

# Action Classification and Highlighting in Videos

Atousa Torabi  
Disney Research Pittsburgh

atousa.torabi@disneyresearch.com

Leonid Sigal  
Disney Research Pittsburgh

lsigal@disneyresearch.com

## Abstract

Inspired by recent advances in neural machine translation, that jointly align and translate using encoder-decoder networks equipped with attention, we propose an attention-based LSTM model for human activity recognition. Our model jointly learns to classify actions and highlight frames associated with the action, by attending to salient visual information through a jointly learned soft-attention networks. We explore attention informed by various forms of visual semantic features, including those encoding actions, objects and scenes. We qualitatively show that soft-attention can learn to effectively attend to important objects and scene information correlated with specific human actions. Further, we show that, quantitatively, our attention-based LSTM outperforms the vanilla LSTM and CNN models used by state-of-the-art methods. On a large-scale youtube video dataset, ActivityNet [4], our model outperforms competing methods in action classification.

## 1. Introduction

Activity recognition, classification and understanding are important problems in computer vision. The popularity of these problems, empowered by recent availability of large-scale on-line video resources (e.g., YouTube, Google, Bing) and datasets (e.g., ActivityNet [4]), stems from the vast number of applications that such technology will enable that include surveillance, video understanding, intelligent search and retrieval, automated metadata tagging and human-computer interactions. However, activity recognition has also proven to be an immensely challenging visual task as granularity, viewpoint, appearance and duration of activities can vary greatly and both spatial and temporal context appear to be critically important.

Neural network-based architectures, particularly Convolutional Neural Network (CNN), and end-to-end learning have recently illustrated great success in advancing state-of-the-art in variety of problems in computer vision, including image classification [17, 39, 34], object detection [8, 11] and most recently action recognition [33, 40]. However,

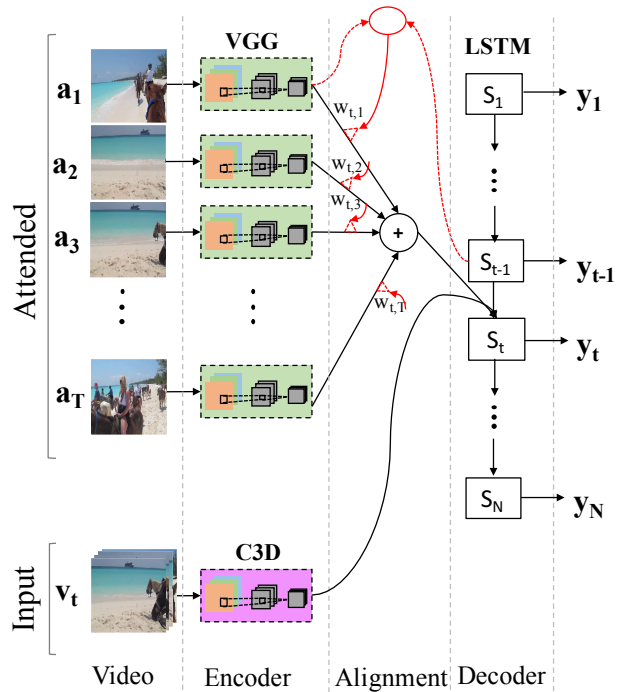


Figure 1. We use an encoder-decoder framework for action recognition. Our LSTM decoder is equipped with an attention/alignment model that allows it focus on salient context/attended data. Attended-data is encoded using VGG network [34] and LSTM input-data is encoded using C3D network [40].

such CNN models largely ignore the temporal progression of the action and typically learn frame-based or short 3-D spatio-temporal filters [40] and aggregate information using simple maximum or average pooling across fixed-length segments [33, 40] or the entire video. Very recently, in less than a handful of papers, Recurrent Neural Network (RNN), that have long history in speech recognition and language translation, have been re-discovered as effective means of more explicitly modeling temporal context in action recognition [37, 26, 47]. One significant limitation of RNN is inability to backpropagate error through long-range temporal interval (a problem known as vanishing gradient effect). A class of models that was specifically developed to be immune to such effect, through the use of simple memory-cell-

like gates, was first introduced by Hochreiter and Schmidhuber [10] and termed Long Short-Term Memory (LSTM).

LSTMs coupled with CNNs have many properties that make them appealing and notably appropriate for temporal perception tasks (hierarchical representation learning, temporal context with memory, possible end-to-end learning). However, such models effectively distill each frame into static learned representation, using one or more CNNs, and aggregate this(these) representation(s) across time using LSTM. One seemingly important property of human perception that is glaringly missing is *attention*, *i.e.*, tendency not to process a whole video at once and rather attend selectively to parts of video that are most important for a given (*e.g.*, activity) recognition task. We note that such attention mechanism is particularly important for videos that contain visual clutter. By attending, parts of the video that are most discriminative can be *highlighted* and used to better inform recognition model during learning and inference.

In this paper, we explore the use of end-to-end encoder-decoder LSTM framework with the built-in attention mechanism. The most distinguishing feature of the proposed approach is that it does not attempt to utilize (encode) the entire input video for prediction of the action class (decoder). Instead it encodes a video into a temporal sequence of visual representations and chooses an adaptively weighted subset of that sequence for prediction. Each time model generates an action class prediction it soft-searches for a set of temporal positions within the video where most salient relevant information is concentrated. This is a much more flexible mechanism, as compared to, for example, earlier structured latent variable models [28], that required much more restrictive conditional independence constraints, or fixed low-level attention schemes [25, 43] which failed to significantly improve the performance or scale in practice.

**Contributions:** To the best of our knowledge, we are the first to show significant performance improvements that can come from temporal attention in action classification task. While in [50] attention-based LSTM is also proposed, the improvements due to attention network are shown to be minimal ( $< 0.5\%$ ). In contrast, we show that by attending on i) much longer sequences and on ii) complementary semantic contextual features (*e.g.*, objects, actions and scenes), much larger improvements in performance can be made (of up to 20% or 8 percentage points). Further, we show that individual models, trained to attend on these different features, are themselves complementary and can be combined for additional performance gains. Conceptually, attention network allows our model to automatically learn to highlight salient spatial/spatio-temporal cues while integrating that information with LSTM model for time-step action classification. Qualitatively, we show that attention network learns to consistently attend to important semantics of the actions, despite the noise present in the input sequences,

which adds to interpretability of the model. Quantitatively, we show improved performance compared to the state-of-the-art on the large scale ActivityNet [4] dataset.

## 2. Related Work

Recognition and understanding of human activities is a well studied problem in computer vision. As such, complete overview is beyond the scope of this paper and we refer reader to surveys [1, 27] for a more historical perspective. Here we focus on recent and most relevant approaches.

