

Video Event Understanding using Natural Language Descriptions

Vignesh Ramanathan* Percy Liang[†] Li Fei-Fei[†]

*Department of Electrical Engineering, Stanford University

[†]Computer Science Department, Stanford University

{vigneshr, pliang, feifeili}@cs.stanford.edu

Abstract

Human action and role recognition play an important part in complex event understanding. State-of-the-art methods learn action and role models from detailed spatio-temporal annotations, which requires extensive human effort. In this work, we propose a method to learn such models based on natural language descriptions of the training videos, which are easier to collect and scale with the number of actions and roles. There are two challenges with using this form of supervision: First, these descriptions only provide a high-level summary and often do not directly mention the actions and roles occurring in a video. Second, natural language descriptions do not provide temporal annotations of actions and roles. To tackle these challenges, we introduce a topic-based semantic relatedness (SR) measure between a video description and an action and role string, and train the model using posterior regularization. Our event recognition system based on these action and role models matches the state-of-the-art method on the TRECVID-MED11 event kit, despite weaker supervision.

1. Introduction

The ability to differentiate complex events is a key step towards video understanding and has warranted significant research in recent years [17, 8, 23]. Complex events can be thought of as a composition of atomic actions performed by people holding different roles. In this work, we provide a method to learn these action and role models based on easily available natural language descriptions of event videos as shown in Fig. 1. We rely entirely on these descriptions and do not require separate ground truth annotations.

The use of action and/or role models trained with extensive spatio-temporal annotations has shown to boost event recognition performance in videos [8, 12]. Such detailed annotations require expensive human effort and severely restrict the scalability with inclusion of more actions and roles. Complex event datasets like TRECVID-MED11 [1] and MPII Cooking [22] are accompanied by natural lan-



Figure 1. Our method relies on natural language video descriptions to train action and role models. Sample videos along with their text description are shown. The descriptions of videos containing the action “play instrument” are bounded in green. Note that this supervision is not available during training.

guage descriptions, which are easy to obtain and require one-time annotation during the collection of a dataset. Internet repositories such as YouTube already have accompanying descriptions, and require no additional annotation.

In our setting, natural language descriptions provide a high-level summary of the event occurring in a video without temporal synchronization. The challenges associated with their use are two-fold: 1) the presence of an action and role in the video might not be provided directly, and 2) the temporal extents of actions and roles are unavailable.

As demonstrated by the positive training examples of the action “play instrument” in Fig. 1, the video description might not contain an action string, but still imply its presence (in this case through an “orchestra” or a “marching band”). Bridging this gap is a fundamental and challenging NLP problem. To tackle this, we identify training action and role labels using a semantic relatedness (SR) measure between an action or role string, and the textual descriptions. We construct a language topic model specific to the task and use it to define the SR measure. We also handle outliers introduced by this measure through self-paced learning.

The lack of temporal annotations is handled through a Posterior Regularization (PR) [7] based training, which lets us represent a video as a bag of human tracklets with latent action and role assignments to these tracklets. PR allows us to learn a structured model, while enforcing constraints regarding the presence or absence of action and role tracklets in the training videos.

We evaluated our approach on action and role classifi-

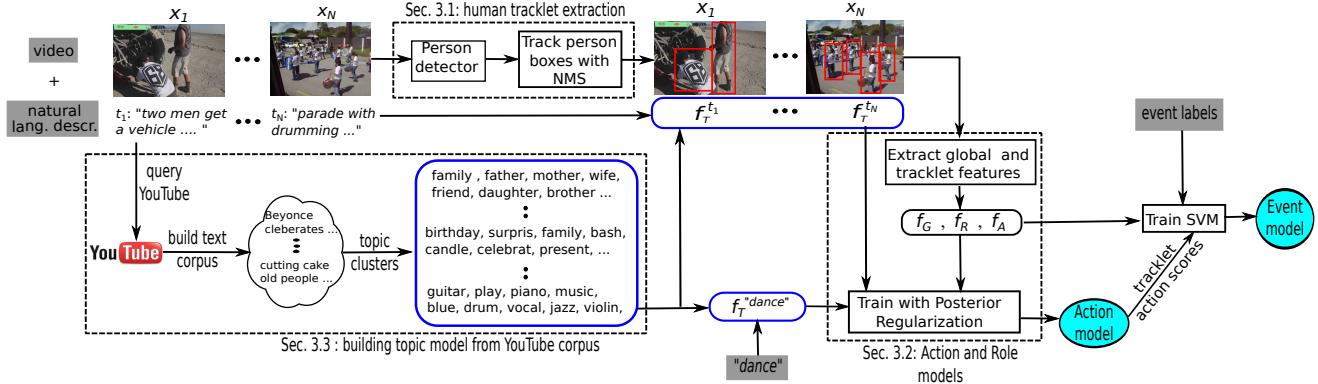


Figure 2. An overview of the system. Input to the system are shaded in grey.

cation, as well as event recognition using our action and role models which are trained with natural language descriptions. On the TRECVID-MED11 event kits, our event recognition model is shown to match the state-of-the-art method [8], which requires detailed spatio-temporal annotations of atomic actions.

