

# Exploiting Long-term Temporal Dynamics for Video Captioning

Yuyu Guo · Jingqiu Zhang · Lianli Gao

Received: date / Accepted: date

**Abstract** Automatically describing videos with natural language is a fundamental challenge for computer vision and natural language processing. Recently, progress in this problem has been achieved through two steps: 1) employing 2-D and/or 3-D Convolutional Neural Networks (CNNs) (e.g. VGG, ResNet or C3D) to extract spatial and/or temporal features to encode video contents; and 2) applying Recurrent Neural Networks (RNNs) to generate sentences to describe events in videos. Temporal attention-based model has gained much progress by considering the importance of each video frame. However, for a long video, especially for a video which consists of a set of sub-events, we should discover and leverage the importance of each sub-shot instead of each frame. In this paper, we propose a novel approach, namely temporal and spatial LSTM (TS-LSTM), which systematically exploits spatial and temporal dynamics within video sequences. In TS-LSTM, a temporal pooling LSTM (TP-LSTM) is designed to incorporate both spatial and temporal information to extract long-term temporal dynamics within video sub-shots; and a stacked LSTM is introduced to generate a list of words to describe the video. Experimental results obtained in two public video captioning benchmarks indicate that our TS-LSTM outperforms the state-of-the-art methods.

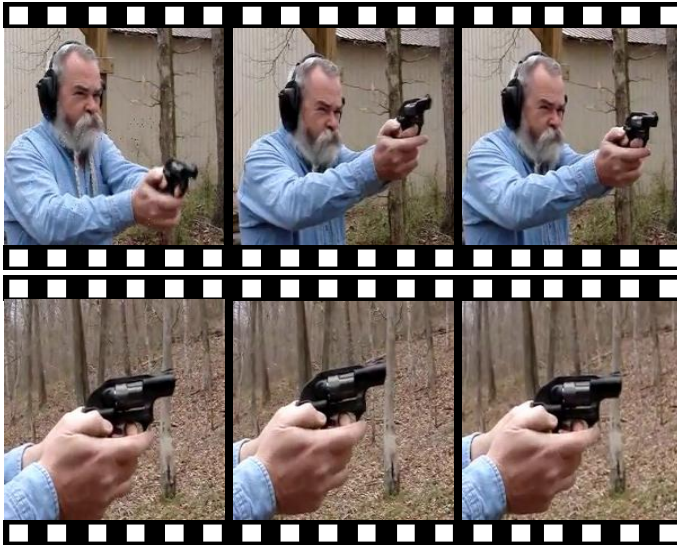
**Keywords** RNNs · Video Captioning · long-term temporal dynamics

---

Lianli Gao  
E-mail: lianli.gao@uestc.edu.cn

Yuyu Guo, Jingqiu Zhang and Lianli Gao are with the Future Media Center and School of Computer Science and Engineering, The University of Electronic Science and Technology of China.

Lianli Gao is the corresponding author



**Fig. 1** We extract six frames from a video, which contains distinct appearance information. However, the top three frames do not contain much variance, so as the bottom three frames.

## 1 Introduction

With the development of multimedia and information technology, huge amounts of videos are uploaded and downloaded on the Internet every day, thus it has spawned a great deal of research into videos or images, such as video/image classification [59], video/image retrieval [58, 42, 39, 60], video segmentation [41, 14], video annotation [13, 40] and video captioning [12, 38] etc. As a bridge between computer vision and natural language, video captioning has become a hot research topic in recent years. Moreover, describing video contents with natural language becomes a key component for improving human-robot interaction and artificial intelligence. To date, extracting features from videos and then translating them into natural language sentences is the main trend, therefore researchers [19, 22, 45, 20] are focusing on solving two sub-problems: 1) how to efficiently extract video features; and 2) how to accurately translate video features into sentences with Recurrent Neural Networks (RNNs).

Specifically, [19, 45] firstly identified video semantic contents and then generated sentences based on some templates. [22, 10, 32] disposed the problem with probabilistic graphical model (PGM) and utilized Markov random field (MRF) or Conditional random field (CRF) to find the relationship between visual content and natural language. Great success in image classification achieved by deep convolutional neural networks (e.g. GoogLeNet [44], VGG

[37] and ResNet [17]), these networks provide researchers with powerful tools to extract image/video features on different fields, such as video captioning [28,27] and video action recognition [24,51]. In general, the basic video captioning framework adopts pre-trained deep CNNs (e.g. ResNet or C3D) to extract spatial and/or temporal features, and then applies an RNN network (e.g. LSTM[35], GRU[6] or their extensions.) to generate words.

Furthermore, Mnih *et al.* [26] proposed a novel recurrent neural network model to extract information from images by adaptively selecting a sequence of regions or locations and only processing the selected regions at high resolution. The experimental results showed that it significantly outperforms the convolutional neural network baseline on a dynamic visual control problem. This strategy is named as visual attention. In fact, the basic idea of the attention mechanism is to selectively focus on the important information and maximumly ignore the unimportant information in the meantime. Therefore, the first step of attention is to estimate which part is important and assign a higher weight to it. Inspired by its great success, a variety of visual attention models have been proposed [38,53]. For example, [38,53] introduced a temporal attention to enhance video captioning by setting frames with different weights to select the most relevant temporal segments by training. However, for a video, the temporal variances are existing between sub-shots instead of adjacent frames. From Fig. 1, we can see that all of the frames are important for describing the video, but many of the frames are duplicate. The top three frames describe a man who is holding a gun, while the bottom three frames describe the shooting action. Weighting each frame would incur excessive computational cost and result in low accuracy.