**Traditional action recognition:** Traditional action recognition approaches largely focused on global video representations that achieved good results on smaller and simpler datasets, *e.g.*, KTH [30], HMDB51 [18] and UCF101 [36]. Such approaches focused on patch-based local motion and appearance information in the form of Histogram of Oriented Gradients (HOG), Histogram of Optical Flow (HOF) [20], Motion Boundary Histogram (MBH) [46] or dense-trajectories [44, 45]. Classification was typically achieved by further aggregating these local representations across videos or using spatio-temporal pyramids [22] with Bag-of-Words(BoW) [20] or Fisher vector based encodings [45] followed by traditional classifiers (*e.g.*, SVMs).

**CNN-based approaches:** Availability of larger datasets and recent advances in representation learning for recognition, in the form of CNN, have led to the use of CNNs in action classification and understanding. The benefits of CNNs, *e.g.*, over BoW, stem from their ability to learn relevant and discriminative representations, at a hierarchy of spatio-temporal granularities. Among the earliest, [16, 14] used 3D CNN to encode few video frames, showing that encoding multiple frames using raw pixels is marginally better than single-frame CNN models for video classification [16]. Simonyan *et al.* [33] incorporated stack of 10-15 optical flow images, instead of raw pixels, in CNN model trained to encode motion information. Frame-level spatial information was learned using a separate CNN stream, and combination the two CNN models was shown to outperform traditional approaches. More recently, Tran *et al.* [40] learned generic video feature, in the form of 3D CNN, and applied it to variety of video understanding tasks. To incorporate some longer-term temporal information, Zha *et al.* [51] experimented with different spatial and temporal pooling over CNN features, applying the resulting model to TRECVID multimedia event detection challenge; similarly, Ng *et al.* [26] study variety of CNN pooling strategies that process up to 120 frames and obtain impressive results on Sport-1M video classification.

**Temporal models and LSTM:** The core idea that temporal progression of an activity is important, has a long history in the field. A variety of early models have been utilized for modeling of temporal context, including Hidden

semi-Markov Models (HSMM) [6], CRFs [35], and finite-state-machines [12]. Most recently, RNNs have also been explored; LSTM, a form of RNN, has recently become particularly popular due to its ability to deal with vanishing and exploding gradients during training. RNN/LSTM have a number of notable benefits with respect to, for example, HMMs or CRFs, including fewer conditional independence assumptions and possible end-to-end training when combined with CNN. Several recent works incorporate spatial/optical-flow CNN features with vanilla LSTM models for global temporal modeling of videos [37, 5, 26, 47, 42]. These works showed improvements with respect to action recognition [5, 22, 47] or video description tasks [49, 42]. Similar to these methods, we also incorporate CNN features with LSTM for action recognition. However, we use a new variant of LSTM model equipped with an attention network that allows it to focus and highlight discriminative frames and use object and scene features as context.

**Temporal and spatial attention:** There have been previous attempts to ask whether information from all frames is equally important for video recognition/understanding [29]. To this end a number of keyframe-based representations have been proposed. For example, [52] rank all frames based on holistic information theoretic measure; [23] rely on spatio-temporal localization and AdaBoost to select keyframes; [28] automatically select compact keyframe representation using a latent structured SVM formulation. In complement, low-level spatial attention models, trained from eye-gaze in video, have been studied in [25, 43], with relatively limited success. More recently, Shapovalova *et al.* [31] use eye-gaze information as weak-supervision for action localization. Our model also has an attention mechanism, but compared with [28] or [31] that suffer from complexities imposed by discrete latent variables, avoids scalability issues by employing soft-attention coupled with LSTM. Our *dynamic* soft-attention mechanism is also in stark contrast with static low-level attention models of [25] and [43].

**LSTM models with attention:** Models that combine neural network architectures with attention are relatively recent (perhaps the earliest mention is in the context of Boltzmann Machine [21]). In vision, Xu *et al.* [48] use spatial soft-attention that learns to fix gaze on salient image parts while generating description for images. Yao *et al.* [49] also use soft-attention mechanism on video frames for video description. Perhaps in the closest work to ours, Yeung *et al.* [50] propose a multiLSTM for action detection, where as input-data their model aggregates few frames CNN features using soft-attention and generates multiple outputs at each time step corresponding to multiple frames. The results show that adding soft-attention for weighted averaging of few input frames only marginally ( $< 0.5\%$ ) improves

the performance over uniform average. In contrast, we apply soft-attention mechanism on temporally longer contextual information (*e.g.*, object, scene or complementary action features) rather than LSTM inputs. We show that in this way our jointly trained model can better focus on, or highlight, important frames, which substantially improves action classification performance and gives some semantic interpretation for results.

### 3. Approach

In our model, activity recognition is defined as sequential problem where inputs are video features,  $\mathbf{v}_t$ , over time and goal is to predict a single label from a fixed action vocabulary  $\langle \mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_t \rangle \rightarrow y$ . We make use of the encoder-decoder framework, where the input video features are encoded using CNN models and the prediction is done using LSTM as a decoder. Figure 1 shows the overview.

In particular, we build on attention-based LSTM [2] for machine translation, and propose a model that takes two separate sets of input video features: video-level spatio-temporal features that we call *Input-data* and frame-level semantic image features (*e.g.*, encoding objects and scenes) that we call *Attended-data*. In our model, Input-data and Attended-data are encoded using different CNN models and are described in detail in Sections 3.1 and 3.2 respectively. We use soft attention mechanism that learns to highlight (or weight) frames with relevant visual semantics (*i.e.*, objects or scenes important for an action label prediction)<sup>1</sup>. We describe our soft attention model for action recognition task in Section 3.3 and the final LSTM decoder in Section 3.4.

#### 3.1. Input-data: CNN Encoder

We assume the input-data to our encoder is a fixed temporal-length window of  $F$  frames. We extract video features from this window using spatio-temporal CNN implemented in Caffe [15]. As a result, input window is represented using a matrix

$$\mathbf{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_N\} \in \mathbb{R}^{N \times d_v} \quad (1)$$

of  $N$  video feature vectors, where each vector  $\mathbf{v}_i \in \mathbb{R}^{d_v}$  is a  $d_v$ -dimensional representation corresponding to a chunk of  $F/N$  frames. Note, that  $N$  will also correspond to the length of the LSTM decoder sequence; *i.e.*, for every video feature vector we will get an activity label prediction.

