

Spatio-Temporal Attention Models for Grounded Video Captioning

Mihai Zanfir^{2*}, Elisabeta Marinoiu^{2*}, Cristian Sminchisescu^{1,2}
mihai.zanfir@imar.ro, elisabeta.marinoiu@imar.ro,
cristian.sminchisescu@math.lth.se

¹Department of Mathematics, Faculty of Engineering, Lund University

²Institute of Mathematics of the Romanian Academy

Abstract. Automatic video captioning is challenging due to the complex interactions in dynamic real scenes. A comprehensive system would ultimately localize and track the objects, actions and interactions present in a video and generate a description that relies on temporal localization in order to ground the visual concepts. However, most existing automatic video captioning systems map from raw video data to high level textual description, bypassing localization and recognition, thus discarding potentially valuable information for content localization and generalization. In this work we present an automatic video captioning model that combines spatio-temporal attention and image classification by means of deep neural network structures based on long short-term memory. The resulting system is demonstrated to produce state-of-the-art results in the standard YouTube captioning benchmark while also offering the advantage of localizing the visual concepts (subjects, verbs, objects), with no grounding supervision, over space and time.

1 Introduction

In this work, we consider the problem of automatic video captioning, where given an input video, a learned model should describe its content with one or more sentences. This is important considering the increasing rate at which multimedia content is uploaded on the Internet, which in turn requires automatic understanding and description for the retrieval of meaningful content. Automatic video captioning would also be beneficial to human-computer interaction, surveillance and monitoring, and as an aid to the blind and visually-impaired.

However, the problem of translating from the visual domain to a textual one is challenging, as it ideally involves understanding the key actors, objects and their interaction in the scene, followed by the construction of a both semantically and grammatically correct natural language description. The first part is made difficult by the large number of semantic categories, which exhibit a high inter-class variability. Objects may be of different sizes, shapes and colors, or only partially visible. The lack of available, rich annotated data, and the absence

* Authors contributed equally.

of localization information for the key elements present in a video makes the problem even harder. Data is difficult to collect due to the tedious and time-consuming process of annotating individual video frames. Progress has been made in the related problem of automatic image captioning, where annotations are plentiful [1,2]. Even with a complete understanding of the video content, there still remains the problem of delivering a sufficiently relevant digest, at different levels of abstraction, required by specific tasks. A bird enthusiast may require a specific video to be described as "Sudan golden sparrows are bathing in water", whereas a regular person may be satisfied with a simpler description like "Birds playing in water".

Very recent work has focused on attention mechanisms that ground textual elements into the video timeline [3,4]. However, the visual domain is usually represented only at a coarse frame level without explicitly revealing the spatio-temporal structure pertaining to a textual element. We believe that a video captioning model can benefit from the work in spatio-temporal segmentation in video [5,6,7], that could offer plausible proposals for localization of the textual elements. Our contributions can be listed as follows: 1) an attention mechanism that links textual elements to spatio-temporal object proposals and is able to provide localized visual support of the words in the generated sentence, with no grounding supervision, 2) integration of high-level semantic representations obtained both from classifiers learned on YouTube dataset [8] (subjects, verbs, objects) and from pre-trained models with state-of-the-art recurrent neural networks, and 3) competitive or better than state-of-the-art results on three different metrics on the challenging YouTube video description dataset. An overview of our model is given in figure 1. Illustrations of the detailed textual and visual output produced by our model appear in figure 5.

2 Related Work

Previous approaches to video and image captioning follow broadly two main lines of work: 1) intermediate concept prediction in the form of subject, verb, object or place (S,V,O,P) followed by a template sentence generation step, or 2) full sentence generation using recurrent neural networks, mainly long short-term memory units (LSTMs).

Concept Discovery. Earlier work on video captioning has attempted to first detect a subject, a verb and an object for each video and then form a sentence using a template model and a learned language model. In [9] the authors first mine (S,V,O) triplets from the video descriptions. Then, they learn a semantic hierarchy over subjects, verbs and objects and use a multi-channel SVM to predict an (S,V,O) triplet over the learned hierarchies by trading-off specificity and semantic similarity. Once an (S,V,O) triplet is obtained for each video, a sentence is formed using a template-based approach. In [10] a factor graph model is proposed to combine visual detections with language statistics in order to learn an (S,V,O,P) tuple for each video. The sentence generation step is similar to that of [9]. In [11] the authors propose a framework to jointly model language

and vision. The language model is a compositional one, learned over (S,V,O) triplets while the vision model is a two-layer neural network built on top of deep features. Treating concept discovery and sentence generation in separate steps has the advantage of solving two potentially easier and better specified problems instead of a harder, less constrained one. However, the downside is that the resulting sentences can be rigid and may fail to capture the richness of human descriptions.

Recurrent Networks for Image and Video Captioning. Inspired by the recent success of recurrent neural networks in automatic language translation [12,13], a series of papers made use of similar models in "translating" from a visual input to a textual output where the visual information is usually encoded using convolutional neural networks (CNN). In the case of image captioning, significant progress has been made in recent work [14,15,16,17]. The authors in [15] use an attention mechanism on top of a CNN and extract features from a lower convolutional layer in order to obtain correspondences between the feature vectors and regions of the 2D image.[16] use external region proposals and learn an alignment model between image regions and sentences. In [17], using a rich annotated dataset [2] of image regions and corresponding textual descriptions, the problem of localization and description is addressed jointly. They propose a fully convolutional localization network (FCLN) for dense captioning.