2. Related Work

Natural language processing for vision Recent works attempting to leverage the vast amount of textual data available with Internet images have developed vision-specific semantic relatedness measures [25, 24] to identify the link between part-based object attributes and image classes. However such measures are derived from a generic corpus and are less suited for human actions/roles specific to a set of events. Other attempts to use textual descriptions in conjunction with attribute recognition were presented in [2, 19]. [2] is restricted to simple part-based attributes directly mentioned in the image description. [19] involves humans in the loop to actively describe a group of images through visual attributes. Another line of work [9, 28] jointly considers multiple modalities including text descriptions to perform image annotation, retrieval or segmentation. [23] transfers composite action videos to an attribute space enabling comparison with textual corpus.

[13, 16, 3, 4] uses time synchronized movie scripts or closed captions to identify video segments corresponding to specific actions. Again, these methods rely on presence of the action string in the script or use a pre-trained NLP classifier [13] and require temporal annotations. [30] performs tag-prediction by using meta-data provided along with YouTube videos. [18] processes descriptions of action segments to automatically discover a set of action classes.

In contrast to the above methods, we learn models based on natural language descriptions which may not contain the action and role string. In particular, we construct a topic-model based measure specific to our task.

Action, role and event recognition [8, 12] showed signifi-

cant improvement in event recognition by using atomic action and role detectors as a part of their event recognition model. Both methods required spatio-temporal annotation of action and roles in the training videos to learn the models. Other works which have investigated the use of social roles in video understanding include [29, 5]. [14] uses attributes to perform action recognition in videos.

Weakly supervised action models Discriminative spatio-temporal regions in videos or images to localize the actions in [26, 21, 27]. Similar in spirit to these works, we try to localize the human actions and roles. However, we develop a model with latent action and role assignments to different human tracklets in a video.

3. Our Approach

An overview of our system is shown in Fig. 2. We first use natural language video descriptions to train action and role models. The scores from the model are then used to train event recognition models.

In our setup, each training video is accompanied by a natural language description. This description might or might not contain the prominent actions present in the video. Formally, we denote our training dataset by $(\langle x_1, t_1 \rangle, \dots, \langle x_n, t_n \rangle)$, where x_i represents the videos belonging to different event classes and t_i the corresponding textual descriptions. Note that no textual descriptions are assumed to be present with the test data.

We assume a fixed set of actions \mathcal{A} and roles \mathcal{R} and define additional variables $\langle y_i, z_i \rangle$ for each x_i . Here, $y_i^a \in \{-1, 0, 1\}$ indicates whether the label of the video corresponding to the action a is negative, unknown or positive. We define $z_i^r \in \{-1, 0, 1\}$ similarly for the role r . Note that these variables are not observed in the training data.

3.1. Human tracklet extraction

Complex event videos are composed of many atomic actions and roles, confined to smaller spatio-temporal regions. We attempt to incorporate this locality by representing a

video as a bag of human tracklets. We assume that the action or role occurring in a video would then correspond to one or more of these tracklets. As illustrated in the corresponding section of Fig. 2, we obtain tracklets by running a human detector [6] across different segments in a video and tracking the resulting bounding boxes within a temporal window of 100 frames. In our experiments, we uniformly partition a video into 20 different segments and obtain 5 tracklets in each segment based on non-maximal suppression. We choose the top 50 tracklets based on their detection scores.

3.2. Action and role model

We define a conditional random field (CRF) encompassing the action and role relation between different tracklets in a video similar in spirit to [12]. However, we do not assume perfect human tracking, and complete person-wise action and role labels for training.

We assume that each video x_i has a set of human tracklets given by \mathcal{H}_i . The potential $\Phi(x, h, a, r)$ of making action assignment a and role assignment r to the tracklet h in video x is given in Eq. 1.

$$\begin{aligned} \Phi(x, h, a, r) = & w_g(a, r) \cdot f_g^x + w_{in.}(a, r) \\ & + w_{ac.}(a) \cdot f_{ac.}^h + w_{ro.}(r) \cdot f_{ro.}^h, \end{aligned} \quad (1)$$

where f_g^x is the d_g dimensional global video feature of x . The $d_{ac.}$ dimensional feature $f_{ac.}^h$, and $d_{ro.}$ dimensional feature $f_{ro.}^h$ are the action and role features for the human tracklet h respectively. The global weight is denoted by $w_g \in \mathbb{R}^{|\mathcal{A}| \times |\mathcal{R}| \times d_g}$, where $w_g(a, r) \in \mathbb{R}^{d_g}$ gives the global weight for action a and role r . Similarly, $w_{in.} \in \mathbb{R}^{|\mathcal{A}| \times |\mathcal{R}|}$ is the weight for joint action and role assignment to a track, with $w_{in.}(a, r) \in \mathcal{R}$ corresponding to action a and role r . The weight $w_{ac.} \in \mathbb{R}^{|\mathcal{A}| \times d_{ac.}}$ is the action-weight and $w_{ro.} \in \mathbb{R}^{|\mathcal{A}| \times d_{ro.}}$ is the role-weight. The action-weight corresponding to a is given by $w_{ac.}(a) \in \mathbb{R}^{d_{ac.}}$ and the role-weight for role r is given by $w_{ro.}(r) \in \mathbb{R}^{d_{ro.}}$. The weights to be learned are then represented by $w = (w_g, w_{in.}, w_{ac.}, w_{ro.})$.

