$W_{\mathcal{T}}$, $f_{\mathcal{V}}(v_i; W_{\mathcal{V}})$ is the projection of a video v_i onto the same space \mathcal{Z} parameterized by $W_{\mathcal{V}}$ and h_i a target variable that equals to 1 when the i -th video is relevant to the query description t_i and 0 otherwise. Naturally, to optimize eq. (2), we first need to define the projections $f_{\mathcal{T}}(\cdot; W_{\mathcal{T}})$ and $f_{\mathcal{V}}(\cdot; W_{\mathcal{V}})$ in eq. (3).

Textual Embedding. The textual embedding component of our model, $f_{\mathcal{T}}(\cdot; W_{\mathcal{T}})$, is illustrated in figure 3 (top). This component is dedicated to learn a projection of a textual input –including any future event queries t – on to the unified space \mathcal{Z} . Before detailing our model $f_{\mathcal{T}}$, however, we note that that the textual embedding can be employed not only with event article descriptions, but also with any other textual information that might be associated to the dataset videos, such as textual metadata. Although we expect the video title not to be as descriptive as the associated article, they may still be able to offer some discriminative information as previously shown [16, 17] which can be associated to the event category.

We model the textual embedding as a shallow (two layers) multi-layer perceptron (MLP). For the first layer we employ a ReLU nonlinearity. The second layer serves a dual purpose. First, it projects the article description of an event on the unified space \mathcal{Z} . This projection is *category-specific*, namely different videos that belong to the same event will share the projection. Second, it can project any *video-specific* textual metadata into the unified space. We, therefore, propose to embed the title metadata k_i , which is uniquely associated with a video, not an event category. To this end, we opt for softmax nonlinearity for the second layer, followed by an additional logistic loss term to penalized misprediction of titles m_i with respect to the video’s event label y_i^j , namely

$$\mathcal{L}_{log} = \sum_{i=1}^N \sum_{j=1}^M -y_i^j \log f_{\mathcal{T}}^j(k_i; W_{\mathcal{T}}). \quad (4)$$

Overall, the textual embedding $f_{\mathcal{T}}^j$ is trained with a dual loss in mind. The first loss term, see eq. (2) (3) takes care that the final network learns event-relevant textual projections. The second loss term, see eq. (4), takes care that the final network does not overfit to the particular event article descriptions. The latter is crucial because the event article descriptions in \mathcal{D}^z will not overlap with the future event queries, since we are in a zero-exemplar retrieval setting. As such, training the textual embedding to be optimal only for these event descriptions will likely result in severe overfitting. Our goal and hope is that the final textual embedding model $f_{\mathcal{T}}$ will capture both event-aware and video-discriminative textual features.

Visual Embedding. The visual embedding component of our model, $f_{\mathcal{V}}(\cdot; W_{\mathcal{V}})$, is illustrated in figure 3 (bottom).

This component is dedicated to learn a projection from the visual input, namely the videos in our zero-exemplar dataset \mathcal{D}^z , into the unified metric space \mathcal{Z} . The goal is to project the videos belonging to semantically similar events; project them into a similar region in the space. We model the visual embedding $f_{\mathcal{V}}(v_i; W_{\mathcal{V}})$ using a shallow (two layers) multi-layer perceptron with tanh nonlinearities, applied to any visual feature for video v_i .

End-to-End Training. At each training forward-pass, the model is given a triplet of data inputs, an event description t_i , a related video v_i and video title k_i . From eq. (3) we observe that the visual embedding $f_{\mathcal{V}}(v_i; W_{\mathcal{V}})$ is encouraged to minimize its distance with the output of the textual embedding $f_{\mathcal{T}}(t_i; W_{\mathcal{T}})$. In the end, all the modules of the proposed model are differentiable. Therefore, we train our model in an end-to-end manner by minimizing the following objective

$$\begin{aligned} \arg \min_{W_{\mathcal{V}}, W_{\mathcal{T}}} \mathcal{L}^{\mathcal{U}}, \\ \mathcal{L}^{\mathcal{U}} = \mathcal{L}_{con} + \mathcal{L}_{log}. \end{aligned} \quad (5)$$

For the triplet input (v_i, t_i, k_i) , we rely on external representations, since our ultimate goal is zero-exemplar search. Strictly speaking, a visual input v_i is represented as CNN [35] feature vector, while textual inputs t_i, k_i are represented as LSI [45] or Doc2Vec [43] feature vectors. However, given that these external representations rely on neural network architectures, if needed, they could also be further fine-tuned. We choose to freeze CNN and Doc2Vec modules to speed up training. Finally, in this paper we refer to our main model with unified embedding, as **model^U**.