In video captioning, the authors of [18] use a stack of two recurrent sequence models (LSTMs) to tackle the problems of activity recognition, image and video description. For image description, features extracted from a pre-trained CNN are fed directly into the LSTMs, while in the case of video description, first, a CRF model is used to obtain a distribution over subjects, verbs and objects. Then, the CRF responses are fed to the LSTMs to form a full sentence. Similarly, [19] uses a two stack LSTMs model for video description, but the visual information is encoded as mean-pooled CNN feature over the video frames. They also show improvement by transfer learning from the image domain (where more training data is available) to video. Approaches have also been pursued using semantic classifiers responses for subjects, verbs and places in combination with LSTMs[20]. Despite the fact that such approaches have achieved a significant improvement (under BLEU and METEOR metrics) compared to previous template based approaches, they do not fully exploit the underlying structure of the video, nor do they attempt to explicitly identify (localize) the main actors or objects that correspond to the textual descriptions produced. [4] incorporate a spatio-temporal 3D CNN trained on video action recognition and use a temporal attention mechanism to select the most relevant temporal segments. [3] are interested in generating multiple sentences and discuss temporal and spatial attention mechanisms. The spatial elements are obtained by sampling image patches around a central actor, on datasets where this assumption holds. They use deep convolutional features like VGG [21] and C3D [22] to represent an image frame. The purpose of these frame-level attention approaches is to guide the model towards different frames of the video at each time step (when a word is produced). Unlike these, our attention model focuses on a pool of spatio-

temporal proposals and learns to choose the best spatial-temporal support for every word of the sentence.

3 Methodology

Our approach to video captioning has two main components: first revealing the spatio-temporal visual support of words in video and then guiding the sentence generation process by including semantic information in the learning process. We integrate a soft-attention mechanism, operating over a pool of spatio-temporal proposals, into a state-of-the-art recurrent network. The joint model learns to produce semantically meaningful sentences while attending to different parts of the video. The semantic information is obtained in two ways: (a) by learning to predict subjects, verbs and objects (S,V,O) and (b) by using pre-trained state-of-the-art image classification and object detection models. An overview of our modeling and computational pipeline is shown in figure 1. Section §3.1 briefly introduces the recurrent model while section §3.2 describes the attention mechanism. Learning the semantic concepts is explained in section §3.3 and the experimental details and results are given in section §4.

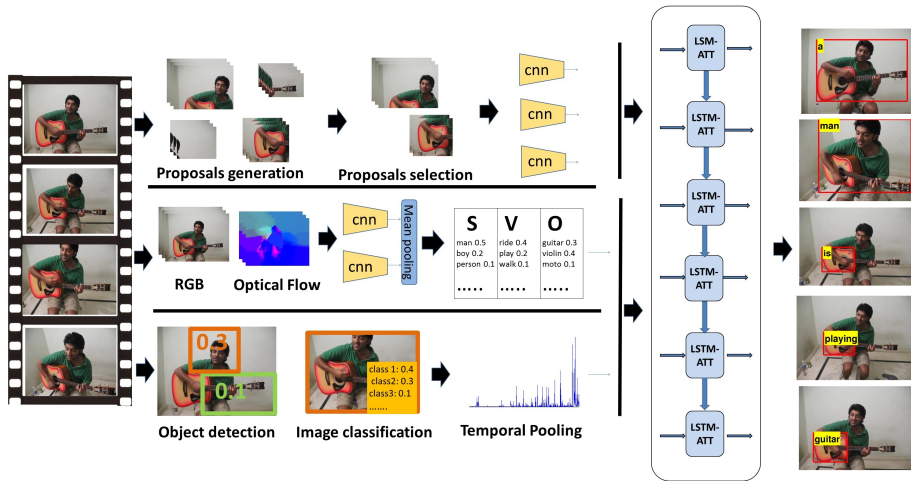


Fig. 1. Overview of our approach. We build an attention mechanism on top of spatio-temporal object proposals and integrate it with state-of-the-art image classifiers, object detectors and recurrent neural networks (LSTM). The image classifiers together with learned high-level semantic features in the form of (Subject, Verb, Object) are provided as contextual features to the Attention-LSTM. Our model is able to visually ground each of the words from the sentence it generates, spatially and temporally, with no additional supervision.

3.1 Recurrent Networks for Video Captioning

Recurrent Neural Networks (RNNs) make use of sequential information and learn temporal dynamics by mapping a sequence of inputs to hidden states and then learn to decode the hidden states into a series of outputs. Their major drawback, however, is the *vanishing gradient* which makes it difficult to learn long-range dependencies that exist in the input sequences [23]. A solution to this issue is to incorporate explicit unit memories, controlled by gates deciding at each step which information should be passed on and which one should be forgotten. Those units, known as long short-term memory (LSTM) units [24], have proven to perform well for machine translation [12] and have recently been used for both image [15,18] and video captioning[25,26,19]. A schema of the LSTM unit (introduced in [27]) used in our experiments is shown in figure 2. The LSTM unit consists of a memory cell, c_t , that encodes the information transmitted from previous units up to current step and *gates* deciding how the information in the memory cell is updated and what the output should be. The input i_t , forget f_t , and output o_t gates are sigmoid functions that decide how much to consider from the current input (u_t), what to retain from the previous cell memory (c_{t-1}) and how much information from the memory cell to be transferred to the hidden state (h_t). The updates at time step t given textual input u_t , visual input z_t , the previous hidden state h_t and the previous memory cell c_{t-1} are given by the following equations:

$$\begin{aligned}
 i_t &= \sigma(W_{xi}u_t + W_{hi}h_{t-1} + W_{zi}z_t + b_i) \\
 f_t &= \sigma(W_{xf}u_t + W_{hf}h_{t-1} + W_{zf}z_t + b_f) \\
 o_t &= \sigma(W_{xo}u_t + W_{ho}h_{t-1} + W_{zo}z_t + b_o) \\
 g_t &= \phi(W_{xg}u_t + W_{hg}h_{t-1} + W_{zg}z_t + b_g) \\
 c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\
 h_t &= o_t \odot \phi(c_t)
 \end{aligned}
 \tag{1}$$