Here, we argue that long-range temporal structure plays an important role in understanding dynamics in video captioning. However, mainstream video captioning frameworks [38,53] usually focus on appearances and short-term motions, which lack the capacity to incorporate long-range temporal structure. In this paper, we aim to study the following problem: How to design an effective and efficient video-level framework for learning video representation that is able to capture long-range temporal structure for improving video captioning. Moreover, in terms of long-range temporal structure modeling, a key observation is that consecutive frames are highly invariant [51,24] (Fig. 1), thus it is unnecessary to directly set dense temporal sampling for LSTMs. Therefore, we propose a novel framework, namely temporal and spatial LSTM (TS-LSTM), which firstly uses a temporal pooling (TP) layer to keep the temporal invariance in a short video shot, then a Long Short-Term Memory (LSTM) [35] to exploit temporal dynamics between long-range video shots. In addition, a stacked Long Short-Term Memory (Stack-LSTM) is adopted to generate words in the final stage. This framework employs representations from spatial and temporal features to enhance video captioning. The contributions of this paper are as follows:

- Given spatial and motion feature representations over time, we propose to integrate a temporal pooling and a LSTM to learn both temporal invari-

- ance and variance. This mechanism fuses high-level spatial and temporal features to learn long-range temporal dynamics over the whole video.
- We introduce a TS-LSTM video captioning framework, which integrates TP-LSTM with a mean pooling and a stacked LSTM to automatically generate words for describing a video. Specifically, the mean pooling is applied on a concatenation of visual features, motion features and long-term dynamics to extract useful information for the decode process. In addition, inspired by the two-stream framework [36,11] which has achieved great results in video action recognition, we adopt a fine-tuned Resnet-152 [17] to extract the temporal features. Compared with C3D [46] features, using Resnet-152 features can achieve better results.
  - We perform experiments on two video captioning datasets, namely MSVD [4] and MSR-VTT [52], to verify the effectiveness of our method. The experimental results show that our method outperforms existing approaches.

## 2 Related Work

### 2.1 Deep Convolutional Neural Network

In the field of deep learning, deep convolutional neural networks (CNNs) have been widely applied to explore visual information, such as image recognition [21], object detection [31] and image retrieval [42] etc. From LeNet [21] to ResNet [17], the performances of such models have greatly improved on the task of image classification. Specifically, ResNet-152 [17] achieves better results than human beings. As a result, many researchers employ these networks to improve the performance of their tasks. For example, Feichtenhofer *et al.* [36] fine-tuned the VGG [37] to improve the performance on video action recognition task, and Yao *et al.* [53] used a pre-trained GoogLeNet [44] to extract features for video captioning. Motivated by the previous works [23,38], we use the ResNet-152 to extract features both for spatial and temporal information. In addition, all the above mentioned deep CNNs contain pooling layers, which are always used to reduce the spatial size and solve the over-fitting problem. Besides, Scherer *et al.* [34] showed that pooling layers have potential to obtain spatial invariance, thus we integrate a temporal pooling layer to explore the temporal invariance in a video short snip in this paper.

### 2.2 Recurrent Neural Networks

Compared with CNNs, Recurrent Neural Networks (RNNs) are good at modeling sequential data, thus they have been widely utilized in natural language processing and achieved great success [9,18]. At each time step, an RNN observes an element and updates its internal states. In the field of speech recognition, the RNN Language Model (RNNLM) [25] models the output distribution by adding a softmax layer onto the hidden states. In order to learn

the RNNLM model’s parameters, it maximizes the log-likelihood by using the gradient descent method.

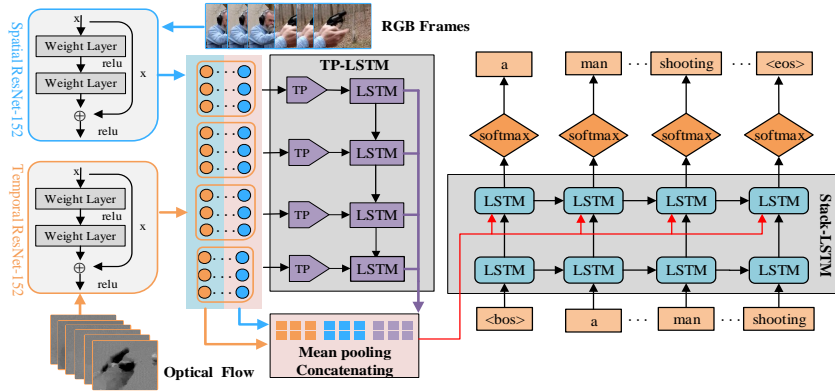
However, the above mentioned RNNs are suffering from the “long-term dependencies” problem [2]. LSTM [35] is designed for learning long-term dependencies. It solves the “long-term dependencies” problem by adding some gates that explicitly allow the RNN to learn when to forget previous hidden states with “forget gate” and when to update hidden states given new inputs. Previous studies showed that LSTM is capable of modeling data sequences, especially for encoding sentences and video features. Therefore, in this paper, we choose LSTM as our basic component for video captioning.

### 2.3 Video Captioning

As a bridge connecting computer vision and natural language processing, video captioning has attracted great attention in both areas. How to auto-generate descriptions of images or videos is an old topic in computer vision [10, 20, 22]. For example, Kojima *et al.* [20] firstly detected human postures, including head positions, head directions and hands positions, and then several predicts and objects are selected with domain knowledge. Finally, they filled these syntactic components into a case frame, and translated the case frames into sentences with some syntactic rules. In addition, same strategy is utilized to enhance other multimedia applications, such as [19, 16].

Later on, some researchers tried to describe videos/images with probabilistic graphical model [22, 10, 32]. For instance, Farhadi *et al.* [10] constructed three spaces: image space, sentence space and meaning space. In order to find the relationship between images and the corresponding sentences, they projected both image and sentence spaces into a common space: the meaning space. Specifically, the meaning space was represented by a triplet indicating as  $\langle \text{object, action, scene} \rangle$ . Mapping the image space to the meaning space was reduced to predicting the triplets from images, while mapping the sentence space into a meaning space was conducted by extracting triplets from sentences and then computing the similarity between two triplets. In addition, Rohrbach *et al.* [32] explored the relationship between visual contents and semantic representations with Conditional Random Field (CRF). However, all of these methods are highly dependent on the templates of sentences, which is insufficient to model the richness of natural language.