Inference. After training, we fix the parameters $(W_{\mathcal{V}}, W_{\mathcal{T}})$. At test time, we set our function $f(\cdot)$ from eq. (1) to be equivalent to the distance function from eq. (??). Hence, at test time, we compute the Euclidean distance in the learned metric space \mathcal{Z} between the embeddings (z^v, z^t) of test video v and novel event description t , respectively.

4. Experiments

4.1. Datasets

Before delving into the details of our experiments, first we describe the external knowledge sources we use.

Training dataset. We leverage videos and articles from publicly available datasets. EvenNet [14] is a dataset of ~90k event videos, harvested from YouTube and categorized into 500 events in hierarchical form according to the events’ ontology. Each event category contains around 180 videos. Each video is coupled with a text title, few tags and related event’s ontology.

We exploit the fact that all events in EventNet are harvested from WikiHow [15] – a website for *How-To* articles covering a wide spectrum of human activities. For instance: “How to Feed a Dog” or “How to Arrange Flowers”. Thus, we crawl WikiHow to get the articles related to all the events in EventNet.

Test dataset. As the task is zero-exemplar, the test sets are different from the training. While EventNet serves as the training, the following serve as the test: TRECVID MED-13 [1] and MED-14 [1]. In details, they are datasets of videos for events. They comprise 27k videos. There are two versions, MED-13 and MED-14 with 20 events for each. Since 10 events overlap, the result is 30 different events in total. Each event is coupled with short textual description (title and definition).

4.2. Implementation Details

Video Features. To represent a video v , we uniformly sample a frame every one second. Then, using ResNet [35], we extract `pool5` CNN features for the sampled frames. Then, we average pool the frame-level features to get the video-level feature x^v . We experiment different features from different CNN models: ResNet (`prob`, `fc1000`), VGG [37] (`fc6`, `fc7`), GoogLeNet [47] (`pool5`, `fc1024`), and Places365 [48] (`fc6`, `fc7`, `fc8`) except we find ResNet `pool5` to be the best. We only use ResNet `pool5` and we don’t fuse multiple CNN features.

Text Features. We choose topic modeling [44, 45], as it is well-suited for long (and sometimes noisy) text articles. We train LSI topic model [45] on Wikipedia corpus [49]. We experiment different latent topics ranging from 300 to 6000, expect we found 2500 to be the best. Also, we experiment other textual representations as LDA [44], SkipThoughts [50] and Doc2Vec [43]. To extract a feature from an event article k or video title t , first we preprocess the text using standard MLP steps: tokenization, lemmatization and stemming. Then, for k, t we extract 2500-D LSI features y^k, y^t , respectively. The same steps apply to MED text queries.

Model Details. Our visual and textual embeddings $f_V(\cdot), f_T(\cdot)$ are learned on top of the aforementioned visual and textual features (x^v, y^k, y^t) . $f_T(\cdot)$ is a 1-hidden layer MLP classifier with ReLU for hidden, softmax for output, logistic loss and 2500–2500–500 neurons for the input, hidden, and output layers, respectively. Similarly, $f_V(\cdot)$ is a 1-hidden layer MLP regressor with ReLU for hidden, contrastive loss and 2048–2048–500 neurons for the input, hidden, and output layers, respectively. Our code is made public¹ to support further research.

¹github.com/noureldien/unified_embedding

4.3. Textual Embedding

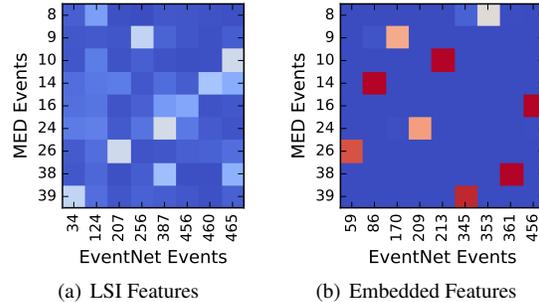


Figure 4. Our textual embedding (b) maps MED to EventNet events better than LSI features. Each dot in the matrix shows the similarity between MED and EventNet events.