3.2 Attention-based LSTM

Soft-Attention Mechanism. We incorporate a soft-attention mechanism into the LSTM in order to allow the model to selectively focus on different parts of the video each time it produces a word. Inspired by the attention mechanism that exploits the spatial layout of an image [15], a few recent methods have attempted to exploit the temporal structure of the video by learning how to assign different weights to frames in a video sequence [4,3]. These methods, however, do not localize objects in images, and thus the same frame can offer support to very different words in the output. This approach can work when a video selectively focuses on individual objects at a time and thus in a single frame very few objects of interest are present. However, in most videos, there are multiple actors and objects present in a frame. To address this problem, we allow the sequence model (LSTM) to choose where to focus among a pool of

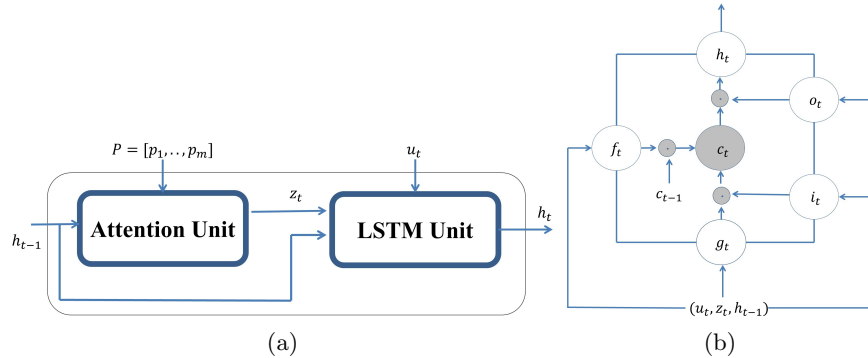


Fig. 2. a) Integration of an attention mechanism with LSTM in our model: at each timestep t , given the previous hidden state h_{t-1} and m spatio-temporal proposals, the Attention Unit outputs a weighted mean-pooled visual feature z_t . The LSTM Unit receives the visual feature z_t , the previous hidden state h_{t-1} and the previously generated word u_t , and outputs the current hidden state h_t . b) Schema of LSTM Unit described in equation 1.

coherent spatio-temporal proposals at each time step t . Thus, the model is able to indicate what is the *localized* visual support used to produce a particular word from the video description.

Considering $P = [p_1, p_2, \dots, p_m]$ the temporal feature vector where m is the number of proposals and p_i the descriptor for the i -th proposal, the LSTM learns at each time step a series of m weights β_{ti} such that the final encoding of the visual input is

$$z_t = \sum_{i=1}^m \beta_{ti} p_i, \quad \text{with} \quad \sum_{i=1}^m \beta_{ti} = 1 \quad (2)$$

The weights β_{ti} represent the importance of the i -th proposal for generating the word at the current time step, given the previously generated words. First, for each proposal we learn a score ϵ_{ti} based on its visual feature p_i and the previous hidden state h_{t-1} , which is given by

$$\epsilon_{ti} = W_{ph} \phi(W_p p_i + W_h h_{t-1} + b_{ph}) \quad (3)$$

where ϕ is the hyperbolic tangent and W_{ph}, W_p, W_h, b_{ph} are parameters to be learned. Those scores are then normalized to obtain the β_{ti} weights:

$$\beta_{ti} = \frac{e^{\epsilon_{ti}}}{\sum_{i=1}^m e^{\epsilon_{ti}}} \quad (4)$$

A schematic view of the Attention-LSTM model is shown in figure 2. At testing time, by inspecting the weights in decreasing order, we can interpret the visual information preferred (selected) by the model when producing a particular word in the textual description of the video.

3.3 High-Level Semantic Description

Generating a textual description of a video requires identifying the actors and their interactions and then constructing a grammatically well-formed sentence. For this purpose, in order to generate a human-like textual description of a video, we first represent a video in the form of a Subject(S), a Verb(V) and an Object(O) (similarly to earlier works [9]). We then integrate this representation with state-of-the-art recurrent models, along with spatio-temporal localization processes and object detection and classification information.