With a slight abuse of notation, we let $a_i \in \mathcal{A}^{|\mathcal{H}_i|}$ be the action-labels assigned to the tracks in video x_i , and $a_i^h \in \mathcal{A}$ denote the action-label of the track h in the video. Similarly, we let $r_i \in \mathcal{R}^{|\mathcal{H}_i|}$ be the role-labels associated with x_i and $r_i^h \in \mathcal{R}$ be the role-label of track h . The probability $p(a_i, r_i; w)$ of this assignment to video x_i is given by Eq. 2.

$$p(a_i, r_i; w) = \frac{1}{Z_i} \exp \left(\sum_{h \in \mathcal{H}_i} \Phi(x_i, h, a_i^h, r_i^h) \right), \quad (2)$$

, where Z_i is the partition function for the video x_i .

The log-likelihood of making action and role assignments \mathbf{a}, \mathbf{r} respectively across all videos is given by Eq. 3

$$L(\mathbf{a}, \mathbf{r}; w) = \sum_i \log p(a_i, r_i; w) \quad (3)$$

Features: The global video feature uses multiple channels through HOG3D [10], ASR, OCR, MFCC [20] and SIFT [15] features. The features $f_{ac.}$ and $f_{ro.}$ are bag of words HOG3D features extracted from the tracklet h .

Training with posterior regularization: Here, we present a method to learn the model by minimizing L from Eq. 3, assuming the labels $\langle y_i, z_i \rangle$ are given. We will later use natural language annotations to derive these labels in Sec. 3.3. We wish to learn model weights while making latent action and role assignments to each tracklet in the video. The setup is close to the Multi Instance Multi Label framework of [31]. However, to facilitate learning of a structured model with action-role relations, we adopt the more general posterior regularization framework [7]. This enables us to optimize the likelihood subject to soft constraints on the predicted action and role distribution. Formally, let $Q(\mathbf{a}, \mathbf{r})$ be the distribution of action and role assignments in the training videos. We wish to ensure that, in a video tagged as positive for a specific action, the number of tracklets corresponding to the action is at least one. Similarly, in negative videos, the number of tracklets corresponding to the action should be zero. The same argument follows for roles as well. We use these constraints to learn a model by solving the optimization in Eq. 4.

$$\begin{aligned} \min_{\substack{w, Q \\ \delta \geq 0, \eta \geq 0}} \quad & \frac{\|w\|^2}{2} - \mathbb{E}_Q[L] + \sum_i \left\{ \sum_a \delta_i^a + \sum_r \eta_i^r \right\} - H_Q \\ \text{subject to} \quad & \mathbb{E}_Q[N_i(a)] \geq 1 - \delta_i^a, \quad \forall y_i^a = 1 \\ & \mathbb{E}_Q[N_i(a)] \leq \delta_i^a, \quad \forall y_i^a = -1 \\ & \mathbb{E}_Q[M_i(r)] \geq 1 - \eta_i^r, \quad \forall z_i^r = 1 \\ & \mathbb{E}_Q[M_i(r)] \leq \eta_i^r, \quad \forall z_i^r = -1, \end{aligned} \quad (4)$$

where $N_i(a) = \sum_{h \in \mathcal{H}_i} \mathbf{1}(a_i^h = a)$, $M_i(r) = \sum_{h \in \mathcal{H}_i} \mathbf{1}(r_i^h = r)$ and H_Q is the entropy of distribution Q .

We optimize Eq. 4 using a modified Expectation Maximization algorithm shown in Sec. A of the supplementary document.

3.3. Using natural language video descriptions

The natural language description of a video contains rich information regarding the event context and can help infer the presence of specific actions and roles in the video. For instance, Fig. 1 provides examples of descriptions which do not contain the action string “play instrument”, while making this association is easy for humans.

Three measures based on WordNet, World Wide Web (WWW) and Wikipedia were introduced in [25] to deter-

mine the semantic relatedness between class names and attributes. The WordNet metric is a poor indicator of similarity between concepts not linked by a hypernym hierarchy. For instance, in Fig. 1 it would be unable to recognize the relation between “marching band” and “playing instruments” which do not fall under the same subtree. The resulting poor performance of this measure was also noted in [25], making it less useful for our purpose. The WWW metric is tailored to measure the similarity of only a pair of terms based on co-occurrence in Internet repositories but does not offer a concept based similarity. While the Wikipedia SR measure uses Wikipedia concepts, it relies on a generic knowledge base and provides no dimensionality reduction specific to the task. We address this issue by building a language topic model specific to our current task and using it to define the SR measure. We also provide a method to handle outliers introduced by this measure using self-paced learning.

Topic model based SR: A natural source for video descriptions is the vast collection of user provided descriptions along with YouTube videos. Hence, as shown in Fig. 2, we build a text corpus by querying YouTube for frequent terms from the training data descriptions. We generate a topic model from this corpus with 200 topics. Since the text corpus was obtained based on training data descriptions, the generated topic clusters often relate to frequent actions and roles in the data. Sample topic clusters are shown in Fig. 2.