Here, we qualitatively demonstrate the benefit of the textual embedding $f_T(\cdot)$. Figure 4 shows the similarity matrix between MED and EventNet events. Each dot represents how a MED event is similar to EventNet events. It shows that our embedding (right) is better than LSI (left) in mapping MED to EventNet events. For example, LSI wrongly maps “9: getting a vehicle unstuck” to “256: launch a boat” while our embedding correctly maps it to “170: drive a car”. Also, our embedding maps with higher confidence than LSI, as in “16: doing homework or study”.

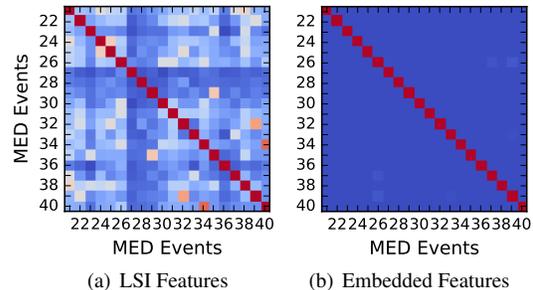


Figure 5. For 20 events of MED-14, our textual embedding (right) is more discriminant than the LSI feature representation (left). Each dot in the matrix shows how similar an event to all the others.

Figure 5 shows the similarity matrix for MED events, where each dot represents how related any MED event to all the others. Our textual embedding (right) is more discriminant than on the LSI feature representation (left). For example, LSI representation shows high semantic correlation between events “34: fixing musical instrument” and “40: tuning musical instrument”, while our embedding discriminate them.

Next, we quantitatively demonstrate the benefit of the textual embedding $f_T(\cdot)$. In contrast to the main model, see section 3, we investigate baseline **model^V**, where we discard the textual embedding $f_T(\cdot)$ and consider only the visual embedding $f_V(\cdot)$. We project a video v on the LSI