SVO representation and vocabulary construction. In order to learn a semantic high-level representation for each video, we represent a sentence in a compact and simplified manner that preserves its main idea by extracting a (S,V,O) tuple - e.g. the sentence *A cat plays with a toy* is represented as (cat, play, toy). We initially used the SVO vocabulary proposed in [9], but found it to be too small (only 45 subjects, 218 verbs and 241 objects) and too semantically restrictive (e.g. no different words for *man* and *woman*, as it only contains *person*). We mine the intermediate concepts differently from [9], such that our vocabulary is richer and with fewer constraints. The important changes we made in the way we build the vocabulary are: 1) considering both the indirect and direct objects when parsing the sentences as opposed to only the direct objects, 2) not grouping words into very general classes, and 3) an S, V or O is added to the final vocabulary if it is mentioned in at least two different sentences in any given video. We use the parser available from [28] to extract from each sentence a subject, verb and an object. Our final vocabulary set is a tuple $\mathcal{D} = \{\mathcal{S}, \mathcal{V}, \mathcal{O}\}$ of the corresponding vocabularies for each sentence part and it has a considerably larger size: 246 subjects, 459 verbs and 801 objects.

SVO Classification. We treat the three vocabularies separately and use Least Squares Support Vector Machine (LS-SVM) as a classifier in a one-vs-all approach. Note that an input video can have multiple labels from each vocabulary (e.g. a video can have labels 'cat', 'animal', 'kitten' for the subject class). We use LS-SVM because it provides a closed form solution both for the leave-one-out prediction and the prediction error via the block inversion lemma [29]. We represent a video as a classifier response vector for all the classes in the combined vocabulary \mathcal{D} . Training videos use the leave-one-out prediction and testing videos use the prediction based on the classifiers learned on the whole training set. This is different from [20], where classifiers scores are learned in the same way for training and testing. The dataset used in [20] has around 56k training examples, with roughly the same vocabulary size as ours. Given that our dataset has a much smaller number of training videos ($\approx 1,300$), we argue that LS-SVM is a better option; we can tune parameters without relying on a separate validation set (further decreasing the amount of label data) and we can better simulate the testing conditions. Our choice is also supported by the increase in classification accuracy when compared to other methods using the same vocabulary (see §4.3).

4 Experimental Details

Dataset Description. We perform our experiments on the YouTube dataset [8] which consists of 1,967 short videos (between 10s and 25s length) collected from YouTube that usually depict only one main activity. Each video has approximately 40 human-generated English descriptions collected through Amazon Mechanical Turk. We use the same splitting into train (1,197 videos), validation (100 videos) and test (670 videos) subsets as previous methods, so that our results are directly comparable to them.

Evaluation Measures. We report our results under BLEU [30] and METEOR [31] metrics which were originally proposed for the evaluation of automatic translation approaches and have also been adopted by previous works in video and image captioning. BLEU@n computes the geometric average of the n-gram precision between generated and reference sentences. METEOR computes an alignment score between sentences by taking into account the exact tokens, the stemmed ones and semantic similarities between them. We use the evaluation software provided by [1] which we adapt to our dataset.

4.1 Spatio-Temporal Object Proposals

We use the method from [6] to gather a pool of spatio-temporal object proposals. We split each video into parts using a shot boundary detection method [32]. Around 1,000 spatio-temporal proposals are extracted separately for each sub-video and together they form the pool of proposals for the whole video. We filter out the proposals that have a small spatial or temporal extent. To diversify the pool and to eliminate very similar proposals, we keep only those that have low IoU scores with each other. From the pool of proposals, we are interested mainly in those that could be attached a semantic meaning. Thus, we sort the proposals according to a semantic measure based on two scores and retain the top m . In our experiments we set $m = 20$. The first score is obtained by running the image classification CNN VGG-19 from [21] (trained on 1.3M images from ImageNet Large-Scale Visual Recognition Challenge (ILSVRC) [33]) on every bounding box of each proposal, retaining the maximum activation in each frame among the 1,000 classes and averaging across all frames. The second score is obtained by running the 20-class object detector from [34] on every frame of each video. For each frame of the proposal, we compute the maximum detection score (multiplied by the IoU between the bounding box of the detection and the spatial extent of the proposal) and then we average the scores across all proposal frames. The final proposal score is the average of the image classification and detection scores.

Given a proposal, for each of its bounding boxes in the video frames, we extract the output of the fc_7 layer of the VGG-19. The feature descriptor for a proposal is obtained by mean-pooling over all the bounding boxes. We represent a video by m such descriptors corresponding to best scoring m spatio-temporal proposals and we refer to this $m \times D$ descriptor as the *temporal visual feature*.

In practice, for some videos, the number of duplicate proposals is very large, and the final number of proposals can be less than m . Since the network takes as input a fixed sized array, for the videos that do not have at least m spatio-temporal proposals we pad the feature matrix to obtain a fixed size descriptor. We mark the padding proposals so that they will be ignored in the learning and testing processes. Examples of the selected proposals can be seen in figure 3.



Fig. 3. Examples of selected object proposals. For each video we generated a large pool of spatio-temporal object proposals and then learned to automatically select those that are most likely to overlap with easily recognisable semantic categories.

4.2 Attention LSTM with Spatio-Temporal Object Proposals

Vocabulary. We use all the words in the training sentences without performing any pre-processing step. This results in a vocabulary of size 9,070. Similarly to previous works, we represent the words as one-hot vectors, set the maximum sentence length to 20 and mark with special characters the beginning and end of a sentence.

Training. We implemented our attention mechanism using the Caffe [35] framework, integrating it on top of the LSTM provided by [18]. We refer to our proposed LSTM model which uses an attention mechanism over spatio-temporal object proposals as LSTM-ATT. The training phase in LSTM-ATT is a sequential one: at each time step t the unit is given the temporal feature P (representing the m proposals), the embedded vector u_t , corresponding to the previous ground-truth word w_{t-1} , and the previous hidden state h_{t-1} . The output hidden-state

h_t (see figure 2) is then used to predict a distribution $P(w_t)$ over the words in the vocabulary. We use the softmax loss and a dropout of 0.5 to avoid over-fitting. We train our models for a maximum of 128 epochs and use the validation set to choose the best iteration for each of the two metrics, BLEU and METEOR.