In particular, we make use of C3D network pre-trained on Sport-1-million dataset and fine-tuned it on trimmed ActivityNet videos for activityNet experiments. C3D network, proposed in [40], has 8 convolutional layers and 2 fully connected layers, followed by a softmax output layer. All 3D

<sup>1</sup>In the context of attention-based LSTM models this is often called *soft alignment*, where in this case alignment is between the output action label and the input frames.

convolution kernels are  $3 \times 3 \times 3$  with stride 1 in both spatial and temporal dimensions. We extract first Fully-Connected layer (FC6) as a feature vector representation, therefore  $d_v = 4096$ .

### 3.2. Attended-data: CNN Encoder

For attended-data, we extract video features using spatial CNN in Caffe [15]. As a result, attended-data is represented using a matrix,

$$\mathbf{A} = \{\mathbf{a}_1, \dots, \mathbf{a}_T\}, \mathbf{a}_i \in \mathbb{R}^{T \times d_a}, \quad (2)$$

of  $T$  video feature vectors where each vector  $\mathbf{a}_i \in \mathbb{R}^{d_a}$  is a  $d_a$ -dimensional representation corresponding to a single frame. In our model each attended frame is sampled every  $F/N$  video frames, therefore  $T = N$ . We make use of VGG-19 network [34] trained for three task of object, action, and scene recognition:

- $VGG_{obj}$ : VGG trained, on ImageNet, for 1000 class object classification [34];
- $VGG_{act}$ : VGG-19 network fine-tuned on ActivityNet data for action recognition of 203 action classes;
- $VGG_{sce}$ : VGG-19 network fine-tuned on MIT-Scenes for scene recognition of 205 scene classes [53].

VGG-19 proposed in [34], has 16 convoluntal layers and 3 fully connected layers, followed by a softmax output layer. For some experiments we extract first Fully-Connected layer (FC6), therefore  $d_a = 4096$ , and for others we extract last Fully-Connected layer (FC8) as a frame representation. For  $VGG_{obj}$  FC8 layer  $d_a = 1000$ , for  $VGG_{act}$  FC8 layer  $d_a = 203$ , and for  $VGG_{sce}$  FC8 layer  $d_a = 205$ .

In general, size of attended-data  $T$  could be different from size of input-data  $N$  ( $T \neq N$ ). For example, one may want to attend to objects/scenes/actions at a coarser or finer temporal granularity compared to the granularity of label predictions. In experiments, we also consider effect of finer grained attention, *e.g.*, sampled every  $F/(2N)$  such that  $T = 2N$ . We also evaluated significance of attended-data size, comparing  $T = 12$  and  $T = 24$ . We only observed a negligible improvement using bigger attended-data size.

### 3.3. Action Attention and Highlighting Model

Attention model is parametrized as a feedforward neural network that is jointly trained with LSTM decoder (see Section 3.4). In other words, the gradients of LSTM loss function are backpropagated not only through LSTM network but also through soft-attention network that allows joint training of both networks. This model has been originally proposed by Bahdanau *et al.* [2] for machine translation. We adopted this model for action highlighting with

attended-data that we described in the previous section. Attention model generates a score for each attended frame  $\mathbf{a}_i$  at time step  $t$ . This score represents how well a sequentially modeled activity video up to time  $t - 1$  and attended  $\mathbf{a}_i$  are semantically matched. In other words, matching score  $m_{t,i}$  measures the relevance of  $\mathbf{a}_i$  and LSTM decoder hidden state at time  $t - 1$  and is defined as

$$m_{t,i} = \Phi(\mathbf{h}_{t-1}, \mathbf{a}_i), \quad (3)$$

where  $\mathbf{h}_{t-1}$  is LSTM decoder hidden state at time  $t - 1$  that contains all information related to a temporally modeled video sample up to time step  $t - 1$ . For model  $\Phi$ , every state  $\mathbf{h}_{t-1}$  is summed with every attended-data  $\mathbf{a}_i$  to obtain a matching-vectors. Then matching-vector is transformed into a single number named matching-score  $m_{t,i}$ . Note, the matching-vector dimension is one of the network parameters and, in general, could have different value than LSTM hidden-state. Finally, attention weights are computed by normalizing matching-scores in a softmax-like function. Attention weight  $w_{t,i}$  for attended frame  $i$  at time  $t$  is computed by:

$$w_{t,i} = \frac{\exp(m_{t,i})}{\sum_{j=1}^T \exp(m_{t,j})}. \quad (4)$$

A higher attention weight reflects more saliency attributed to a specific attended frame  $i$  with respect to the specific task (action classification).

**Action Attention Context:** In order to integrate attended-data with LSTM decoder, we compute the weighted average of the attended-data using the attention weights obtained above. We term this procedure - *action attention context*:

$$\mathbf{k}_t(A) = \sum_{i=1}^T w_{t,i} \mathbf{a}_i. \quad (5)$$

### 3.4. Attention-Based LSTM Decoder

Long Short-Term Memory (LSTM) is a specific type of a recurrent neural network that recently has been widely used in various domain for sequential modeling, due to its capabilities to cope with "vanishing gradients" problem and its "internal memory" property. In high-level definition, our attention-based LSTM decoder  $\psi$  sequentially updates LSTM state  $\mathbf{h}_t$  and internal memory  $\mathbf{c}_t$  given previous state  $\mathbf{h}_{t-1}$ , current input-data  $\mathbf{v}_t$ , and action attention context vector  $\mathbf{k}_t(A)$ . Also, simultaneously, it predicts action label  $\mathbf{y}_t$  at time  $t$ . In abstract, our LSTM decoder is defined as following:

$$\begin{bmatrix} \mathbf{h}_t \\ \mathbf{y}_t \\ \mathbf{c}_t \end{bmatrix} = \psi(\mathbf{h}_{t-1}, \mathbf{v}_t, \mathbf{k}_t(A)). \quad (6)$$