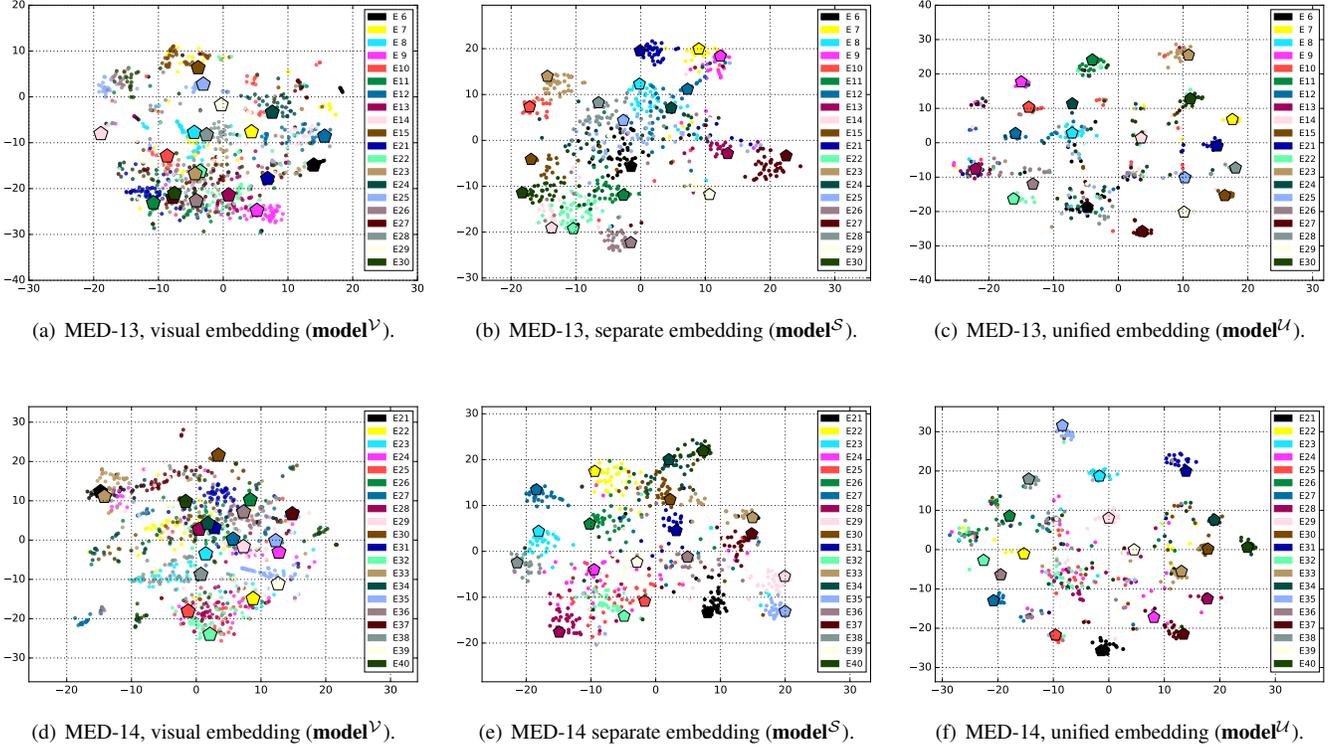


Figure 6. We visualize the results of video embedding using the unified embedding model^U and baselines model^V , model^S . Each sub-figure shows how discriminant the representation of the embedded videos. Each dot represents a projected video, while each pentagon-shape represents a projected event description. We use t-SNE to visualize the result.

representation y of the related event t . Thus, this baseline falls in the first family of methods, see figure 2(a). It is optimized using mean-squared error (MSE) loss \mathcal{L}_{mse}^V , see eq. 6. The result of this baseline is reported in section 5, table 1.

$$\mathcal{L}_{mse}^V = \frac{1}{N} \sum_{i=1}^N \|y_i - f_V(v_i; W_V)\|_2^2. \quad (6)$$

Also, we train another baseline model^C , which is similar to the aforementioned model^V except instead of using MSE loss \mathcal{L}_{mse}^V , see eq. (6), it uses contrastive loss \mathcal{L}_{con}^C , as follows:

$$\mathcal{L}_{con}^C = \frac{1}{2N} \sum_{i=1}^N h_i \cdot d_i^2 + (1 - h_i) \max(1 - d_i, 0)^2, \quad (7)$$

$$d_i = \|y_i - f_V(v_i; W_V)\|_2.$$

4.4. Unified Embedding and Metric Learning

In this experiment, we demonstrate the benefit of the unified embedding. In contrast to our model presented in section 3, we investigate baseline model^S , where this baseline does not learn joint embedding. Instead, it separately learns visual $f_V(\cdot)$ and textual $f_T(\cdot)$ projections. We model these

projections as a shallow (2-layer) MLP trained to classify the data input into 500 event categories, using logistic loss, same as eq. (4).

We conduct another experiment to demonstrate the benefit of learning metric space. In contrast to our model presented in section 3, we investigate baseline model^N , where we discard the metric learning layer. Consequently, this baseline learns the visual embedding is a shallow (2 layers) multi-layer perceptron with \tanh non linearities. Also, we replace the contrastive loss \mathcal{L}_c , see eq. (2) with mean-squared error loss \mathcal{L}_{mse} , namely

$$\mathcal{L}_{mse}^N = \frac{1}{N} \sum_{i=1}^N \|f_T(t_i; W_T) - f_V(v_i; W_V)\|_2^2. \quad (8)$$

During retrieval, this baseline embeds a test video v_i and novel text query t_i as features z^v, z^t onto the common space \mathcal{Z} using textual and visual embeddings $f_T(\cdot), f_V(\cdot)$, respectively. However, in a post-processing step, retrieval score s_i for the video v_i is the cosine distance between (z^v, z^t) . Similarly, all test videos are scored, ranked and retrieved. The results of the aforementioned baselines model^S and model^N are reported in table 1.

Comparing Different Embeddings. In the previous ex-

periments, we investigated several baselines of the unified embedding (model^U), namely visual-only embedding (model^V), separate visual-textual embedding (model^S) and non-metric visual-textual embedding (model^N). In a qualitative manner, we compare the results of such embeddings. As shown in figure 6, we use these baselines to embed event videos of MED-13 and MED-14 datasets into the corresponding spaces. At the same time, we project the textual description of the events on the same space. Then, we use t-SNE [51] to visualize the result on 2D manifold. As seen, the unified embedding, see sub-figures 6(e), 6(f) learns more discriminant representations than the other baselines, see sub-figures 6(a), 6(b), 6(c) and 6(d). The same observation holds for both for MED-13 and MED-14 datasets.

4.5. Mitigating Noise in EventNet

Based of quantitative and qualitative analysis, we conclude that EventNet is noisy. Not only videos are unconstrained, but also some of the video samples are irrelevant to their event categories. EvenNet dataset [14] is accompanied by 500-category CNN classifier. It achieves top-1 and top-5 accuracies of 30.67% and 53.27%, respectively. Since events in EventNet are structured as an ontological hierarchy, there is a total of 19 high-level categories. The classifier achieves top-1 and top-5 accuracies of 38.91% and 57.67%, respectively, over these high-level categories.