All video descriptions t_i can now be represented by a 200 dimensional vector $f_d^{t_i}$ providing the distribution of the topics in the description. An action a , can be represented by f_d^a giving the topic distribution over the action name. Similarly, f_d^r is defined for a role r . The cosine similarity $\text{sim}(f_d^a, f_d^{t_i})$ provides the proximity of a video x_i to the action a . We will refer to this measure as the *topic model SR*. Each training video can now be assigned a training label $y_i^a \in \{-1, 0, 1\}$ based on a threshold τ as shown in Eq. 5. z_i^r is defined similarly for a role r .

$$y_i^a = \begin{cases} 1 & \text{if } \text{sim}(f_d^a, f_d^{t_i}) \geq (1 - \tau) \text{ or} \\ & t_i \text{ contains action string } a \\ -1 & \text{if } \text{sim}(f_d^a, f_d^{t_i}) \leq \tau \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

Handling outliers Discovering semantic relatedness in textual space is challenging and these measures are not robust. While y_i^a, z_i^r from Eq. 5 can now be directly used in Eq. 4, it would result in significant outliers. We handle this problem by defining a pool of potential positive and negative examples according to Eq. 5 as shown in the first step of Fig. 3 and letting the model gradually choose more examples from this pool in successive iterations as illustrated in the third and fourth step of Fig. 3. This is achieved through a self-paced learning scheme introduced in [11].

We further modify the self-paced method to treat $f_d^{t_i}$ as an additional global feature for x_i in the initial iterations

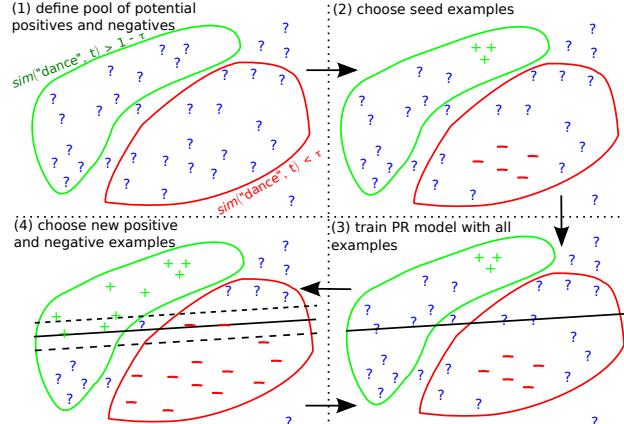


Figure 3. An overview of our self-paced approach shown for one action or role. The green and red boundaries indicate the positive and negative pool of samples chosen using the topic model SR.

but gradually reduce it to zero across the iterations. Intuitively, we are leveraging the textual information present along with videos to choose good examples in the initial phase of the training. However, as the model grows confident with more iterations, the textual features are ignored resulting in a model which only uses video features. The complete details are shown in Sec. B of the supplementary.

3.4. Training event model

We use the action and role detection scores to perform video event classification. The expected number of tracklets corresponding to different actions and roles are used as additional features along with the global video features to train a linear ensemble SVM. We use the same set of global video features from Sec. 3.2. Similar to [8], we first train separate event classifiers for each individual feature mentioned in Sec. 3.2 and finally treat the event classification score from these classifiers as global video features. Since only a small set of actions and roles are usually related to an event, we add an additional L_1 regularization term for the action and role feature weights to encourage only the relevant action and/or role scores to be selected.

4. Experiments

We test our event, action and role classification models on the TRECVID-MED11 event kit videos. The dataset contains videos belonging to 15 complex event classes. Each video is accompanied by a synopsis describing the events in the video, and few of them mention the atomic actions and objects present in the video. We use the same training and testing splits as [8].

4.1. Implementation details

We define *crude* action labels \tilde{y}_i and role labels \tilde{z}_i for each video x_i based on simple text processing. We set

$\tilde{y}_i^a = 1$ if t_i contains the action string a , $\tilde{y}_i^a = -1$ if none of the textual descriptions in the event class of x_i contain the action string a ; otherwise, we set $\tilde{y}_i^a = 0$. We define $\tilde{z}_i^r \in \{-1, 0, 1\}$ similarly for the video x_i and role r . These crude labels are used to train baseline models not using the complete video description as well as to initialize the self-paced scheme in Sec. 3.3. The value of τ is set to consider the top 300 and 30 videos closest to the action and role description respectively as potential positives.

The CRF allows any arbitrary set of actions and roles. In our experiments, we train separate models for each action and role. While training an action model, we consider the relation of the action to all the roles including a null role. Thus, the action model makes a latent role assignment to a tracklet while identifying the action in the tracklet. Similarly, each role model performs latent action assignment to the tracklets while identifying presence of the role in the tracklet. In practice, this makes the learning more tractable and also performs better than training a single model considering all actions and roles together.

4.2. Action and role classification

A set of 62 atomic events were used in [8]. Some of these events were non-human actions like vehicle movement. We select a subset of 46 classes which involve one or more humans. We choose only these action classes which are directly mentioned at least once in the training data descriptions. We consider a set of 13 roles appearing in different events, as listed in Tab. 2. Each video in the test set is annotated with the actions and roles present in it for evaluation.