In order to fuse input-data and attended-data, the information from *action attention context* are distributed across the input-data  $\mathbf{v}_t$ . We used the version of LSTM in [9], where  $\mathbf{h}_t$  is updated by the following intermediate functions:

$$\mathbf{i}_t = \sigma(W_i \mathbf{v}_t + U_i \mathbf{h}_{t-1} + \mathbf{b}_i) \quad (7)$$

$$\mathbf{f}_t = \sigma(W_f \mathbf{v}_t + U_f \mathbf{h}_{t-1} + \mathbf{b}_f) \quad (8)$$

$$\mathbf{o}_t = \sigma(W_o \mathbf{v}_t + U_o \mathbf{h}_{t-1} + \mathbf{b}_o) \quad (9)$$

$$\mathbf{g}_t = \tanh(W_c \mathbf{v}_t + U_c \mathbf{h}_{t-1} + \mathbf{b}_c) \quad (10)$$

$$\mathbf{c}_t = \mathbf{f}_t \mathbf{c}_{t-1} + \mathbf{i}_t \mathbf{g}_t \quad (11)$$

$$\mathbf{h}_t = \mathbf{o}_t \tanh(\mathbf{c}_t), \quad (12)$$

where  $\mathbf{i}_t$ ,  $\mathbf{f}_t$ ,  $\mathbf{o}_t$ , and  $\mathbf{c}_t$  are respectively input gate, forget gate, output gate, and cell gate (*i.e.*, internal memory) at time  $t$ . In brief  $\mathbf{i}_t$ ,  $\mathbf{f}_t$ , and  $\mathbf{o}_t$  control influence of current input  $\mathbf{g}_t$  and previous memory state  $\mathbf{c}_{t-1}$  on generating new hidden state  $\mathbf{h}_t$ . The terms  $W_i, W_f, W_o, W_c, U_i, U_f, U_o, U_c$  encode the various weight matrices and  $\mathbf{b}_i, \mathbf{b}_f, \mathbf{b}_o, \mathbf{b}_c$  corresponding biases. The prediction  $\mathbf{y}_t$  is computed by linear transform on hidden state  $\mathbf{h}_t$  followed by a softmax:

$$\mathbf{y}_t = \text{softmax}(W_y \mathbf{h}_t + \mathbf{b}_y), \quad (13)$$

where  $\mathbf{y}_t \in \mathbb{R}^C$  is probability distribution over fixed number of  $C$  classes.

**Action classification:** We average the scores of  $N$  time-steps of LSTM for each video chunk of  $F$  frames (more details about classification are in Section 4.1).

## 4. Experiments

We use ActivityNet dataset. ActivityNet [4] is a recent large-scale fine-grained human activity dataset with diverse 203 classes organized in an ontology. ActivityNet includes 190 Youtube video samples per class, on average, and video samples average 1382.3 frames (frame rate is not identical for all the videos, however, for most videos it is 29 frame per second). Samples correspond to trimmed activities from longer videos. An action label is assigned to each trimmed video clip. The trimmed data and labels are provided as part of standard activity classification task defined in [4]. However, the trimmed videos are very noisy and contain various visual concepts which are not directly related to the action label. The dataset includes three splits corresponding to training, validation, and test sets. The test set is not available (withheld by the authors). In our experiment we use video samples with minimum 192 frames (about 90% of videos). To this end, we split the original 25,124 training split into 22,624 training samples and 2,500 validations samples for training our algorithm. We used ActivityNet validation split, defined in [4], consisting of 8,403 videos, as our test set.

### 4.1. Implementation Details

We implemented our attention-enabled LSTM decoder in Theano-based framework named BLOCKS [3, 41]. Our data pipeline has been implemented using FUEL framework [41]. Our LSTM models have  $N = 12$  time-steps and we set hidden-state dimension to 512. The matching-vector dimension is also 512. Each video sample has fixed size of  $F = 192$  frames. The input-data is encoded by C3D network every  $F/N$  frames which here is 16 frames (see Section 3.1 for details). The choice of 16 frame inputs for C3D network was made in light of experimental findings in [40]. Attended-data is encoded by VGG networks applied to frames sampled every 16 frames, from the same video of length 192 (see Section 3.2 for details). We added a mlp fully-connected layer before LSTM. This layer has a sigmoid non-linearity applied on input-data without changing input-data dimension. For ActivityNet classification, we use softmax loss function.

**Training:** We train with stochastic gradient descent and apply dropout with probability 0.5 to all LSTM non-recurrent connections and fully-connected layer.

For ActivityNet, we use a fixed 0.01 learning rate for 2,512 iterations on minibatches of size 200. After that, the learning rate is decreased every 125 iterations by a fixed ratio. We also used a fixed 0.9 momentum and use early stopping strategy on validation set loss. We augment the dataset by adding horizontally-flipped frames to training data. We perform an on-the-fly temporal cropping of training videos which means that at each iteration we randomly crop a fixed-size video chunk of 192 frames from each video sample in the mini-batch.

**Testing:** For each test video, we randomly crop 40 video chunks of 192 frames and 40 video chunks of 192 horizontally flipped frames. The class scores for a test video clip are then computed by averaging the scores across all 80 chunks.

**Evaluation Metrics:** For ActivityNet, we compute the average accuracy over all the trimmed validation samples (used as test set; see above). We also report mean Average Precision (mAP) over all trimmed validation samples and over all 203 classes.

### 4.2. Qualitative Results

In order to show how well our attention networks align action with contextual data, we present multiple visualizations. Figure 2 shows the video frames with the highest (green rectangle) and lowest (red rectangle) attention weights for 3 different samples of each class. The top row shows examples for class *setting the table* and *playing guitar* respectively. For these examples contextual data (attended-data) are extracted from  $VGG_{obj}$ . You can see that for first example, images containing *table, chair, cup,*



Figure 2. Three examples show frames with highest (green) and lowest (red) attention weights for each of the classes (left-to-right and top-to-bottom): *setting the table*, *playing guitar*, *sailing*, *mountain climbing*, *curling*, and *playing violin*. For *setting table* and *playing guitar* examples attention is on *objects* context ( $aLSTM-VGG_{obj}$  (fc8)), for *sailing* and *climbing mountain* examples attention is on *scenes* context ( $aLSTM-VGG_{sce}$  (fc8)), and for *curling* and *playing violin* examples attention is on *action* context ( $aLSTM-VGG_{act}$  (fc8)).

*dish*, *glass*, *candle* and *napkin* have maximum attention weights and image with *big couch*, very zoomed image, or image of a *woman sitting in a bedroom* have the minimum attention weights. In next example, the alignment model focus on *guitar* images and occluded images with text or an interviewer picture have the minimum attention weight for class *playing guitar*. Middle row shows 3 video examples for class *sailing* and *mountain climbing* respectively. For these examples the attended features are extracted using  $VGG_{sce}$ . In the first example, images related to the scene including *sea*, *sky* and *sailing boat* have the highest attention weights and images that contain no scene salient features, related to sailing, have the minimum attention weights. Also in the next example, scenes show-

Table 1. ActivityNet: Our model action recognition results.