Based on these observations, we prune EventNet to remove noisy videos. To this end, first we represent each video as average pooling of ResNet pool5 features. Then, we follow the conventional 5-fold cross validation with 5 rounds. For each round, we split the dataset into 5 subsets, 4 subsets \mathcal{V}_t for training and the last \mathcal{V}_p for pruning. Then we train a 2-layer MLP for classification. After training, we forward-pass the videos of \mathcal{V}_p and rule-out the misclassified ones.

The intuition behind pruning is that we rather learn salient event concepts using less video samples than learn noisy concepts with more samples. Pruning reduced the total number of videos by 26%, from 90.2k to 66.7k. This pruned dataset is all what we use in our experiments.

4.6. Latent Topics in LSI

When training LSI topic model on Wikipedia corpus, a crucial parameter is the number of latent topics K the model constructs. We observe improvements in the performance directly proportional to increasing K . The main reason that the bigger the value of K , the more discriminant the LSI feature is. Figure 7 confirms our understanding.

5. Results

Evaluation metric. Since we are addressing, in essence, an information retrieval task, we rely on the average precision (AP) per event, and mean average precision (mAP) per

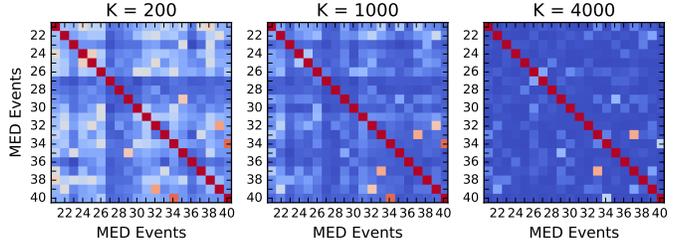


Figure 7. Similarity matrix between LSI features of MED-14 events. The more the latent topics (K) in LSI model, the higher the feature dimension, and the more discriminant the feature.

dataset. We follow the standard evaluation method as in the relevant literature [1, 2, 52].

Comparing against model baselines. In table 1, we report the mAP score of our model baselines, previously discussed in the experiments, see section 4. The table clearly shows the marginal contribution of each of novelty for the proposed method.

Baseline	Loss	Metric	$f_V(\cdot)$	$f_T(\cdot)$	MED13	MED14
model ^V	\mathcal{L}_{mse}^V (6)	✗	✓	✗	11.90	10.76
model ^C	\mathcal{L}_{con}^C (7)	✓	✓	✗	13.29	12.31
model ^S	\mathcal{L}_{log}^S (4)	✗	✓	✓	15.60	13.49
model ^N	\mathcal{L}_{mse}^N (8)	✗	✓	✓	15.92	14.36
model ^U	\mathcal{L}^U (5)	✓	✓	✓	17.86	16.67

Table 1. Comparison between the unified embedding and other baselines. The unified embedding model^U achieves the best results on MED-13 and MED-14 datasets.

Comparing against related work. We report the performance of our method, the unified embedding model^U on TRECVID MED-13 and MED-14 datasets. When compared with the related works, our method improves over the state-of-the-art by a considerable margin, as shown in table 2 and figure 8.

Method		MED13	MED14
TagBook [18]	ToM '15	12.90	05.90
Discovery [7]	ICAI '15	09.60	–
Composition [8]	AAAI '16	12.64	13.37
Classifiers [9]	CVPR '16	13.46	14.32
VideoStory [†] [17]	PAMI '16	15.90	05.20
VideoStory* [17]	PAMI '16	20.00	08.00
This Paper (model ^U)		17.86	16.67

Table 2. Performance comparison between our model and related works. We report the mean average precision (mAP%) for MED-13 and MED-14 datasets.

It is important to point out that VideoStory[†] uses only object feature representation, so its comparable to our method.