Inference. Inference is also performed in a sequential manner: given m proposals and the previous emitted word at time $t - 1$, sampled from $P(w_{t-1})$, the model generates the current word until the special character for end of sentence is met or the maximum length of a sentence is reached. We perform the sampling using beam search with beam size 20. Because we noticed that standard beam search implementations sometimes tends to end sentences too early, we modified it to force longer sentences to be produced (at least 4 words).

4.3 Integrating Contextual Semantic Features

Semantic SVO Representation. In order to obtain the SVO responses, we use the LS-SVM described in section §3.3 and consider 3 different classification problems, one for S, one for V and one for O. Each video can then be described by concatenating the responses to these three classification problems. We use different features, depending on the part of sentence we want to classify. For the subject and object classes, we use the VGG-19 of [21], extract feature responses from the fc7 layer for each frame of the video and then perform mean pooling. For the verb class we use two types of motion features: the trajectory features of [36] and the motion-CNN features of [37], again followed by mean pooling. For S, V and O we obtained the following classification accuracies, respectively: 62.5%, 40.9% and 28.30%. Among the 3 classification problems, the results on O are the lowest since the number of classes in the object vocabulary is the largest: 801 compared with 264 for S and 459 for V. Also, the objects have a smaller spatio-temporal extent in video as they usually represent the objects manipulated by a person or animal (e.g *onion*, *ball*, etc). Since we considerably augmented the vocabulary, our classification results on this vocabulary do not compare directly with previous methods. However, we run our method on the initial proposed vocabulary [9] and show results in table 1, against the most common (S,V,O) triple found in human annotations for each video. We show state-of-the-art results for S and O and a slightly lower accuracy than state-of-the-art for V. Notice that the methods marked with (*) generate a full sentence using a recurrent neural network and then extract the S, V and O using a dependency parser. Our aim is to use these intermediate semantic concepts as features to guide and ground our LSTM attention model (LSTM-ATT model) in the sentence generation process.

Semantic Representations. Apart from the SVO responses obtained by training using only the YouTube dataset, we also extract high level semantic features using state-of-the-art image classification and detection models. More precisely, we run the VGG CNN from [21] in each frame of the video and obtain the 1,000 dimensional score vector representing the classification responses over the 1,000

Model	S%	V%	O%
HVC [38]	76.5	22.2	11.9
FGM [38]	76.4	21.3	12.3
JointEmbed [11]	78.2	24.4	11.9
(*) LSTM-E (VGG+C3D) [39]	80.4	29.8	13.8
(*) LSTM-YT-coco [19]	76.0	23.3	14.0
(*) LSTM-YT-coco+flicker [19]	75.61	25.3	12.4
LS-SVM(ours)	83.6	28.1	23.1

Table 1. Binary SVO accuracy computed against the most common (S,V,O) triple provided by humans. Entries marked with (*) first obtain a sentence describing the whole video and then mine the (S,V,O), whereas the others perform a classification over the S, V and O vocabularies.

classes from the ImageNet classification dataset. To obtain a semantic representation of the video, we experimented with both average and max pooling over the frames and noticed that average pooling performs better in our experiments. We also run the 20 class object detector of [34] in each frame of every video and compute a 20-dimensional descriptor. For each class, we retain the detection response scores in every frame then perform temporal pooling across a window of 25 frames. The final score for a class is the maximum of the temporal pooled scores for that class. The temporal pooling ensures that the detection observed is stable and lasts for at least 1 second. The maximum over such detection scores represents the confidence in having seen a particular object in video.

Integration with LSTM-ATT. We have experimented with two methods to integrate the high-level semantic features - SVO classification, object detection (DET) and image classification scores (CLS) - with the LSTM-ATT. In the first method, both the temporal visual features P and semantic features s are provided as input to the LSTM-ATT. In the second one, we stack a LSTM, that processes only the semantic features s , on top of LSTM-ATT which receives the temporal visual feature P as input. We refer to these methods as LSTM-ATT(SEM) and LSTM2-ATT(SEM), respectively, where (SEM) stands for different subsets of semantic features. A schematic view of the two models we use is shown in figure 4.

4.4 Experimental Results

Quantitative Results. Results obtained with the proposed models are shown in table 2. We first check whether the attention mechanism provides an advantage over mean pooling the temporal visual feature, as quantified by the currently most used metrics, BLEU and METEOR. Using a simple LSTM that receives as input the mean pooled temporal visual feature, we obtain a score of 45.4% on BLEU@4 and 31.2% on METEOR. With LSTM-ATT, the results are considerably higher: 48.7% on BLEU@4 and 31.9% on METEOR, which demonstrates

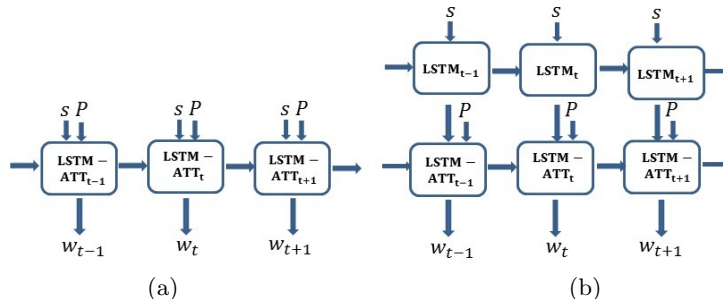


Fig. 4. Integration of semantic features with LSTM-ATT. a) LSTM-ATT(SEM): both semantic features s and temporal features P are processed by the LSTM-ATT unit. b) LSTM2-ATT(SEM): we stack a LSTM unit that processes only the semantic input s on top of the LSTM-ATT.

that the attention mechanism not only provides a way to focus selectively on the input video but also improves results. This is also true when adding semantic features both to the standard LSTM (with mean pooled temporal visual feature) and to the LSTM-ATT.