Model: CNN			
Features	AC	mAP	
$C3D$	(fc8)	40.9%	40.0%
Model: LSTM			
Features	AC	mAP	
$LSTM-C3D$	(fc6)	40.3%	40.6%
$aLSTM-C3D$	(fc8)	42.7%	42.5%
$aLSTM-VGG_{obj}$	(fc6)	45.7%	46.2%
$aLSTM-VGG_{obj}$	(fc8)	47.4%	48.0%
$aLSTM-VGG_{sce}$	(fc6)	44.0%	43.9%
$aLSTM-VGG_{sce}$	(fc8)	45.3%	44.9%
$aLSTM-VGG_{act}$	(fc6)	47.0%	47.3%
$aLSTM-VGG_{act}$	(fc8)	<b>48.1%</b>	<b>48.6%</b>

ing *mountain*, *climbers*, *sky*, and *snow* have the highest attention weights and close-up images or text, with no features related to *mountain climbing*, have minimum attention weights. Finally, last row shows 3 video examples for class *curling* and *playing guitar*. For these examples the attended features are extracted using  $VGG_{act}$ . As one can see, most images with maximum attention weight mainly contain *curling* and *playing guitar*. Besides showing that attention concentrates on relevant objects, actions and scenes, the figure also shows that there is reasonable alignment/visual consistency across the attended frames for each action.

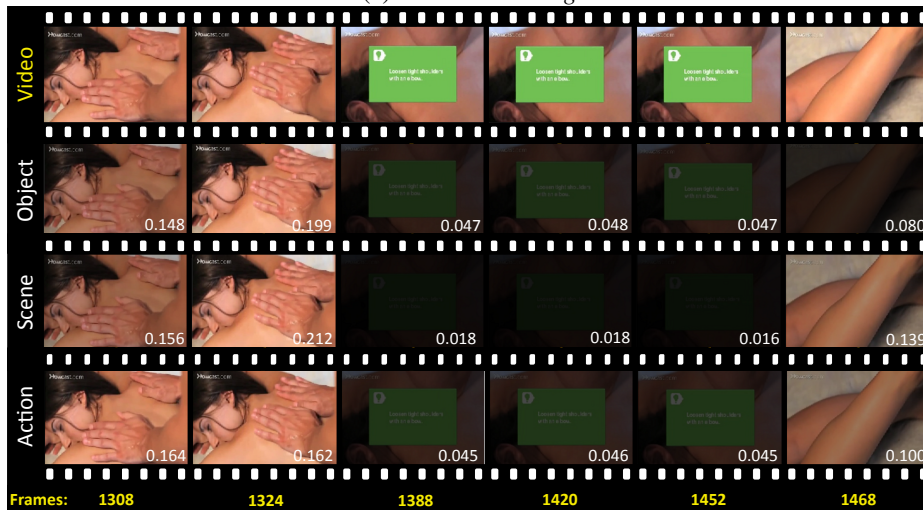
Figure 3 shows 2 video samples for 2 action classes. In this figure we show 6 sampled frames out of 12 from original video that are passed to  $aLSTM$  for context; weighted video frames for attention-based model focus on salient object features, salient scene features and salient action features respectively. For each example, we picked 6 frame samples that show variety of attention weights. The darker images represent lower attention weights, therefore less significant features, and the lighter images have higher attention weights, therefore containing more salient aspects related to corresponding action class. For all examples attention mechanism mainly focus on salient and discriminative video frames. However, interestingly, they are some variations between what different models attend to. For instance, in *having massage* example, for  $aLSTM$  model attending on scene, images containing big text box have minimal weight, since those images have no salient features related to natural scenes. For  $aLSTM$  model attending on object, in *yatching* example, last image has relatively high attention weight. In fact, if you search for *yatching* on Google, the first category that comes up is life style and top example photos are images of people in swimwear on boats.

### 4.3. Comparative Study of Our Model Variants

Table 4 shows average accuracy (acc) and mean average precision over all classes (mAP) for the ActivityNet.



(a) Action: *Yatching*



(b) Action: *Having Massage*

Figure 3. Visualization of our *aLSTM-VGG* (fc8) models attention weights, attending on object, scene, and action for ActivityNet dataset. The lighter video frames have higher attention weight and darker frames represent lower attention weight.

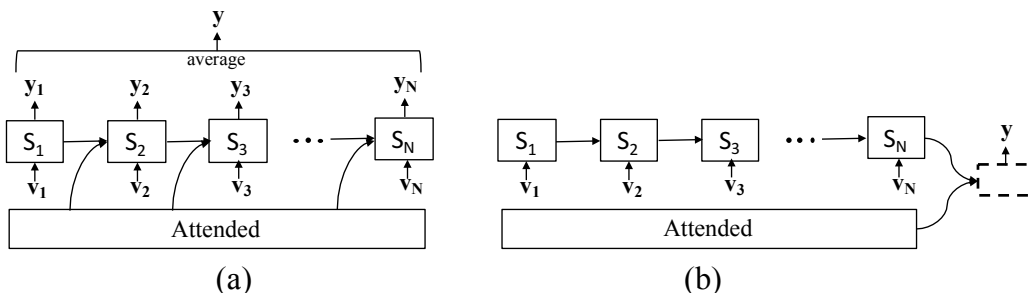


Figure 4. Two variations of our attention-enabled LSTM architectures.

*C3D(fc8)* shows results for 3D CNN model softmax classification directly. *LSTM-C3D* (FC6) is vanilla LSTM model using *C3D* model fc6 layer representation as input-data. *aLSTM* stands for our attention-based LSTM model. We have evaluated our model for all different possible variations of attended-data using VGG model (note that

for all three variations input-data is *C3D(fc8)* features): 1) *aLSTM - VGG<sub>obj</sub>*, attends on objects (*VGG<sub>obj</sub>*), 2) *aLSTM - VGG<sub>sce</sub>*, attends on scenes (*VGG<sub>sce</sub>*), and 3) *aLSTM - VGG<sub>act</sub>* attends on actions (*VGG<sub>act</sub>*). For each of these three models we tried both fc6 and fc8 layer representations. A few observations can be made based on Table

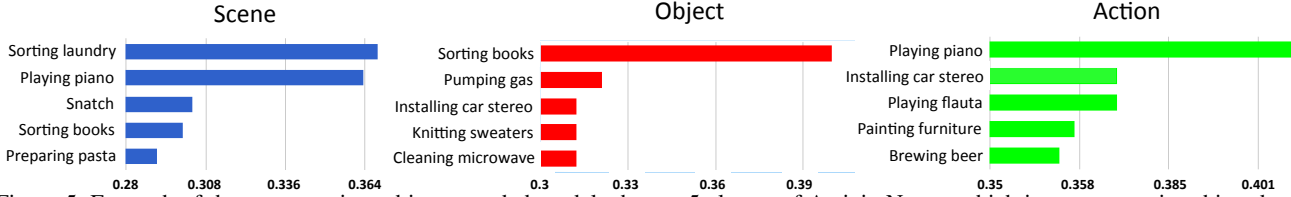


Figure 5. For each of the scene, action, object attended models the top 5 classes of ActivityNet on which improvement is achieved with respect to vanilla LSTM are shown. Note that on some classes improvement is as much as 40% in classification.