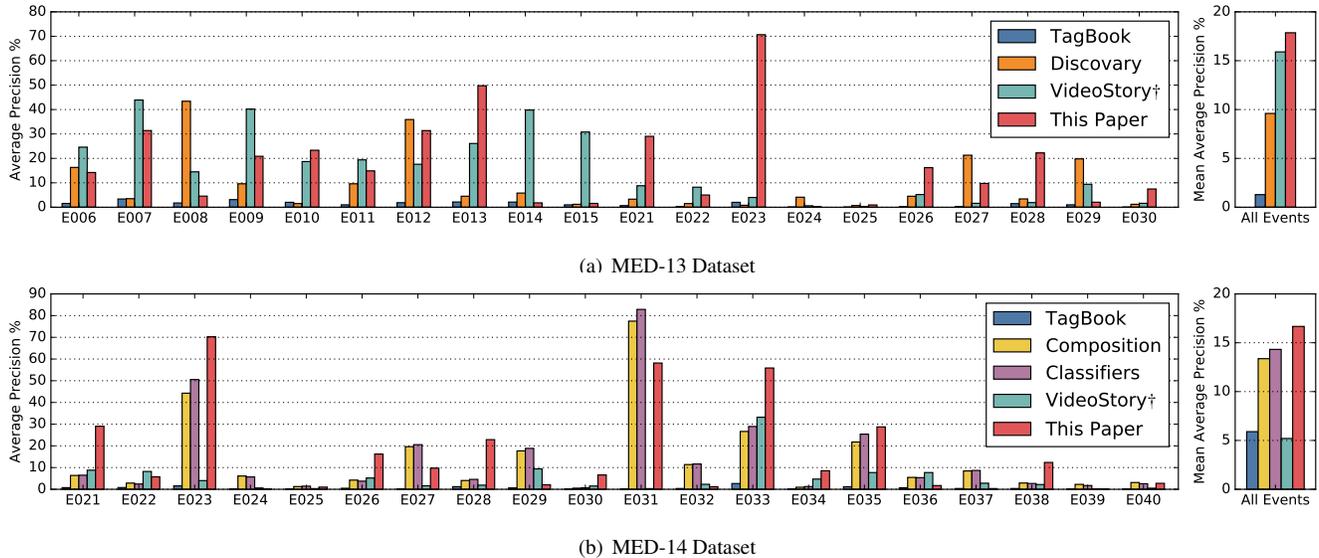


Figure 8. Event detection accuracies: per-event average precision (AP%) and per-dataset mean average precision (mAP%) for MED-13 and MED-14 datasets. We compare our results against TagBook [18], Discovery [7], Composition [8], Classifiers [9] and VideoStory [17].

However, VideoStory* uses motion feature representation and expert text query (i.e. using term-importance matrix H in [17]). To rule out the marginal effect of using different datasets and features, we train VideoStory and report results in table 3. Clearly, CNN features and video exemplars in the training set can improve the model accuracy, but our method improves against VideoStory when trained on the same dataset and using the same features. Other works (Classifiers [9], Composition [8]) use both image and action concept classifiers. Nonetheless, our method improves over them using only object-centric CNN feature representations.

Method	Training Set	CNN Feat.	MED14
VideoStory	VideoStory46k [17]	GoogleNet	08.00
VideoStory	FCVID [53]	GoogleNet	11.84
VideoStory	EventNet [14]	GoogleNet	14.52
VideoStory	EventNet [14]	ResNet	15.80
This Paper	EventNet [14]	ResNet	16.67

Table 3. Our method improves over VideoStory when trained on the same dataset and using the same feature representation.

6. Conclusion

In this paper, we presented a novel approach for detecting events in unconstrained web videos, in a zero-exemplar fashion. Rather than learning separate embeddings from cross-modal datasets, we proposed a unified embedding where several cross-modalities are jointly projected. This enables end-to-end learning. On top of this, we exploited the fact that zero-exemplar is posed as retrieval task and

proposed to learn metric space. This enables measuring the similarities between the embedded modalities using this very space.

We experimented the novelties and demonstrated how they contribute to improving the performance. We complemented this by improvements over the state-of-the-art by considerable margin on MED-13 and MED-14 datasets.

However, the question still remains, how can we discriminate between these two MED events “34: fixing musical instrument” and “40: tuning musical instrument”. We would like to argue that temporal modeling for human actions in videos is of absolute necessity to achieve such fine-grained event recognition. In future research, we would like to focus on human-object interaction in videos and how to model it temporally.

Acknowledgment

We thank Dennis Koelma, Masoud Mazloom and Cees Snoek² for lending their insights and technical support for this work.