Recently, inspired by the great success of deep learning, many researchers [15, 28, 53, 48] applied deep neural networks to solve the video captioning problem. Specifically, Venugopalan *et al.* [48] employed a stacked LSTM for generating good descriptions effectively. The first LSTM encodes the visual features from pre-trained CNNs and the second LSTM generates words. Pan *et al.* [28] leveraged the semantics, both from entire sentence and video content, to learn a visual-semantic embedding model. Some works [29, 23] showed that semantic attributes make a significant contribution to video captioning. Pan *et al.* [29] adopted the Multiple Instance Learning (MIL) to learn the semantic at-



**Fig. 2** The framework of our model. TP-LSTM explores the invariance and variance in the video, while a Stack-LSTM is applied to generate words for describing the video.

tributes from videos, then utilized the generated attributes to improve the performance of their models. Compared with mean pooling, [38, 53, 56] were interested in tackling video captioning with attention mechanisms. Yao *et al.* [53] introduced a temporal soft attention mechanism into video captioning to automatically select the most relevant frames. Yu *et al.* [56] introduced a supervised spatial attention mechanism to guide the model to learn the relevant spatial information for video captioning. Different with above works, we are focusing on further extracting informative features for videos in terms of exploiting a long-range temporal structure.

### 3 The Proposed Approach

In this section, we introduce our approach for video captioning. Firstly, we define the terms and notations. Next, we describe our proposed network. Finally, we introduce the loss function of our model.

#### 3.1 Terms and Notations

Given a video  $\mathbf{V}$ , we extract its features  $\mathbf{V} = \{v_1, v_2, \dots, v_i, \dots, v_{N_v}\} \in \mathbb{R}^{D_v \times N_v}$ ,  $D_v$  denotes the dimension of visual features,  $N_v$  denotes the number of sampled frames from the video. A sentence  $\mathbf{S} = \{s_1, s_2, \dots, s_i, \dots, s_{N_s}\} \in \mathbb{R}^{D_s \times N_s}$  consisting of  $N_s$  words for describing the video, and  $s_i$  is a one-hot vector.  $D_s$  is the size of dictionary. And we denote  $\langle \text{BOS} \rangle$  as the start of a sentence. Our framework is shown in Fig. 2. This framework consists of six major components. The first component is a Spatial ResNet-152 network which takes RGB frames as inputs and extracts visual features from each video frame,

while the second component is the Temporal ResNet-152 network which takes optical flows as input and produces temporal features for each frame. Next, the third component is a concatenation that concatenates the outputs of Spatial ResNet-152 and Temporal ResNet-152 networks. Then, TP-LSTM takes a set of concatenations as inputs with a temporal pooling strategy. Finally, the second concatenation integrates visual features, temporal features and the outputs of TP-LSTM into a new video representation. The last component is a stacked LSTM, which takes the new video representation and words to produce a natural language sentence.

### 3.2 Temporal Pooling LSTM

How to extract effective visual features is an important problem for analyzing videos. Due to the rapid development of deep convolutional neural networks (CNNs), which have made a great success in image classification [17], object detection [31] and video action recognition [36], it is common to apply deep CNNs to extract visual features. In this work, we use the ResNet-152 pre-trained on the ImageNet to extract video frame visual features. In addition, a video contains not only spatial information but also temporal information. Therefore, we utilize another fine-tuned ResNet-152, which takes optical flow images as inputs, to extract video temporal features. After that, we concatenate above two features together. In order to model the invariance and variance of the input video, we propose a temporal pooling LSTM to dispose the fused new feature. More specifically, we divide the new features into  $N_e$  parts along the temporal dimension, thus each part has  $N_k = N_v/N_e$  features. Next, we average the features from same part. This process is expressed as follow:

$$e_i = \frac{\sum_{j=(i-1) \times N_k}^{i \times N_k} v_j}{N_k} \quad i \in \{1, 2, \dots, N_e\} \quad (1)$$

$\mathbf{E} = \{e_1, e_2, \dots, e_i, \dots, e_{N_e}\} \in \mathbb{R}^{D_v \times N_e}$  is generated after the temporal pooling.

In the next step, we aim to extract long-term dynamics across video by applying a Recurrent Neural Network (RNN) on  $\mathbf{E}$ . As mentioned above, we employ LSTM to model the long-term temporal dynamics of  $\mathbf{E}$ . The structure of LSTM is described below:

$$\begin{aligned} f_t &= \sigma(W_{xf}e_t + W_{hf}h_{t-1} + b_f) \\ i_t &= \sigma(W_{xi}e_t + W_{hi}h_{t-1} + b_i) \\ o_t &= \sigma(W_{xo}e_t + W_{ho}h_{t-1} + b_o) \\ g_t &= \phi(W_{xg}e_t + W_{hg}h_{t-1} + b_g) \\ c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\ h_t &= o_t \odot \phi(c_t) \end{aligned} \quad (2)$$

where  $\sigma(\cdot)$  denotes the sigmoid function,  $\phi(\cdot)$  denotes the hyperbolic tangent function, and  $\odot$  denotes the element-wise multiplication.  $c_t$  is a cell state

vector, and  $h_t$  is an hidden state vector.  $W_*$  is a set of parameters, and  $b_*$  is a set of bias values. For convenience, we define the function as:

$$h_t, c_t = LSTM(e_t, h_{t-1}, c_{t-1}; W, b) \quad t \in \{1, \dots, N_e\} \quad (3)$$

where  $e_t$  is the input at  $t$ -th time step, and  $h_0, c_0$  are initialized vectors. In our model, we use  $h_t$  as the output of the LSTM. After  $N_e$  time steps, we get  $\mathbf{H} = \{h_1, h_2 \dots, h_i, \dots, h_{N_e}\} \in \mathbb{R}^{D_h \times N_e}$ .  $D_h$  is the output dimension of the LSTM. Next, we average the outputs of LSTM and visual features, respectively. See below:

$$\begin{aligned} \bar{v} &= \frac{\sum_{i=1}^{N_v} v_i}{N_v} & v_i &\in \mathbf{V} \\ \bar{h} &= \frac{\sum_{i=1}^{N_e} h_i}{N_e} & h_i &\in \mathbf{H} \end{aligned} \quad (4)$$

Then, we concatenate them  $y = [\bar{v}, \bar{h}]$  and feed them into our Stack-LSTM.