4. First, LSTM on top of  $C3D$  features has minor improvement over  $C3D$  CNN for ActivityNet. All  $aLSTM$  models outperform vanilla LSTM and CNN models. Specifically,  $aLSTM$  on  $C3D$  improves the performance by approximately 2 percentage points for ActivityNet. Further for ActivityNet, attending on (contextual) complementary trained VGG19 models for scene, object, and action improves the performance by a large margin, between 4 – 5 percentage points for the same fc6 features; using fc8 features is uniformly better, even further improving performance.

Among individual models, the best performance is achieved by  $aLSTM-VGG_{act}$  (fc8) model for ActivityNet (likely because  $VGG_{act}$  has been fine-tuned for ActivityNet), which is consistent with recent results that show importance of objects [13].

Figure 4 shows two possible variants of integrating attention with LSTM for action classification. In the first variant, Figure 4 (a), attention scores are computed at every time step  $t$  and attended data is distributed into LSTM state  $S_{t-1}$  to compute state  $S_t$ . Also, class scores  $y_t$  are computed at every time step and final score  $y$  is the average of all  $y_t$ s. This is the model we focused on in the approach section. The second variant, Figure 4 (b), has a somewhat different architecture. In particular, attention score is only computed and integrated into the last LSTM state  $S_N$ . The classification score  $y$  is then computed directly (once) at the end. We compare both architectures, using identical setup, and we observed architecture (a) that we used for our experiment has better performance.

#### 4.4. Comparison With Competitive Models

Table 5 shows comparison of our model to state-of-the-art on both ActivityNet. In [4],  $DF$  stand for Deep features,  $MF$  stands for motion features, and  $SF$  stands for spatial features (see [4] for more details).  $SVM - DF$  (fc8) is one-vs-all SVM model over  $DF$  fc8 layer representation described in [4].

For comparison, we have combined different variations of our models described above. Model combination has been implemented by averaging the softmax output scores over multiple  $aLSTM$  models. For example,  $aLSTM-VGG_{act+obj}$  stands for combining  $aLSTM-VGG_{act}$  and  $aLSTM-VGG_{obj}$  and so on. Also  $SVM_{MF+DF}$  stands for SVM model over  $MF$  and  $DF$  feature concatenation described in [4]. Table 5 shows that for ActivityNet our com-

Table 2. ActivityNet: Comparative results.

Model: CNN			
Features		AC	mAP
CNN- $C3D$	(fc8)	40.9%	40.0%
SVM- $DF$ [4]	(fc8)	-	38.1%
Model: Combined			
Features		AC	mAP
SVM $_{MF+DF}$ [4]		-	49.5%
SVM $_{MF+SF}$ [4]		-	48.9%
SVM $_{MF+DF+SF}$ [4]		-	50.5%
$aLSTM-VGG_{act+obj}$	(fc6)	50.0%	51.0%
$aLSTM-VGG_{act+obj}$	(fc8)	50.7%	51.4%
$aLSTM-VGG_{act+sce}$	(fc6)	50.4%	50.7%
$aLSTM-VGG_{act+sce}$	(fc8)	50.7%	51.2%
$aLSTM-VGG_{act+obj+sce}$	(fc6)	51.4%	52.4%
$aLSTM-VGG_{act+obj+sce}$	(fc8)	51.5%	52.8%
$aLSTM-VGG_{act+obj+sce}$	(all)	<b>52.8%</b>	<b>54.2%</b>

bin models improve both AC and mAPs up to  $\sim 11.5\%$  (5.6 percentage points) compared to our best individual model  $aLSTM-VGG_{act}$  (fc8). This reveals that variations of  $aLSTM$  that attend on *object*, *scene*, and *action* contextual data to some extent have complementary information, where the combined model  $aLSTM-VGG_{act+obj+sce}$  (all) has the best performance.

Figure 5 shows, the top 5 classes of ActivityNet with highest accuracy improvement with respect to vanilla LSTM- $C3D$  for  $aLSTM-VGG_{sce}$  (fc8),  $aLSTM-VGG_{obj}$  (fc8), and  $aLSTM-VGG_{act}$  (fc8) models respectively. For *object* a number of classes (e.g., *pumping gas* and *installing car stereo*) that contain easily detectable objects, like a *car*, are most significantly improved (by over 30%).

#### 4.5. Conclusions

We have presented  $aLSTM$  model for action recognition. With this model, we are able to attend on and highlight most relevant frames and perform fine-grained action classification simultaneously. We have shown that by attending on features extracted from complementary trained VGG19 models for scene, object, and action (contextual data), we get a significant gain compared to  $aLSTM$  on  $C3D$ . Furthermore, our model variations that attend on *object*, *scene*, and *action* are complementary which allows their combina-



tion to have the best performance.