References

- [1] Paul Over, George Awad, Jon Fiscus, Greg Sanders, and Barbara Shaw. Trecvid 2013—an introduction to the goals, tasks, data, evaluation mechanisms, and metrics. In *TRECVID Workshop*, 2013. 1, 2, 5, 7
- [2] Paul Over, Jon Fiscus, Greg Sanders, David Joy, Martial Michel, George Awad, Alan Smeaton, Wessel Kraaij, and Georges Quénot. Trecvid 2014—an overview of the goals, tasks, data, evaluation mechanisms and metrics. In *TRECVID Workshop*, 2014. 1, 2, 7

²{kolema, m.mazloom, cgmsnoek}@uva.nl

- [3] Lu Jiang, Shou-I Yu, Deyu Meng, Teruko Mitamura, and Alexander G Hauptmann. Bridging the ultimate semantic gap: A semantic search engine for internet videos. In *ICMR*, 2015. 1, 2
- [4] Amirhossein Habibian, Thomas Mensink, and Cees GM Snoek. Composite concept discovery for zero-shot video event detection. In *ICMR*, 2014. 1
- [5] Amirhossein Habibian, Thomas Mensink, and Cees GM Snoek. Discovering semantic vocabularies for cross-media retrieval. In *ICMR*, 2015. 1
- [6] Masoud Mazloom, Efstratios Gavves, and Cees G. M. Snoek. Conceptlets: Selective semantics for classifying video events. In *IEEE TMM*, 2014. 1
- [7] Xiaojun Chang, Yi Yang, Alexander G Hauptmann, Eric P Xing, and Yao-Liang Yu. Semantic concept discovery for large-scale zero-shot event detection. In *IJCAI*, 2015. 1, 2, 7, 8
- [8] Xiaojun Chang, Yi Yang, Guodong Long, Chengqi Zhang, and Alexander G Hauptmann. Dynamic concept composition for zero-example event detection. In *arXiv*, 2016. 1, 2, 7, 8
- [9] Xiaojun Chang, Yao-Liang Yu, Yi Yang, and Eric P Xing. They are not equally reliable: Semantic event search using differentiated concept classifiers. In *IEEE CVPR*, 2016. 1, 2, 7, 8
- [10] Yi-Jie Lu. Zero-example multimedia event detection and recounting with unsupervised evidence localization. In *ACM MM*, 2016. 1, 2, 3
- [11] Lu Jiang, Shou-I Yu, Deyu Meng, Yi Yang, Teruko Mitamura, and Alexander G Hauptmann. Fast and accurate content-based semantic search in 100m internet videos. In *ACM MM*, 2015. 1
- [12] Thomas Mensink, Efstratios Gavves, and Cees Snoek. Costa: Co-occurrence statistics for zero-shot classification. In *IEEE CVPR*, 2014. 1
- [13] E. Gavves, T. E. J. Mensink, T. Tommasi, C. G. M. Snoek, and T. Tuytelaars. Active transfer learning with zero-shot priors: Reusing past datasets for future tasks. In *IEEE ICCV*, 2015. 1
- [14] Guangnan Ye, Yitong Li, Hongliang Xu, Dong Liu, and Shih-Fu Chang. Eventnet: A large scale structured concept library for complex event detection in video. In *ACM MM*, 2015. 1, 4, 7, 8
- [15] Wikihow. <http://wikihow.com>. 1, 5
- [16] Amirhossein Habibian, Thomas Mensink, and Cees GM Snoek. Videostory: A new multimedia embedding for few-example recognition and translation of events. In *ACM MM*, 2014. 2, 3, 4
- [17] Amirhossein Habibian, Thomas Mensink, and Cees GM Snoek. Videostory embeddings recognize events when examples are scarce. In *IEEE TPAMI*, 2016. 2, 3, 4, 7, 8
- [18] Masoud Mazloom, Xirong Li, and Cees Snoek. Tagbook: A semantic video representation without supervision for event detection. In *IEEE TMM*, 2015. 2, 7, 8
- [19] Shuang Wu, Sravanthi Bondugula, Florian Luisier, Xiaodan Zhuang, and Pradeep Natarajan. Zero-shot event detection using multi-modal fusion of weakly supervised concepts. In *IEEE CVPR*, 2014. 2
- [20] Mohamed Elhoseiny, Jingen Liu, Hui Cheng, Harpreet Sawhney, and Ahmed Elgammal. Zero-shot event detection by multimodal distributional semantic embedding of videos. In *arXiv*, 2015. 2
- [21] Tomas Mikolov, Quoc V Le, and Ilya Sutskever. Exploiting similarities among languages for machine translation. In *arXiv*, 2013. 2
- [22] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *NIPS*, 2013. 2, 3
- [23] Lu Jiang, Deyu Meng, Teruko Mitamura, and Alexander G Hauptmann. Easy samples first: Self-paced reranking for zero-example multimedia search. In *ACM MM*, 2014. 3
- [24] Lu Jiang, Teruko Mitamura, Shou-I Yu, and Alexander G Hauptmann. Zero-example event search using multimodal pseudo relevance feedback. In *ICMR*, 2014. 3
- [25] Arnav Agharwal, Rama Kovvuri, Ram Nevatia, and Cees GM Snoek. Tag-based video retrieval by embedding semantic content in a continuous word space. In *IEEE WACV*, 2016. 3
- [26] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. C3d: generic features for video analysis. In *arXiv*, 2014. 3
- [27] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In *NIPS*, 2014. 3
- [28] Heng Wang, Alexander Kläser, Cordelia Schmid, and Cheng-Lin Liu. Action recognition by dense trajectories. In *IEEE CVPR*, 2011. 3
- [29] J Uijlings, IC Duta, Enver Sangineto, and Nicu Sebe. Video classification with densely extracted hog/hof/mbh features: an evaluation of the accuracy/computational efficiency trade-off. In *IJMR*, 2015. 3
- [30] B. Fernando, E. Gavves, J. Oramas, A. Ghodrati, and T. Tuytelaars. Rank pooling for action recognition. In *IEEE TPAMI*, 2016. 3
- [31] Lindaswa Muda, Mumtaj Begam, and I Elamvazuthi. Voice recognition algorithms using mel frequency cepstral coefficient (mfcc) and dynamic time warping (dtw) techniques. In *arXiv*, 2010. 3
- [32] Anurag Kumar and Bhiksha Raj. Audio event detection using weakly labeled data. In *arXiv*, 2016. 3
- [33] Liping Jing, Bo Liu, Jaeyoung Choi, Adam Janin, Julia Bernd, Michael W Mahoney, and Gerald Friedland. A discriminative and compact audio representation for event detection. In *ACM MM*, 2016. 3
- [34] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. In *IEEE*. 3