Our LSTM-ATT model achieves competitive results compared to other methods. Adding semantic features on top of this model improves the state-of-the-art results on the BLEU@n metric, while also performing well on METEOR. We show results with both LSTM-ATT(SEM) and LSTM2-ATT(SEM). The contributions of the SVO semantic features alone and in conjunction with DET and CLS features are also presented. In the case of SVO features alone, the best results are obtained with LSTM2-ATT(SVO) method for both evaluation metrics (BLEU@4 52.0%, METEOR 32.3%), while when using the full semantic features, our best performing method under BLEU is LSTM-ATT(SVO,DET,CLS) (50.6%) and under METEOR is LSTM2-ATT(SVO,DET,CLS) (32.4%).

Qualitative Results. Our attention mechanism, built on top of spatio-temporal object proposals, allows for a *visual explanation* of what the model ranked as the most relevant visual support for emitting a particular word. This can be done by inspecting the learned weights β (see equation 4) and their associated proposals. In figure 5 we show the proposal with the highest associated β weight (a random frame from it) that was used in generating a particular word. For display purposes, we ignore linking words and articles that do not have a visual grounding in video. Our model correctly indicates the localization of the key video description components and the words it emitted, even for those having a very small spatial extent such as *pepper*, *ball*, *toy*, *gun*. There are cases when a single spatio-temporal proposal is chosen as the best visual explanation for multiple words, as it is with (*girl*, *riding*, *horse*). We also show examples when the obtained sentences are wrong (marked with red in figure 5). In some of the cases, our algorithm correctly identifies parts of the sentences - especially

subjects and verbs such as (*man-cutting, dog-playing*) - but fails to find the correct object. This is due to the large variability in objects appearance and size and also depends on the quality of the spatio-temporal proposals pool.

Method	BLEU@1	BLEU@2	BLEU@3	BLEU@4	METEOR
FGM [38]	-	-	-	13.68	23.9
S2VT[25]	-	-	-	-	29.8
MM-VDN[26]	-	-	-	37.64	29.00
LSTM-YT-coco [19]	-	-	-	33.29	29.07
LSTM-YT-coco+flicker [19]	-	-	-	33.29	28.88
Temporal attention [4]	-	-	-	41.92	29.60
LSTM-E (VGG+C3D) [39]	78.8	66.0	55.4	45.3	31.0
h-RNN[3]	81.5	70.4	60.4	49.9	32.6
HRNE with attention[40]	79.2	66.3	55.1	43.8	33.1
GRU-RCNN [41]	-	-	-	49.63	31.70
LSTM	78.0	66.4	56.7	45.4	31.2
LSTM(SVO)	80.1	68.1	57.5	45.8	31.2
LSTM(DET,CLS)	81.2	68.9	57.9	46.2	31.1
LSTM(SVO,DET,CLS)	80.8	69.3	59.3	48.3	30.7
LSTM-ATT	80.1	68.9	59.4	48.7	31.9
LSTM-ATT(SVO)	81.0	70.5	61.2	50.5	32.3
LSTM-ATT(DET,CLS)	81.9	70.9	60.9	50.5	31.6
LSTM-ATT(SVO,DET,CLS)	82.0	71.6	62.4	51.5	32.0
LSTM2-ATT(SVO)	82.4	71.8	62.5	52.0	32.3
LSTM2-ATT(DET,CLS)	80.6	68.1	57.4	46.0	31.8
LSTM2-ATT(SVO,DET,CLS)	81.5	70.8	61.5	50.6	32.4

Table 2. Comparison with previous works on BLEU@1 - BLEU@4 and METEOR metrics. Values are reported as percentage %.

5 Conclusions

In this paper we have addressed some of the challenges in automatic video captioning by aiming to spatio-temporally ground the semantic video concepts, as an intermediate step, without grounding supervision. In contrast to most existing automatic video captioning systems that map from the raw video to the high level textual description, bypassing localization, we aim at aggregating additional, potentially valuable information, by relying on spatio-temporal video proposals and image classification responses for content localization and improved generalization, fused using deep neural network attention models, based on long short-term memory. Our resulting system produces competitive, state-of-the-art results in the standard YouTube captioning benchmark and offers the additional advantage of localizing the concepts (subjects, verbs, objects), with no grounding supervision, over space and time.