The action and role classification performance is evaluated by computing the average precision on the testing data as shown in Tab. 1, 2. The expected number of tracklets performing an action in a video is treated as the corresponding action score for the video. Similarly, the expected number of tracklets holding a role in a video provides the role score.

Full model refers to the complete algorithm using video descriptions to train PR models in a self-paced setting. The different baselines are explained below. The first three baselines train only with crude labels $\langle \tilde{y}_i, \tilde{z}_i \rangle$.

- global only: uses global video features to train a SVM.
- simple PR: train action or role models without considering joint action-role relation in Eq. 1.
- full PR: uses action-role relation in addition to tracklet features to train the PR model.
- wiki SR [25]: trains full PR model by identifying positives and negative training examples based on the Wikipedia SR using a threshold as defined in Sec. 4.1.
- topic SR: Our full model without outlier handling through self paced learning.

Comparing the performance of global only and simple PR baselines in Tab. 1, 2, we observe that identifying human tracklets in the videos improves the overall action as well as

Action	global only	simple PR	full PR	wiki SR [25]	topic SR	full model
bending	0.0604	0.0708	0.0689	0.0688	0.0586	0.0601
blowing candles	0.4616	0.4485	0.5088	0.5222	0.4934	0.5134
carving	0.2131	0.0229	0.0918	0.0794	0.0359	0.2348
casting	0.0046	0.0125	0.0118	0.0119	0.0141	0.0135
clapping	0.1433	0.1865	0.2720	0.2615	0.2236	0.2408
cleaning	0.0262	0.0047	0.0047	0.0048	0.0048	0.0240
cutting	0.1928	0.0794	0.0764	0.0776	0.0760	0.1906
cutting cake	0.0885	0.1361	0.1764	0.2803	0.1208	0.1764
cutting fabric	0.1896	0.0152	0.1541	0.1526	0.1557	0.1351
dancing	0.5941	0.5556	0.6189	0.6052	0.6357	0.6261
drilling	0.0570	0.0145	0.0142	0.0157	0.0661	0.0910
drinking	0.0258	0.0347	0.0445	0.0556	0.0421	0.0322
eating	0.0532	0.0522	0.0613	0.0558	0.0598	0.0569
falling	0.1081	0.1697	0.1523	0.1390	0.1513	0.1512
flipping	0.3995	0.4316	0.4554	0.2636	0.4364	0.4524
hammering	0.0794	0.0057	0.0056	0.0057	0.2743	0.2741
jacking car	0.0734	0.0185	0.0172	0.0164	0.0185	0.0373
jumping	0.5572	0.5184	0.5443	0.5734	0.5203	0.5586
kissing	0.1499	0.5232	0.4976	0.5318	0.4716	0.4976
laughing	0.0853	0.1508	0.1624	0.1605	0.1753	0.1611
lighting candle	0.0218	0.0437	0.0805	0.0772	0.0513	0.0631
open door	0.1276	0.0846	0.0692	0.0692	0.1285	0.0989
petting	0.0253	0.0103	0.0103	0.0103	0.0103	0.0115
planing	0.0525	0.0162	0.0140	0.0084	0.0449	0.0555
play instrument	0.1335	0.2424	0.2059	0.2705	0.2083	0.2059
pointing	0.0159	0.0437	0.0466	0.0398	0.0336	0.0238
polishing	0.0015	0.0015	0.0015	0.0015	0.0017	0.0025
pouring	0.0051	0.0061	0.0103	0.0038	0.0088	0.0026
pushing	0.2768	0.2871	0.1922	0.2865	0.1783	0.1824
reeling	0.4603	0.4675	0.4669	0.4973	0.4665	0.4788
rolling	0.0533	0.0074	0.0072	0.0065	0.0078	0.0091
sawing	0.0416	0.0305	0.0628	0.0667	0.0750	0.2390
sewing	0.3073	0.4089	0.2801	0.2839	0.2660	0.2588
shake	0.0067	0.0062	0.0062	0.0058	0.0064	0.0101
singing	0.0384	0.0900	0.0732	0.0721	0.0901	0.0742
sliding	0.0438	0.0811	0.0776	0.0750	0.0761	0.0806
stir	0.0398	0.0208	0.2037	0.2071	0.1975	0.1967
surfing	0.1039	0.1442	0.1494	0.1510	0.1302	0.1382
turning wrench	0.0788	0.0234	0.0233	0.0232	0.0232	0.0573
using knife	0.0015	0.0016	0.0017	0.0023	0.0016	0.0017
using tire tube	0.0675	0.0145	0.0141	0.0136	0.0138	0.0351
walking	0.1771	0.2562	0.2520	0.2110	0.2697	0.2557
washing	0.1329	0.0309	0.0307	0.0309	0.0307	0.0307
waving	0.0965	0.1302	0.1555	0.1413	0.1473	0.1700
wiping	0.0428	0.0253	0.0253	0.0253	0.0405	0.0369
writing	0.0400	0.1309	0.1080	0.1010	0.1413	0.1856
mean	0.1295	0.1356	0.1415	0.1427	0.1453	0.1616

Table 1. Action classification results. The highest score for a class is shown in bold font. The first three columns do not use the complete video descriptions, but train with labels $\langle \tilde{y}_i, \tilde{z}_i \rangle$.

role classification. The effect is even more prominent for roles, since roles are governed by the humans holding the roles, whereas most actions like reeling, cutting, jacking car, turning wrench can be determined by object manipulations in the scene. We notice that actions like kissing, walking, playing instrument which can be determined by observing the complete or the upper human body, benefit more from the human tracklet representation compared to actions like cleaning, cutting, petting, turning wrench which are often shown as close-up shots of the hand. Similarly, human detectors fail to detect people when they are not upright, leading to a drop in performance.