### 3.3 Stacked LSTM

In order to reduce the dimension of the one-hot vector and explore the semantic information from the one-hot vector, we follow previous works [29, 55, 38] to embed the one-hot vector into a low-dimensional vector as follow:

$$\mathbf{M} = W_s \mathbf{S} \quad (5)$$

where  $W_s \in \mathbb{R}^{D_m \times D_s}$  is a parameter matrix. After embedding, we obtain an embedding matrix  $\mathbf{M} = \{m_1, m_2 \dots, m_i, \dots, m_{N_s}\} \in \mathbb{R}^{D_m \times N_s}$ .

Then we use LSTM layers to explore semantic information from both sentences and videos. Donahue *et al.* [8] suggested that two LSTM layers are better than one or four layers for image captioning. Compared with their two LSTM layers, our first LSTM layer is used to encode sentence information, while the second LSTM layer is applied to fuse both sentence and visual information for achieving semantic features. More specifically, at first we use a standard LSTM to explore the relationship between words:

$$q_t, u_t = LSTM(m_t, q_{t-1}, u_{t-1}; W_q, b_q) \quad t \in \{1, \dots, N_s\} \quad (6)$$

where  $q_0$  and  $u_0$  are initialized vectors.  $W_q$  and  $b_q$  are parameters. After  $N_s$  time steps, we get a series of vectors  $\mathbf{Q} = \{q_1, q_2 \dots, q_i, \dots, q_{N_s}\} \in \mathbb{R}^{D_q \times N_s}$ , which contain temporal information from a sentence. Next, we use a multi-modal LSTM (M-LSTM) which incorporates features from different information sources (i.e., video and words) into a set of higher-level representations. The M-LSTM integrates information from visual and word sources into latent



semantic features by adjusting their weights to improve the video captioning performance. The structure of the multi-modal LSTM is described as follows:

$$\begin{aligned}
f'_t &= \sigma(W'_{xf}q_t + W'_{hf}h'_{t-1} + W'_{yf}y + b'_f) \\
i'_t &= \sigma(W'_{xi}q_t + W'_{hi}h'_{t-1} + W'_{yi}y + b'_i) \\
o'_t &= \sigma(W'_{xo}q_t + W'_{ho}h'_{t-1} + W'_{yo}y + b'_o) \\
g'_t &= \phi(W'_{xg}q_t + W'_{hg}h'_{t-1} + W'_{yg}y + b'_g) \\
c'_t &= f'_t \odot c'_{t-1} + i'_t \odot g'_t \\
h'_t &= o'_t \odot \phi(c'_t)
\end{aligned} \tag{7}$$

where  $W'_*$  and  $b'_*$  are the parameters, which need to be learned.  $y$  is the concatenated feature, mentioned in Eq. 4.  $h'_0 \in \mathbb{R}^{D_{h'} \times 1}$  and  $c'_0 \in \mathbb{R}^{D_{h'} \times 1}$  are initialized vectors. Finally, we use a softmax layer to estimate the conditional probability distribution over  $s_{t+1}$ .

$$P(s_{t+1}|s_{<t}, \mathbf{V}) = \text{softmax}(W_f h'_t + b_f) \quad t \in \{1, \dots, N_s\} \tag{8}$$

where  $W_f \in \mathbb{R}^{D_s \times D_{h'}}$  and  $b_f \in \mathbb{R}^{D_s}$  are the parameters. If the input is represented as  $x \in \mathbb{R}^{D_s \times 1}$ , the softmax function can be expressed as:

$$\text{softmax}(x_i) = \frac{e^{x_i}}{\sum_{j=1}^{D_s} e^{x_j}} \quad i \in \{1, \dots, D_s\} \tag{9}$$

### 3.4 Loss Function

Previous works [49, 48, 55] defined their loss functions based on maximum likelihood estimation (MLE). In this work, we follow them to define our loss function by optimizing the log-likelihood:

$$\begin{aligned}
\mathcal{L} &= \log P(\mathbf{S}|\mathbf{V}) \\
&= \sum_{t=1}^{N_s} s_t^T \log P(s_t|s_{<t}, \mathbf{V})
\end{aligned} \tag{10}$$

By maximizing this loss function, we can estimate the parameters in the whole model. After extracting features with deep CNNs, we simultaneously train the rest of model (i.e. TP-LSTM, Mean Pooling Concatenating and Stack-LSTM in Fig. 2). More specifically, we use back-propagation through time (BPTT) algorithm to compute the gradients and conduct the optimization with adadelta [57].

## 4 Experiments

We evaluate our model on the task of video captioning. We firstly study the performance of different features, and then evaluate the influence of different hyper-parameters. Finally, we compare our model with the state-of-the-art methods.

## 4.1 Datasets

In our experiments, we use two public video captioning benchmarks that have been widely used in many other works.

**The Microsoft Video Description Corpus (MSVD).** This dataset is proposed by Chen *et al.* [4]. There are 1,970 short video clips collected from YouTube and about 80,000 descriptions collected by Amazon Mechanical Turkers (AMT) in this dataset, and an average length of each video clip about 9s. Each video clip has an average of forty descriptions. And this dataset is open-domain and covers a wide range of topics such as people, animals, sports, actions, music, scenarios, landscapes etc. In total, all the descriptions contain nearly 16,000 unique vocabularies. Following previous work [28,27,55], we split this dataset into a training, a validation and a testing dataset with 1200 (60%), 100 (5%) and 670 (35%) video clips, respectively.

**MSR Video to Text (MSR-VTT).** Xu *et al.* [52] collected this dataset by a commercial video search engine. It's a new large-scale and open-domain video captioning benchmark for supporting video understanding, especially for the task of automatically describing videos. There are 10K video clips and 200K descriptions in this dataset, collected by Amazon Mechanical Turkers workers (AMT) same as MSVD dataset, about 20 sentences for each short video. It covers about 20 categories and diverse visual content. The updated version contains many quality sentences, so we implement our experiments on the updated version. This dataset is divided into three subsets: 65% for training, 5% for validating and 30% for testing, corresponding to 6,513, 497 and 2,990 clips.