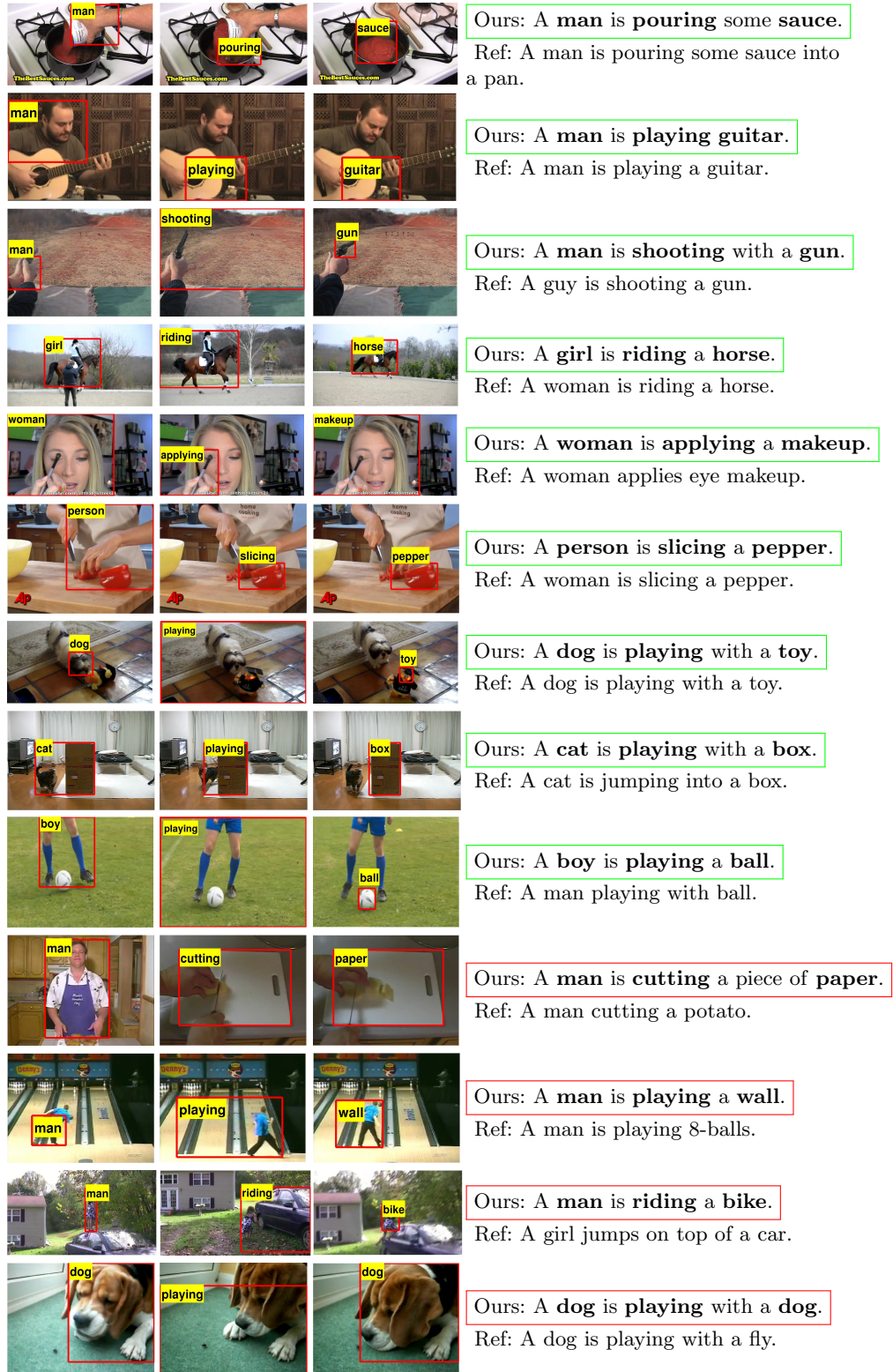


Fig. 5. Highest scoring proposals of our model (according to β eq. 4) for each emitted word in the sentence. We only illustrate the grounding of the main words in the sentence and ignore linking words. The complete sentence is shown in the right column together with the closest reference from the human annotations. For each proposal we show a single, randomly selected frame.

Acknowledgement This work was supported in part by CNCS-UEFISCDI under PCE-2011-3-0438, JRP-RO-FR-2014-16 and NVIDIA through a GPU donation.