## References

- [1] J. Aggarwal and M. Ryoo. Human activities analysis: A review. *ACM Computing Surveys*, 2011.
- [2] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. *CoRR*, abs/1409.0473, 2014.
- [3] F. Bastien, P. Lamblin, R. Pascanu, J. Bergstra, I. J. Goodfellow, A. Bergeron, N. Bouchard, and Y. Bengio. Theano: new features and speed improvements. In *NIPS*, 2012.
- [4] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *CVPR*, 2015.
- [5] J. Donahue, L. A. Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *CVPR*, 2015.
- [6] T. Duong, H. Bui, D. Phung, and S. Venkatesh. Activity recognition and abnormality detection with the switching hidden semi-markov model. In *CVPR*, 2005.
- [7] B. Fernando, E. Gavves, J. M. Oramas, A. Ghodrati, and T. Tuytelaars. Modeling video evolution for action recognition. In *CVPR*, 2015.
- [8] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014.
- [9] A. Graves. Generating sequences with recurrent neural networks. *CoRR*, abs/1308.0850, 2013.
- [10] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computing*, 1997.
- [11] J. Hoffman, S. Guadarrama, E. Tzeng, J. Donahue, R. B. Girshick, T. Darrell, and K. Saenko. LSDA: large scale detection through adaptation. In *NIPS*, 2014.
- [12] N. Ikizler and D. Forsyth. Searching video for complex activities with finite state models. In *CVPR*, 2007.
- [13] M. Jain, J. C. van Gemert, and C. G. M. Snoek. What do 15,000 object categories tell us about classifying and localizing actions? In *CVPR*, 2015.
- [14] S. Ji, W. Xu, M. Yang, and K. Yu. 3d convolutional neural networks for human action recognition. *IEEE PAMI*, 2013.
- [15] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- [16] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. In *CVPR*, 2014.
- [17] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [18] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre. Hmdb: a large video database for human motion recognition. In *ICCV*, 2011.
- [19] Z.-Z. Lan, M. Lin, X. Li, A. G. Hauptmann, and B. Raj. Beyond gaussian pyramid: Multi-skip feature stacking for action recognition. In *CVPR*, 2015.
- [20] I. Laptev, M. Marszalek, C. Schmid, and B. Rozenfeld. Learning realistic human actions from movies. In *CVPR*, 2008.
- [21] H. Larochelle and G. E. Hinton. Learning to combine foveal glimpses with a third-order boltzmann machine. In *NIPS*, 2010.
- [22] S. Lazebnik, C. Schmid, and J. Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *CVPR*, 2006.
- [23] L. Liu, L. Shao, and P. Rockett. Boosted key-frame selection and correlated pyramidal motion-feature representation for human action recognition. *Pattern Recognition*, 2012.
- [24] M. Marszałek, I. Laptev, and C. Schmid. Actions in context. In *IEEE Conference on Computer Vision & Pattern Recognition*, 2009.
- [25] S. Mathe and C. Sminchisescu. Dynamic eye movement datasets and learnt saliency models for visual action recognition. In *ECCV*, 2012.
- [26] J. Y. Ng, M. J. Hausknecht, S. Vijayanarasimhan, O. Vinyals, R. Monga, and G. Toderici. Beyond short snippets: Deep networks for video classification. In *CVPR*, 2015.
- [27] R. Poppe. A survey on vision-based human action recognition. *IVC*, 28(6), June 2010.
- [28] M. Raptis and L. Sigal. Poselet key-framing: A model for human activity recognition. In *CVPR*, 2013.
- [29] K. Schindler and L. V. Gool. Action snippets: How many frames does human action recognition require? In *CVPR*, 2008.
- [30] C. Schldt, I. Laptev, and B. Caputo. Recognizing human actions: A local svm approach. In *ICPR*, 2004.
- [31] N. Shapovalova, M. Raptis, L. Sigal, and G. Mori. Action is in the eye of the beholder: Eye-gaze driven model for spatio-temporal action localization. In *NIPS*, 2013.
- [32] S. Sharma, R. Kiros, and R. Salakhutdinov. Action recognition using visual attention. *arXiv preprint arXiv:1511.04119*, 2015.
- [33] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. In *NIPS*, 2014.
- [34] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [35] C. Sminchisescu, A. Kanaujia, Z. Li, and D. Metaxas. Conditional models for contextual human motion recognition. In *ICCV*, 2005.
- [36] k. Soomro, A. Roshan Zamir, and M. Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. In *CRCV-TR-12-01*, 2012.
- [37] N. Srivastava, E. Mansimov, and R. Salakhutdinov. Unsupervised learning of video representations using lstms. In *ICML*, 2015.
- [38] L. Sun, K. Jia, T. H. Chan, Y. Fang, G. Wang, and S. Yan. Dlsfa: Deeply-learned slow feature analysis for action recognition. In *CVPR*, 2014.
- [39] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015.
- [40] D. Tran, L. D. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features with 3d convolutional networks. In *ICCV*, 2015.

- [41] B. van Merriënboer, D. Bahdanau, V. Dumoulin, D. Serdyuk, D. Warde-Farley, J. Chorowski, and Y. Bengio. Blocks and fuel: Frameworks for deep learning. *CoRR*, abs/1506.00619, 2015.
- [42] S. Venugopalan, M. Rohrbach, J. Donahue, R. J. Mooney, T. Darrell, and K. Saenko. Sequence to sequence - video to text. In *ICCV*, 2015.
- [43] E. Vig, M. Dorr, and D. Cox. Space-variant descriptor sampling for action recognition based on saliency and eye movements. In *ECCV*, 2012.
- [44] H. Wang, A. Kläser, C. Schmid, and L. Cheng-Lin. Action recognition by dense trajectories. In *CVPR*, 2011.
- [45] H. Wang and C. Schmid. Action recognition with improved trajectories. In *ICCV*, 2013.
- [46] Y. Wang and G. Mori. Hidden part models for human action recognition: Probabilistic vs. max-margin. *IEEE PAMI*, 2010.
- [47] Z. Wu, X. Wang, Y. Jiang, H. Ye, and X. Xue. Modeling spatial-temporal clues in a hybrid deep learning framework for video classification. In *ACM Multimedia*, 2015.
- [48] K. Xu, J. Ba, R. Kiros, K. Cho, A. C. Courville, R. Salakhutdinov, R. S. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, 2015.
- [49] L. Yao, A. Torabi, K. Cho, N. Ballas, C. Pal, H. Larochelle, and A. Courville. Describing videos by exploiting temporal structure. In *ICCV*, 2015.
- [50] S. Yeung, O. Russakovsky, N. Jin, M. Andriluka, G. Mori, and L. Fei-Fei. Every moment counts: Dense detailed labeling of actions in complex videos, 2015.
- [51] S. Zha, F. Luisier, W. Andrews, N. Srivastava, and R. Salakhutdinov. Exploiting image-trained CNN architectures for unconstrained video classification. In *BMVC*, 2015.
- [52] Z. Zhao and A. M. Elgammal. Information theoretic key frame selection for action recognition. In *BMVC*, 2008.
- [53] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva. Learning deep features for scene recognition using places database. In *NIPS*, 2014.

## A. Supplementary Results

Table 3 shows average accuracy (AC), mean average precision (mAP) for ActivityNet test data. In table 3, one could observe that our *aLSTM* models attending on  $VGG_{obj}$ ,  $VGG_{sce}$ , and  $VGG_{act}$  features perform better than vanilla LSTM models for the corresponding  $VGG$  features; comparison to LSTM model with  $C3D$  features is in the main paper.

We also use Hollywood2 dataset in our experiments. Hollywood2 (HOHA2) [24] is a human action dataset with 12 classes. HOHA2 includes 1707 video clips from Hollywood movies. Dataset is trimmed, but number of samples per class is not fixed and varies between 24 and 135. HOHA2 consists of 823 training samples and 884 test samples, and each clip can have multiple labels. The much smaller training set makes training LSTM-based models that much more challenging, as is also witnessed in [32]. For HOHA2, we use multilabel classification loss (similar to [50]). We initialize LSTM parameters with the pre-trained LSTM model trained for ActivityNet. We then fine-tune the model for 2,500 iterations on minibatches of size 100. For first 1,000 iterations, we use fixed learning rate of 0.001 and momentum of 0.6; after that, learning rate is decreased every iterations by a fixed ratio. Similar to ActivityNet model training, we also perform an on-the-fly temporal cropping of training videos. For evaluation, we report mean Average Precision (mAP) over all 12 classes.