From simple PR and full PR results in Tab. 1, 2, we no-

Role	global only	simple PR	full PR	wiki SR [25]	topic SR	full model
bride	0.7115	0.7946	0.7880	0.7873	0.8017	0.7877
groom	0.6751	0.7755	0.7805	0.7857	0.7901	0.7901
priest	0.2686	0.4263	0.4094	0.4468	0.4002	0.4058
performer	0.1499	0.1212	0.1708	0.1805	0.1722	0.1737
musician	0.0990	0.2933	0.2468	0.2506	0.2334	0.2643
parent	0.2084	0.2388	0.2123	0.2028	0.2014	0.2245
birthday child	0.7442	0.8350	0.8212	0.8086	0.8359	0.8359
audience/guest	0.3163	0.4260	0.3623	0.4263	0.3807	0.4429
friends	0.2849	0.5619	0.5671	0.4979	0.5591	0.5641
fisherman	0.2677	0.0238	0.2949	0.2617	0.2925	0.2873
craftsman	0.0957	0.0284	0.0281	0.0285	0.0282	0.0265
mechanic	0.0446	0.0268	0.0267	0.0262	0.0266	0.0264
police/soldier	0.2267	0.2346	0.2122	0.2398	0.2384	0.2303
mean	0.3148	0.3682	0.3785	0.3802	0.3816	0.3892

Table 2. Role classification results. The highest score for a class is shown in bold font. The first three columns do not use the complete video descriptions, but train with labels $\langle \tilde{y}_i, \tilde{z}_i \rangle$.

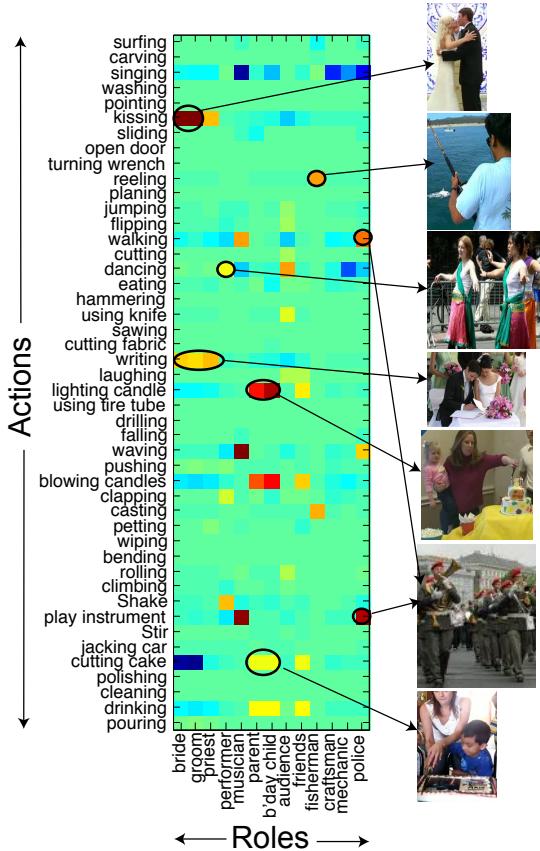


Figure 4. The weights corresponding to different action-role relation $w_I(a, r)$ are shown. Sample frames depicting the action-role relations corresponding to some high weights are also shown.

tice that jointly modeling the action-role relations in the posterior regularized setup increases the performance for action classes like kissing, writing, lighting candle, cutting cake and the roles fisherman, birthday child, performer. In order to analyze the effect of action-role relations, we visualize the joint action-role weights in Fig. 4. As expected,

we see strong correlation between certain action and role classes (highlighted by ovals). These correspondences result in an improved classification accuracy for the respective classes. Sample frames pertaining to some high weights are shown besides the matrix in Fig. 4. We further demonstrate some qualitative results in Fig. 5, where the highest scoring tracklet in a video for a certain action is shown along with the corresponding role assignment.



Figure 5. High scoring tracklets for few action models are shown for six test videos. The role assignments are also shown.

The wiki SR and topic SR models trained by identifying positives and negatives based on a SR measure are seen to only marginally improve the performance. This can be accounted to the addition of false positive and negative training labels, demonstrating the difficulty in processing natural language descriptions. The results are seen to be worse in the case of Wikipedia based SR measure.