## 4.2 Evaluation Metrics

Following previous works [28,48,53], for evaluating the performance of our method, we utilize the following three evaluation metrics: BLUE [30], METEOR [7], and CIDEr [47].

## 4.3 Implementation Details

**Preprocessing.** For preprocessing the descriptions of MSVD dataset, firstly we convert sentences to lower cases, and then use the `wordpunct_tokenizer` in NLTK <sup>1</sup> library to tokenize sentences and remove punctuations. Finally, we obtain a dictionary of 15,903 in size on the training splits.

For preprocessing the descriptions of MSR-VTT dataset, we directly split descriptions with a blank space, because they have been tokenized. As a result, we can obtain a dictionary of 23,662 in size on the training splits. In this experiment, we only take words which appear more than two times as the dictionary. Finally, we get a dictionary of 13,626 in size.

<sup>1</sup> <http://www.nltk.org/index.html>

For the visual features, we use same method to extract features on both two datasets. For the spatial features, thanks to the ResNet-152 achieved the great results in image classification and video captioning [28, 54, 48], we use a per-trained ResNet-152 on ImageNet [33] to extract visual features. At first, we select equally-spaced 30 frames from each video, then feed them into the per-trained ResNet-152 to extract features from the *pool5* layer. Finally, we get a  $2048 \times 30$  feature matrix for each video. For the temporal features, inspired by [11], we first transform RGB images to optical flow images [3] stacking with 10 frames, then we use a fine-tuned ResNet-152 [11] on UCF101 to extract features from *pool5* layer. As a result, we obtain a  $2048 \times 30$  feature matrix. Next, we concatenate spatial and temporal features together and then feed them into our model. In our experiments,  $D_v = 4096$  and  $N_v = 30$ .

**Training Details.** In the training phase, sentences in corpus are varying lengths, thus we add a begin-of-sentence flag <BOS> to start each sentence and an end-of-sentence flag <EOS> to end each sentence. In the testing phase, we input <BOS> flag into our model to trigger the process of sentence generation. Beam search method, a heuristic search algorithm based on greedy algorithm, is utilized to find a sentence, which has the max partial probability. In addition, the width of the beam search is set as 5.

In addition, all the LSTM unit sizes are set as 512 ( $D_h = D_q = D_{h'} = 512$ ) and the word embedding size is set as 512 ( $D_m = 512$ ), empirically. In our experiments, we throw away the sentences whose length is more than 30, thus  $N_s < 30$ . The batch sizes are set as 64 on MSVD dataset and 256 on MSR-VTT dataset. We apply the back-propagation through time (BPTT) algorithm to compute the gradients of the parameters and conduct the optimization with adadelata [57]. In addition, we set the learning rate as  $10^{-4}$  to avoid the gradient explosion. We utilize dropout regularization with the rate of 0.5 in all layers and clip gradients element wise at 10. We stop training our model until 500 epoches are reached or until the evaluation metric does not improve the validation set at the patience of 20. Moreover, we utilize Theano [1] framework to conduct our experiments.

All experiments are conducted on the Ubuntu 14.04 with Intel(R) Xeon(R) CPU E5-2650 v3 @ 2.30GHz and GeForce GTX TITAN X (Pascal) GPU.

#### 4.4 Experiments on MSVD

For verifying the effectiveness of our framework, we design following experiments:

**Effectiveness of Different Features.** C3D features [46] are widely used for video captioning [55, 23, 27]. In this experiment, we evaluate the influence of spatial ResNet-152 feature (*res\_s*) and compare our temporal ResNet-152 (*res\_t*) features with the C3D features. The baseline is our model without the TP-LSTM part. The experimental results are shown in Tab. 1. From Tab. 1, we can see that simply applying spatial ResNet-152 is quite effective for video captioning with B@4 (51.5%), M (33.5%) and C (75.8%). Making use of both

**Table 1** Performances of our model with different features, where `res_s` stands for the the spatial ResNet-152 feature, `res_t` stands for the the temporal ResNet-152 feature, `c3d` stands for the C3D feature. B, M and C are short for BLUE, METEOR, and CIDEr. All values are reported as percentage (%).

| model                                               | B@1         | B@2         | B@3         | B@4         | M           | C           |
|-----------------------------------------------------|-------------|-------------|-------------|-------------|-------------|-------------|
| baseline( <code>res_s</code> )                      | 81.7        | 71.1        | 61.9        | 51.5        | 33.5        | 75.8        |
| baseline( <code>res_s</code> + <code>c3d</code> )   | <b>82.6</b> | <b>72.1</b> | <b>62.8</b> | 52.1        | 32.2        | 63.6        |
| baseline( <code>res_s</code> + <code>res_t</code> ) | 81.6        | 71.1        | 62.2        | <b>52.7</b> | <b>34.3</b> | <b>75.9</b> |

**Table 2** Performances of our model with different  $N_e$ . All models except baseline use spatial ResNet-152 feature and temporal ResNet-152 feature. B, M and C are short for BLUE, METEOR, and CIDEr. All values are reported as percentage (%).

| model                                                            | B@1         | B@2         | B@3         | B@4         | M           | C           |
|------------------------------------------------------------------|-------------|-------------|-------------|-------------|-------------|-------------|
| baseline ( <code>res_s</code> + <code>res_t</code> )             | 81.6        | 71.1        | 62.2        | 52.7        | 34.3        | 75.9        |
| TS-LSTM( $N_e = 1$ )( <code>res_s</code> + <code>res_t</code> )  | 83.0        | 72.1        | 62.6        | 52.8        | 33.7        | 77.2        |
| TS-LSTM( $N_e = 30$ )( <code>res_s</code> + <code>res_t</code> ) | 82.4        | 72.1        | 63.1        | 53.3        | 34.0        | 76.7        |
| TS-LSTM( $N_e = 3$ )( <code>res_s</code> + <code>res_t</code> )  | <b>83.8</b> | <b>73.8</b> | <b>64.5</b> | <b>54.5</b> | <b>34.5</b> | <b>79.3</b> |

`res_s` and `c3d`, B@1, B@2, B@3 and B@4 improves, but M and C drops. In terms of video captioning evaluation, METEOR and CIDE are more reliable than BLEU. Tab. 1 also shows that `res_s` and `res_t` performs best in terms of B@4, M and C. Therefore, we prove that our `res_t` performs better than `c3d` for video captioning. In the following experiments, all the models take both `res_s` and `res_t`.