Figure 4 shows two possible variants of integrating attention with LSTM for action classification. In the first variant, Figure 4 (a), attention scores are computed at every time step  $t$  and attended data is distributed into LSTM state  $S_{t-1}$  to compute state  $S_t$ . Also, class scores  $y_t$  are computed at every time step and final score  $y$  is the average of all  $y_t$ s. This is the model we focused on in the approach section. The second variant, Figure 4 (b), has a somewhat different architecture. In particular, attention score is only computed and integrated into the last LSTM state  $S_N$ . The classification score  $y$  is then computed directly (once) at the end. We compare both architectures, using identical setup, on HOHA2. As can be seen from Table 4 (*aLSTM-VGG<sub>obj</sub>* (fc8) vs. *aLSTM-VGG<sub>obj</sub>* (fc8)(b)), variant (b) has 2.2 percentage points lower mAP compared to model (a). We also evaluated significance of attended-data size, comparing  $T = 12$  and  $T = 24$ . We only observed a slight improvement of 0.2 mAP for HOHA2 using  $T = 24$ . Based on these observations, we use architecture (a) with  $T = 12$  from now on.

Table 4 shows mean average precision over all classes (mAP) for Hollywood2 (HOHA2) test data.  $C3D(fc8)$  shows results for 3D CNN model softmax classification directly. LSTM- $C3D$  (FC6) is vanilla LSTM model using  $C3D$  model fc6 layer representation as input-data. *aLSTM* stands for our attention-based LSTM model. We have eval-

Table 3. ActivityNet: Comparative results of vanilla LSTM and our attention-based LSTM.

Model: LSTM		
Features	AC	mAP
LSTM- $VGG_{obj}$ (fc6)	42.2%	47.6%
<i>aLSTM-VGG<sub>obj</sub></i> (fc6)	<b>45.6%</b>	<b>51.2%</b>
LSTM- $VGG_{sce}$ (fc6)	39.0%	44.4%
<i>aLSTM-VGG<sub>sce</sub></i> (fc6)	<b>44.0%</b>	<b>49.5%</b>
LSTM- $VGG_{act}$ (fc6)	45.9%	51.1%
<i>aLSTM-VGG<sub>act</sub></i> (fc6)	<b>47.0%</b>	<b>52.4%</b>

Table 4. HOHA2: Our model action recognition results.

Model: CNN		
Features		mAP
$C3D$	(fc8)	41.4%
Model: LSTM		
Features		mAP
LSTM- $C3D$	(fc6)	39.9%
<i>aLSTM-C3D</i>	(fc8)	44.4%
<i>aLSTM-VGG<sub>obj</sub></i>	(fc6)	43.6%
<i>aLSTM-VGG<sub>obj</sub></i>	(fc8)	<b>46.3%</b>
<i>aLSTM-VGG<sub>obj</sub></i>	(fc8)(b)	44.1%
<i>aLSTM-VGG<sub>sce</sub></i>	(fc6)	43.6%
<i>aLSTM-VGG<sub>sce</sub></i>	(fc8)	41.6%
<i>aLSTM-VGG<sub>act</sub></i>	(fc6)	44.8%
<i>aLSTM-VGG<sub>act</sub></i>	(fc8)	45.3%

uated our model for different variations of attended-data: 1) *aLSTM - VGG<sub>obj</sub>*, attends on objects ( $VGG_{obj}$ ), 2) *aLSTM - VGG<sub>sce</sub>*, attends on scenes ( $VGG_{sce}$ ), and 3) *aLSTM - VGG<sub>act</sub>* attends on actions ( $VGG_{act}$ ). For each of these three models we tried both fc6 and fc8 layer representations.

A few observations can be made based on Table 4. First, LSTM on top of  $C3D$  features has slightly lower performance for HOHA2 (according to mAP). All *aLSTM* models outperform vanilla LSTM and CNN models. Specifically, *aLSTM* on  $C3D$  improves the performance by approximately 4 percentage points for over vanilla LSTM. Further attending on (contextual) complementary trained VGG19 models for object and action improve the performance. However, attending on trained VGG19 scene model has lower performance compared to attending on  $C3D$  directly. Among individual models, the best performance is achieved by *aLSTM-VGG<sub>obj</sub>* (fc8) which is consistent with recent results that show importance of objects [13].

Table 5 shows comparison of our model to state-of-the-art on HOHA2. For comparison, we have combined different variations of our models described above. Model combination has been implemented by averaging the softmax output scores over multiple *aLSTM* models. For exam-

Table 5. ActivityNet: Comparative results.

<b>Model: CNN</b>		
Features		mAP
CNN- <i>C3D</i>	(fc8)	41.4%
<b>Model: RGB data and/or Hand-Crafted Features</b>		
Features		mAP
Soft attention model [32]		43.9%
DL-SFA[38] [4]		48.1%
$aLSTM-VGG_{act+obj}$	(fc6)	46.2%
$aLSTM-VGG_{act+obj}$	(fc8)	48.5%
$aLSTM-VGG_{act+obj}$	(all)	<b>48.7%</b>
$aLSTM-VGG_{act+sce}$	(fc6)	46.1%
$aLSTM-VGG_{act+sce}$	(fc8)	45.2%
$aLSTM-VGG_{act+obj+sce}$	(fc6)	46.6%
$aLSTM-VGG_{act+obj+sce}$	(fc8)	47.9%
$aLSTM-VGG_{act+obj+sce}$	(all)	48.1%
<b>Model: Hand-Crafted Features</b>		
Features		HOHA2 (mAP)
Multi-skip Feature [19]		68%
VideoDarwin [7]		<b>73.7%</b>

ple,  $aLSTM-VGG_{act+obj}$  stands for combining  $aLSTM-VGG_{act}$  and  $aLSTM-VGG_{obj}$  and so on. Table 5 shows that our combined model for action and object improves the performance, however, adding scene model does not improve compared to our best individual model  $aLSTM-VGG_{act}$  (fc8). Further, we compare our results to *DL – SFA* [38] and, more importantly, to ArXiv results of soft-attention model of [32], which proposes a spatial 3-layer attention-based LSTM. Our best combined model  $aLSTM-VGG_{act+obj}$  (all), outperforms [32] by 11% (or 5 percentage points) and has highest performance compared to methods that use RGB data directly; as opposed to working with combination of optical flow and RGB data [19, 7].