To analyze the utility of our topic model, we run experiments where textual descriptions are assumed to be present both during training and testing. We train two separate PR models which use topic model based textual features and Wikipedia SR based textual features respectively as additional global features both during training and testing. The mean AP of these models correspond to the green bars in Fig. 6. For Wikipedia features, we concatenate the Wikipedia SR measure of a description with each of the action and role strings to form a feature. We treat f_d as the topic model feature. The considerable gain achieved by using topic model features during test time over other methods justifies the use of topic model SR for the current task. The lower performance of Wikipedia-based methods can be explained by the use of generic corpus and lack of dimensionality reduction to consider concepts specific to the task.

Further, note that in our adaptation of the self-paced approach, while textual descriptions are not available during testing, we use features extracted from textual descriptions in the initial iterations of training, and finally anneal their weights to zero. The effectiveness of our textual features as shown in Fig. 6 allows us to handle outliers introduced by the SR measure.

From Tab. 1, 2, we see that our full model, which handles outliers, is able to achieve significant improvement over the other methods. Our method is seen to be particularly ef-

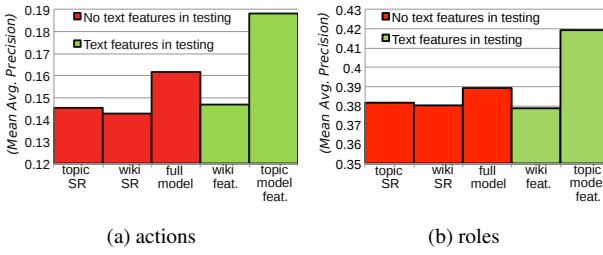


Figure 6. The green bars correspond to the setting where textual descriptions are used at test time. The red bars are from Tab. 1, 2.

fective for action classes like hammering, writing, planing, drilling, where the number of training video descriptions directly containing the action string were below 10. We show sample videos which were added as positives by our full method in Fig. 7, along with the corresponding descriptions. We notice the inclusion of videos whose descriptions do not contain the action directly.

4.3. Event classification

We compare the event classification performance of our model from Sec. 3.4 against baseline methods as well as state-of-the-art results from [8] in Tab. 3. It is to be noted that unlike our method, [8] used extensive spatio-temporal annotation to learn completely supervised atomic action classification models. The scores from these models were finally used to perform event classification. We report two sets of results from [8], one using action classification scores in linear ensemble SVM and the other using them in a joint CRF model. In addition, we demonstrate results against the following baselines

- global only: uses global video features only
- global+actions : uses only action classification features in addition to global video features
- global+roles: uses only role classification features in addition to global video features

From Tab. 3, we observe that our methods using either the action or role features outperform an SVM trained only with global video features. Our full model using both action and role scores achieves the maximum mean AP. Thus, our action and role models trained only with natural language descriptions matches state-of-the-art methods from [8], which uses ground truth spatio-temporal action annotations for training. This confirms the utility of the our action and role models learned with very weak supervision.

5. Conclusion

We have presented a method to learn atomic action and role models based on easily available natural language video descriptions. We proposed a language topic-model based semantic relatedness measure to identify positive and

Event	global only	[8]* SVM	[8]* joint CRF	global + action	global + roles	full model
Boarding trick	0.8766	0.7560	0.7570	0.8276	0.8625	0.8402
Feeding animal	0.4535	0.5820	0.5650	0.4490	0.3958	0.4595
Landing fish	0.6612	0.7410	0.7220	0.6612	0.6811	0.6593
Wedding	0.4729	0.6650	0.6750	0.5942	0.7555	0.7871
Woodworking project	0.2227	0.5760	0.6530	0.3697	0.2086	0.3568
Birthday party	0.9083	0.7090	0.7820	0.9207	0.9041	0.9008
Changing tire	0.5100	0.4650	0.4770	0.5200	0.4977	0.5012
Flash mob	0.9301	0.8590	0.9190	0.9273	0.9248	0.9240
Vehicle unstuck	0.6288	0.6610	0.6910	0.6212	0.5862	0.6173
Grooming animal	0.3881	0.4570	0.5100	0.3914	0.3927	0.5415
Making sandwich	0.5604	0.3560	0.4190	0.5739	0.5442	0.5704
Parade	0.7462	0.6570	0.7240	0.7283	0.6582	0.7335
Parkour	0.5426	0.5340	0.6640	0.6211	0.5681	0.6144
Repairing appliance	0.8025	0.8080	0.7820	0.7989	0.7692	0.7840
Sewing project	0.6579	0.5690	0.5750	0.6563	0.6286	0.6688
mean	0.6241	0.6263	0.6610	0.6441	0.6252	0.6639

Table 3. Event classification results. The highest score for a class is shown in bold font. * Unlike our method, [8] uses extensive ground truth spatio-temporal annotations for training separate action classifiers to aid event classification.

negative training examples. These labels were used to train a CRF model with posterior regularization, while making latent action and role assignments to human tracklets. Outliers introduced by the SR measure were handled through a self-paced scheme. The action and role models were used to achieve state-of-the-art event classification performance on the TRECVID-MED11 event kits. We demonstrated the efficacy of the topic model based SR measure in identifying training labels as well as the gain due to the posterior regularized method in a weakly supervised setting without temporal annotations. In future work, we wish to ground action-object relations in complex event videos based on natural language descriptions and simultaneously learn these relation from both textual and visual data.