- [35] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *arXiv*, 2015. 3, 4, 5
- [36] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 3
- [37] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *arXiv*, 2014. 3, 5
- [38] Chuang Gan, Ting Yao, Kuiyuan Yang, Yi Yang, and Tao Mei. You lead, we exceed: Labor-free video concept learning by jointly exploiting web videos and images. In *IEEE CVPR*, 2016. 3
- [39] Karen Simonyan, Omkar M Parkhi, Andrea Vedaldi, and Andrew Zisserman. Fisher vector faces in the wild. In *BMVC*, 2013. 3
- [40] Relja Arandjelovic and Andrew Zisserman. All about vlad. In *IEEE CVPR*, 2013. 3
- [41] Pascal Mettes, Jan C van Gemert, Spencer Cappallo, Thomas Mensink, and Cees GM Snoek. Bag-of-fragments: Selecting and encoding video fragments for event detection and recounting. In *ICMR*, 2015. 3
- [42] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *NIPS*, 2014. 3
- [43] Quoc V Le and Tomas Mikolov. Distributed representations of sentences and documents. In *ICML*, 2014. 3, 4, 5
- [44] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. In *JMLR*, 2003. 3, 5
- [45] Scott Deerwester, Susan T Dumais, George W Furnas, Thomas K Landauer, and Richard Harshman. Indexing by latent semantic analysis. In *JACS*, 1990. 3, 4, 5
- [46] Sumit Chopra, Raia Hadsell, and Yann LeCun. Learning a similarity metric discriminatively, with application to face verification. In *IEEE CVPR*, 2005. 3
- [47] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *IEEE CVPR*, 2015. 5
- [48] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Antonio Torralba, and Aude Oliva. Places: An image database for deep scene understanding. In *arXiv*, 2016. 5
- [49] Wikipedia, 2016. <http://wikipedia.com>. 5
- [50] Ryan Kiros, Yukun Zhu, Ruslan R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Skip-thought vectors. In *NIPS*, 2015. 5
- [51] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. In *JMLR*, 2008. 7
- [52] Yu-Gang Jiang, Guangnan Ye, Shih-Fu Chang, Daniel Ellis, and Alexander C Loui. Consumer video understanding: A benchmark database and an evaluation of human and machine performance. In *ICMR*, 2011. 7
- [53] Yu-Gang Jiang, Zuxuan Wu, Jun Wang, Xiangyang Xue, and Shih-Fu Chang. Exploiting feature and class relationships in video categorization with regularized deep neural networks. In *IEEE TPAMI*, 2017. 8