References

1. Chen, X., Fang, H., Lin, T., Vedantam, R., Gupta, S., Dollár, P., Zitnick, C.L.: Microsoft COCO captions: Data collection and evaluation server. arXiv preprint arXiv:1504.00325 (2015)
2. Krishna, R., Zhu, Y., Groth, O., Johnson, J., Hata, K., Kravitz, J., Chen, S., Kalantidis, Y., Li, L., Shamma, D.A., Bernstein, M.S., Li, F.: Visual genome: Connecting language and vision using crowdsourced dense image annotations. arXiv preprint arXiv:1602.07332 (2016)
3. Yu, H., Wang, J., Huang, Z., Yang, Y., Xu, W.: Video paragraph captioning using hierarchical recurrent neural networks. In: CVPR. (2016)
4. Yao, L., Torabi, A., Cho, K., Ballas, N., Pal, C., Larochelle, H., Courville, A.: Describing videos by exploiting temporal structure. In: ICCV. (2015)
5. Taralova, E.H., De la Torre, F., Hebert, M.: Motion words for videos. In: ECCV. (2014)
6. Oneata, D., Revaud, J., Verbeek, J., Schmid, C.: Spatio-temporal object detection proposals. In: ECCV. (2014)
7. Fragkiadaki, K., Arbelaez, P., Felsen, P., Malik, J.: Learning to segment moving objects in videos. In: CVPR. (2015)
8. Chen, D.L., Dolan, W.B.: Collecting highly parallel data for paraphrase evaluation. In: ACL. (2011)
9. Guadarrama, S., Krishnamoorthy, N., Malkarnenkar, G., Venugopalan, S., Mooney, R., Darrell, T., Saenko, K.: Youtube2text: Recognizing and describing arbitrary activities using semantic hierarchies and zero-shot recognition. In: ICCV. (2013)
10. Thomason, J., Venugopalan, S., Guadarrama, S., Saenko, K., Mooney, R.: Integrating language and vision to generate natural language descriptions of videos in the wild. In: COLING. (2014)
11. Xu, R., Xiong, C., Chen, W., Corso, J.J.: Jointly modeling deep video and compositional text to bridge vision and language in a unified framework. In: AAAI Conference on Artificial Intelligence. (2015)
12. Sutskever, I., Vinyals, O., Le, Q.V.: Sequence to sequence learning with neural networks. In: NIPS. (2014)
13. Bahdanau, D., Cho, K., Bengio, Y.: Neural Machine Translation by Jointly Learning to Align and Translate. ICLR (2015)
14. Vinyals, O., Toshev, A., Bengio, S., Erhan, D.: Show and tell: A neural image caption generator. In: CVPR. (2015)
15. Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A.C., Salakhutdinov, R., Zemel, R.S., Bengio, Y.: Show, attend and tell: Neural image caption generation with visual attention. In: ICML. (2015)
16. Karpathy, A., Fei-Fei, L.: Deep visual-semantic alignments for generating image descriptions. In: CVPR. (2015)
17. Johnson, J., Karpathy, A., Fei-Fei, L.: Densecap: Fully convolutional localization networks for dense captioning. arXiv preprint arXiv:1511.07571 (2015)
18. Donahue, J., Hendricks, L.A., Guadarrama, S., Rohrbach, M., Venugopalan, S., Saenko, K., Darrell, T.: Long-term recurrent convolutional networks for visual recognition and description. In: CVPR. (2015)

19. Venugopalan, S., Xu, H., Donahue, J., Rohrbach, M., Mooney, R., Saenko, K.: Translating videos to natural language using deep recurrent neural networks. In: NAACL HLT. (2015)
20. Rohrbach, A., Rohrbach, M., Schiele, B.: The long-short story of movie description. In: GCPR. (2015)
21. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: ICLR. (2014)
22. Tran, D., Bourdev, L.D., Fergus, R., Torresani, L., Paluri, M.: Learning spatiotemporal features with 3d convolutional networks. In: ICCV. (2015)
23. Hochreiter, S., Bengio, Y., Frasconi, P., Schmidhuber, J.: Gradient flow in recurrent nets: the difficulty of learning long-term dependencies (2001)
24. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural computation* **9** (1997) 1735–1780
25. Venugopalan, S., Rohrbach, M., Donahue, J., Mooney, R., Darrell, T., Saenko, K.: Sequence to sequence – video to text. In: ICCV. (2015)
26. Xu, H., Venugopalan, S., Ramanishka, V., Rohrbach, M., Saenko, K.: A multi-scale multiple instance video description network. In: arXiv preprint arXiv:1505.05914. (2015)
27. Zaremba, W., Sutskever, I.: Learning to execute. arXiv preprint arXiv:1410.4615 (2014)
28. Manning, C.D., Surdeanu, M., Bauer, J., Finkel, J., Bethard, S.J., McClosky, D.: The Stanford CoreNLP natural language processing toolkit. In: ACL. (2014)
29. Cawley, G.C.: Leave-one-out cross-validation based model selection criteria for weighted ls-svms. In: IJCNN. (2006)
30. Papineni, K., Roukos, S., Ward, T., Zhu, W.J.: Bleu: A method for automatic evaluation of machine translation. In: ACL. (2002)
31. Lavie, A., Agarwal, A.: Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. (2005) 65–72
32. Lienhart, R.W.: Comparison of automatic shot boundary detection algorithms. In: Electronic Imaging'99, International Society for Optics and Photonics (1998) 290–301
33. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. *IJCV* **115** (2015) 211–252
34. Ren, S., He, K., Girshick, R., Sun, J.: Faster R-CNN: Towards real-time object detection with region proposal networks. In: NIPS. (2015)
35. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., Darrell, T.: Caffe: Convolutional architecture for fast feature embedding. In: ACMMM. (2014)
36. Wang, H., Kläser, A., Schmid, C., Liu, C.L.: Action recognition by dense trajectories. In: CVPR, IEEE (2011)
37. Gkioxari, G., Malik, J.: Finding action tubes. In: CVPR. (2015)
38. Thomason, J., Venugopalan, S., Guadarrama, S., Saenko, K., Mooney, R.: Integrating language and vision to generate natural language descriptions of videos in the wild. In: COLING. (2014)
39. Pan, Y., , T.M., Yao, T., Li, H., Rui, Y.: Jointly modeling embedding and translation to bridge video and language. In: arXiv preprint arXiv:1505.01861. (2015)
40. Pan, P., Xu, Z., Yang, Y., Wu, F., Zhuang, Y.: Hierarchical recurrent neural encoder for video representation with application to captioning. arXiv preprint arXiv:1511.03476 (2015)

41. Ballas, N., Yao, L., Pal, C., Courville, A.C.: Delving deeper into convolutional networks for learning video representations. arXiv preprint arXiv:1511.06432 (2015)