Acknowledgements

We thank O. Russakovsky, K. Tang and B. Yao for helpful comments. This research is partially supported by the NFS grant ??????.

References

- [1] Trecvid multimedia event detection track, 2011. 1
- [2] T. L. Berg, A. C. Berg, and J. Shih. Automatic attribute discovery and characteriztion from noisy web data. In *ECCV*, 2010. 2
- [3] T. Cour, C. Jordan, E. Miltzakaki, and B. Taskar. Movie/script: Alignment and parsing of video and text transcription. In *ECCV*, 2008. 2
- [4] M. Everingham, J. Sivic, and A. Zisserman. Hello! my name is... buffy automatic naming of characters in tv video. In *BMVC*, 2006. 2
- [5] A. Fathi, J. K. Hoggins, and J. M. Rehg. Social interactions: A first person perspective. In *CVPR*, 2012. 2



Figure 7. Videos with highest SR measure added as positives by our full method for different actions are shown. The last column shows wrong videos which are added by the method. The video descriptions are shown below it. Note that they do not contain the action string.

- [6] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part based models. *PAMI*, 32(9):1627–1645, 2010. 3
- [7] K. Ganchev, J. Graca, J. Gillevater, and B. Taskar. Posterior regularization for structured latent variable models. *JMLR*, 11, 2010. 1, 3
- [8] H. Izadinia and M. Shah. Recognizing complex events using large margin joint low-level event model. In *ECCV*, 2012. 1, 2, 4, 5, 7
- [9] Y. Jia, M. Salzmann, and T. Darrell. Learning cross-modality similarity for multinomial data. In *ICCV*, 2011. 2
- [10] A. Kläser, M. Marszałek, and C. Schmid. A spatio-temporal descriptor based on 3d-gradients. In *BMVC*, 2008. 3
- [11] M. P. Kumar, B. Packer, and D. Koller. Self-paced learning for latent variable models. In *NIPS*, 2010. 4
- [12] T. Lan, L. Sigal, and G. Mori. Social roles in hierarchical models for human activity recognition. In *CVPR*, 2012. 1, 2, 3
- [13] I. Laptev, M. Marszałek, C. Schmid, and B. Rozenfeld. Learning realistic human actions from movies. In *CVPR*, 2008. 2
- [14] J. Liu, B. Kuipers, and S. Savarese. Recognizing human actions by attributes. In *CVPR*, 2011. 2
- [15] D. G. Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, 60, 2004. 3
- [16] M. Marszałek, I. Laptev, and C. Schmid. Actions in context. In *CVPR*, 2009. 2
- [17] M. Merler, B. Huang, L. Xie, G. Hua, and A. Natsev. Semantic model vectors for complex video event recognition. *IEEE Trans. Multimedia*, 14:88–101, 2012. 1
- [18] T. S. Motwani and R. J. Mooney. Improving video activity recognition using object recognition and text mining. In *ECAI*, 2012. 2
- [19] D. Parikh and K. Grauman. Interactively building a discriminative vocabulary of nameable attributes. In *CVPR*, 2011. 2
- [20] L. Rabiner and B.-H. Juang. *Fundamentals of speech recognition*. Prentice-Hall, Inc., 1993. 3
- [21] M. Raptis, I. Kokkinos, and S. Soatto. Discovering discriminative action parts from mid-level vide representations. In *ECCV*, 2012. 2
- [22] M. Regneri, M. Rohrbach, D. Wetzel, S. Thater, B. Schiele, and M. Pinkal. Grounding action descriptions in videos. *TACL*, 1:25–36, 2013. 1
- [23] M. Rohrbach, M. Regneri, M. Andriluka, S. Amin, M. Pinkal, and B. Schiele. Script data for attribute-based recognition of composite activities. In *ECCV*, 2012. 1, 2
- [24] M. Rohrbach, M. Stark, G. Szarvas, and B. Schiele. Combining language sources and robust semantic relatedness for attribute-based knowledge transfer. In *ECCV-PnA*, 2010. 2
- [25] M. Rohrbach, M. Stark, G. Szarvas, and B. Schiele. What helps where – and why? semantic relatedness for knowledge transfer. In *CVPR*, 2010. 2, 3, 4, 5, 6
- [26] S. Satkin and M. Hebert. Modelling the temporal extent of action. In *ECCV*, 2010. 2
- [27] S. Singh, A. Gupta, and A. A. Efros. Unsupervised discovery of mid-level discriminative patches. In *ECCV*, 2012. 2
- [28] R. Socher and L. Fei-Fei. Connecting modalities: Semi-supervised segmentation and annotation of images using unaligned text corpora. In *CVPR*, 2010. 2
- [29] G. Wang, A. Gallagher, J. Luo, and D. Forsyth. Seeing people in social context: recognizing people and social relationships. In *ECCV*, 2010. 2
- [30] W. Yang and G. Toderici. Discriminative tag learning on youtube videos with latent sub-tags. In *CVPR*, 2011. 2
- [31] Z.-H. Zhou and M.-L. Zhang. Multi-instance multi-label learning with application to scene classification. In *NIPS*, 2007. 3