**The Effect of Segmentation Numbers.** In order to explore the effectiveness of temporal pooling, we study the influence of the segmentation numbers, represented as  $N_e$ . In this experiment, we set  $N_e = 1$  (average whole features),  $N_e = 3$  and  $N_e = 30$  (without average operation) and the experiments are shown in Tab. 2, where baseline stands for our approach without TP-LSTM. From Tab. 2, we can see that when  $N_e = 3$ , our model achieves better results than  $N_e = 1$  and  $N_e = 30$ . When  $N_e = 1$ , the model ignores the temporal variance between long-range video shots. When  $N_e = 30$ , the model ignores the temporal invariance in a short video shot. This proves that reasonable number of segments can improve the performance of video captioning. Compared with the baseline, our model (TS-LSTM  $N_e = 3$ ) achieves better results with 2.2%, 2.7%, 2.3%, 1.8% 0.2% and 3.4% increases on BLUE-1, BLUE-2, BLUE-3, BLUE-4, METEOR and CIDEr, respectively. Therefore, in the following experiments, we set  $N_e = 3$ .

**Comparing with existing methods.** To verify the availability of our model, we compare our results with the following methods:

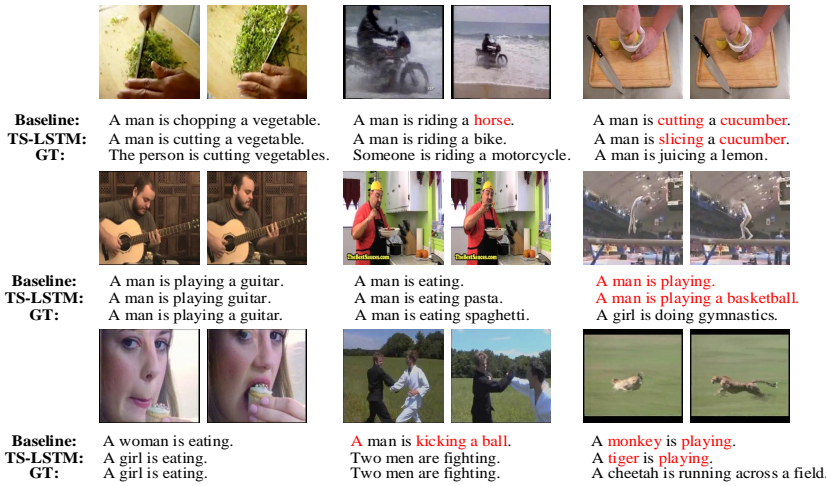
- MP-LSTM. Venugopalan *et al.* [49] used a mean-pooling layer to dispose all extracted frame-level features, then stacked two LSTM layers to explore semantic information.
- SA-LSTM. Yao *et al.* [53] introduced a temporal attention mechanism to automatically select the relevant frames, and combined with the spatial

**Table 3** BLEU@N (B@N), METEOR (M), and CIDEr(C) scores of our model and other state-of-the-art methods. This experiment is conducted on the MSVD dataset. All values are reported as percentage (%).

| Model        | B@1         | B@2         | B@3         | B@4         | M           | C           |
|--------------|-------------|-------------|-------------|-------------|-------------|-------------|
| MP-LSTM[49]  | -           | -           | -           | 33.3        | 29.1        | -           |
| SA-LSTM[53]  | 80.0        | 64.7        | 52.6        | 42.2        | 29.6        | 51.7        |
| LSTM-E[28]   | 78.8        | 66.0        | 55.4        | 45.3        | 31.0        | -           |
| HRNE-AT[27]  | 81.1        | 68.6        | 57.8        | 46.7        | 33.9        | -           |
| h-RNN[55]    | 81.5        | 70.4        | 60.4        | 49.9        | 32.6        | 65.8        |
| M3-LSTM[50]  | 82.5        | 72.4        | 62.8        | 52.8        | 33.3        | -           |
| MFA-LSTM[23] | 82.9        | 72.0        | 62.7        | 52.8        | 33.4        | 68.9        |
| LSTM-TSA[29] | 82.8        | 72.0        | 62.8        | 52.8        | 33.5        | 74.0        |
| hLSTMat[38]  | 82.9        | 72.2        | 63.0        | 53.0        | 33.6        | 73.8        |
| TS-LSTM      | <b>83.8</b> | <b>73.8</b> | <b>64.5</b> | <b>54.5</b> | <b>34.5</b> | <b>79.3</b> |

temporal 3-D convolutional neural network (3D-CNN) features, the model achieve the great results on the video captioning task.

- LSTM-E. Pan *et al.*, [28] assumed that a low-dimensional embedding exists for the representation of video and sentence, thus they mapped the video features and sentence features to the visual-semantic embedding and minimized the relevance loss to adequately explore the semantic information from videos.
- HRNE-AT. Pan *et al.* [27] proposed a Hierarchical Recurrent Neural Encoder (HRNE) structure, which stacks a short LSTM on a long LSTM for adequately exploring the temporal information of a video.
- h-RNN. Yu *et al.* [55] designed a sentence generator and a paragraph generator to generate paragraphs. The paragraph generator is stacked on the sentence generator and receives the state of the sentence generator, then initials the sentence generator.
- M3-LSTM. Wang *et al.* [50] designed a visual and semantic shared memory structure for achieving the long-term visual-semantic dependency to further guide global visual attention. In this way, the model can learn an effective mapping from visual space to language space.
- MFA-LSTM. Long *et al.* [23] selected the most frequent subject and verb across captions of each video, and took these as the semantic attributes and used a multi-modal attention mechanism to explore the semantic information from videos.
- LSTM-TSA. Pan *et al.* [29] introduced the Multiple Instance Learning (MIL). A weakly-supervised method was proposed to learn attribute detectors and great results were achieved.
- hLSTMat. Song *et al.* [38] proposed a adjusted temporal attention mechanism, which can automatically decide whether to depend on the visual features or the semantic information, to improve the attention mechanism on video captioning.



**Fig. 3** Some example sentences on the MSVD dataset. These sentences are generated by our model TS-LSTM and Baseline. GT denotes the ground truth. The imprecise words are labeled with red color.

#### 4.5 Comparison results on MSVD

In this experiment, we firstly compare our method with the existing methods on the MSVD dataset and the results are shown in Tab. 3. From Tab. 3, we can see that our model obtains the best performance. In particular, the BLEU-4 of our model reaches 54.5%, making a great improvement over h-RNN, MFA-LSTM, LSTM-TSA, hLSTMat by 4.6%, 1.7%, 1.7%, 1.5%, respectively. The METEOR of our model is 34.5%, which outperforms h-RNN, MFA-LSTM, LSTM-TSA, hLSTMat by 1.9%, 1.1%, 1.0%, 0.9%, respectively.

In Fig. 3, we show some example sentences generated by our TS-LSTM model and our baseline mentioned in Section 4.4. The first column shows that both TS-LSTM and baseline can generate correct sentences to describe each video. From the second column, we have the following observations: 1) TS-LSTM model can generate sentences with accurate words to describe objects within a video, such as “bike” in the top video. 2) Compared with baseline, TS-LSTM is able to provide more detailed information for describing video contents. For instance, in the middle video, TS-LSTM indicates that a man is “eating pasta” instead of just “eating”. 3) For the bottom video in the second column, it shows that TS-LSTM has ability to calculate the number of objects within a video. In addition, the third column introduces some wrong examples. For the bottom video in the third column, TS-LSTM and baseline generate two sentences: “a monkey is playing” and “a tiger is playing”, respectively. Both of them are incorrect due to the following reason that the MSVD dataset contains few videos about “cheetah”. Therefore, the trained models both encounter an over-fitting problem.

**Table 4** BLEU@N (B@N), METEOR (M), and CIDEr(C) scores of our model and other state-of-the-art methods. This experiment is conducted on MSR-VTT dataset. All values are reported as percentage (%).

| Model        | B@1         | B@2         | B@3         | B@4         | M           | C           |
|--------------|-------------|-------------|-------------|-------------|-------------|-------------|
| MP-LSTM[49]  | 81.7        | 65.1        | 48.5        | 35.8        | 25.3        | -           |
| SA-LSTM[53]  | <b>82.3</b> | <b>65.7</b> | 49.7        | 36.6        | 25.9        | -           |
| M3-LSTM[50]  | 73.6        | 59.3        | 48.3        | 38.1        | 26.6        | -           |
| hLSTMat[38]  | -           | -           | -           | 38.3        | 26.3        | -           |
| MFA-LSTM[23] | -           | -           | -           | 39.2        | 26.6        | <b>44.6</b> |
| TS-LSTM      | 77.6        | 64.0        | <b>51.3</b> | <b>39.9</b> | <b>27.1</b> | 43.8        |

#### 4.6 Comparison results on MSR-VTT

To further illustrate the performance of our model, we compare our model with the state-of-the-art methods on the MSR-VTT dataset, which has the largest number of video-sentence pairs. The experimental results are shown in Tab. 4.

From Tab. 4, we have the following observations. Firstly, our model achieves the best performance on BLEU-3 (51.3%), BLEU-4 (39.9%) and METEOR (27.1%). Compared with MP-LSTM, SA-LSTM, M3-LSTM, hLSTMat, MFA-LSTM, our model improves the BLEU-4 by 4.1%, 3.3%, 1.8%, 1.6%, 0.7%, respectively, and increases the METEOR by 1.8%, 1.2%, 0.5%, 0.8%, 0.5%, respectively.

## 5 Conclusion

In this paper, we present our temporal spatial LSTM network (TS-LSTM), a video level framework that aims to model long-term temporal dynamics and integrates dynamics with spatial and temporal features to improve video captioning. As demonstrated on two challenging datasets, this work outperforms the existing methods while keeping a reasonable computation cost. In this framework, the TP-LSTM is proposed to explore the long-range structure by taking segments of video visual and motion features as inputs, and produces informative long-term dynamics for video captioning. The experimental results show the effectiveness of our proposed approach.

## References

1. Bastien, F., Lamblin, P., Pascanu, R., Bergstra, J., Goodfellow, I.J., Bergeron, A., Bouchard, N., Warde-Farley, D., Bengio, Y.: Theano: new features and speed improvements. CoRR [abs/1211.5590](https://arxiv.org/abs/1211.5590) (2012)
2. Bengio, Y., Simard, P.Y., Frasconi, P.: Learning long-term dependencies with gradient descent is difficult. IEEE Trans. Neural Networks **5**(2), 157–166 (1994)
3. Brox, T., Bruhn, A., Papenberger, N., Weickert, J.: High accuracy optical flow estimation based on a theory for warping. In: ECCV, pp. 25–36 (2004)
4. Chen, D., Dolan, W.B.: Collecting highly parallel data for paraphrase evaluation. In: ACL HLT, pp. 190–200 (2011)

5. Chen, X., Fang, H., Lin, T.Y., Vedantam, R., Gupta, S., Dollar, P., Zitnick, C.L.: Microsoft coco captions: Data collection and evaluation server. *Computer Science* (2015)
6. Chung, J., Gülçehre, Ç., Cho, K., Bengio, Y.: Empirical evaluation of gated recurrent neural networks on sequence modeling. *CoRR* **abs/1412.3555** (2014)
7. Denkowski, M.J., Lavie, A.: Meteor universal: Language specific translation evaluation for any target language. In: *The Workshop on Statistical Machine Translation*, pp. 376–380 (2014)
8. Donahue, J., Hendricks, L.A., Rohrbach, M., Venugopalan, S., Guadarrama, S., Saenko, K., Darrell, T.: Long-term recurrent convolutional networks for visual recognition and description. *IEEE Trans. Pattern Anal. Mach. Intell.* **39**(4), 677–691 (2017)
9. Elman, J.L.: Finding structure in time. *Cognitive Science* **14**(2), 179–211 (1990)
10. Farhadi, A., Hejrati, S.M.M., Sadeghi, M.A., Young, P., Rashtchian, C., Hockenmaier, J., Forsyth, D.A.: Every picture tells a story: Generating sentences from images. In: *ECCV*, pp. 15–29 (2010)
11. Feichtenhofer, C., Pinz, A., Zisserman, A.: Convolutional two-stream network fusion for video action recognition. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27–30, 2016*, pp. 1933–1941 (2016)
12. Gao, L., Guo, Z., Zhang, H., Xu, X., Shen, H.T.: Video captioning with attention-based lstm and semantic consistency. *IEEE Transactions on Multimedia* **19**(9), 2045–2055 (2017). DOI 10.1109/TMM.2017.2729019
13. Gao, L., Song, J., Nie, F., Yan, Y., Sebe, N., Shen, H.T.: Optimal graph learning with partial tags and multiple features for image and video annotation. In: *CVPR*, pp. 4371–4379 (2015)
14. Gao, L., Song, J., Nie, F., Zou, F., Sebe, N., Shen, H.T.: Graph-without-cut: An ideal graph learning for image segmentation. In: *AAAI*, pp. 1188–1194 (2016)
15. Guo, Z., Gao, L., Song, J., Xu, X., Shao, J., Shen, H.T.: Attention-based LSTM with semantic consistency for videos captioning. In: *ACM MM*, pp. 357–361 (2016)
16. Hanckmann, P., Schutte, K., Burghouts, G.J.: Automated textual descriptions for a wide range of video events with 48 human actions. In: *ECCV*, pp. 372–380 (2012)
17. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *CVPR*, pp. 770–778 (2016)
18. Jordan, M.I.: Serial order: A parallel distributed processing approach. **121**, 64 (1986)
19. Khan, M.U.G., Zhang, L., Gotoh, Y.: Human focused video description. In: *ICCV*, pp. 1480–1487 (2011)
20. Kojima, A., Tamura, T., Fukunaga, K.: Natural language description of human activities from video images based on concept hierarchy of actions. *International Journal of Computer Vision* pp. 171–184 (2002)
21. Lecun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proceedings of the IEEE* **86**(11), 2278–2324 (1998)
22. Lee, M.W., Hakeem, A., Haering, N., Zhu, S.: SAVE: A framework for semantic annotation of visual events. In: *CVPR*, pp. 1–8 (2008)
23. Long, X., Gan, C., de Melo, G.: Video captioning with multi-faceted attention. *CoRR* **abs/1612.00234** (2016)
24. Ma, C., Chen, M., Kira, Z., AlRegib, G.: TS-LSTM and temporal-inception: Exploiting spatiotemporal dynamics for activity recognition. *CoRR* **abs/1703.10667** (2017)
25. Mikolov, T., Karafiát, M., Burget, L., Cernocký, J., Khudanpur, S.: Recurrent neural network based language model. In: *INTERSPEECH*, pp. 1045–1048 (2010)
26. Mnih, V., Heess, N., Graves, A., Kavukcuoglu, K.: Recurrent models of visual attention. In: *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8–13 2014, Montreal, Quebec, Canada*, pp. 2204–2212 (2014)
27. Pan, P., Xu, Z., Yang, Y., Wu, F., Zhuang, Y.: Hierarchical recurrent neural encoder for video representation with application to captioning. In: *CVPR*, pp. 1029–1038 (2016)
28. Pan, Y., Mei, T., Yao, T., Li, H., Rui, Y.: Jointly modeling embedding and translation to bridge video and language. In: *CVPR*, pp. 4594–4602 (2016)
29. Pan, Y., Yao, T., Li, H., Mei, T.: Video captioning with transferred semantic attributes. In: *CVPR* (2017)
30. Papineni, K., Roukos, S., Ward, T., Zhu, W.: Bleu: a method for automatic evaluation of machine translation. In: *ACL*, pp. 311–318 (2002)



31. Ren, S., He, K., Girshick, R., Sun, J.: Faster R-CNN: Towards real-time object detection with region proposal networks. In: NIPS (2015)
32. Rohrbach, M., Qiu, W., Titov, I., Thater, S., Pinkal, M., Schiele, B.: Translating video content to natural language descriptions. In: ICCV, pp. 433–440 (2013)
33. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M.S., Berg, A.C., Li, F.: Imagenet large scale visual recognition challenge. *International Journal of Computer Vision* **115**(3), 211–252 (2015)
34. Scherer, D., Müller, A.C., Behnke, S.: Evaluation of pooling operations in convolutional architectures for object recognition. In: Artificial Neural Networks - ICANN 2010 - 20th International Conference, Thessaloniki, Greece, September 15-18, 2010, Proceedings, Part III, pp. 92–101 (2010)
35. S.Hochreiter, J.Schmidhuber: Long short-term memory. *Neural Computation* **9**(8), 1735–1780 (1997)
36. Simonyan, K., Zisserman, A.: Two-stream convolutional networks for action recognition in videos. In: Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada, pp. 568–576 (2014)
37. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: ICLR (2014)
38. Song, J., Gao, L., Guo, Z., Liu, W., Zhang, D., Shen, H.T.: Hierarchical LSTM with adjusted temporal attention for video captioning. In: Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017, pp. 2737–2743 (2017)
39. Song, J., Gao, L., Liu, L., Zhu, X., Sebe, N.: Quantization-based hashing: A general framework for scalable image and video retrieval. *Pattern Recognition* (2017)
40. Song, J., Gao, L., Nie, F., Shen, H.T., Yan, Y., Sebe, N.: Optimized graph learning using partial tags and multiple features for image and video annotation. *IEEE Trans. Image Processing* **25**(11), 4999–5011 (2016)
41. Song, J., Gao, L., Puscas, M.M., Nie, F., Shen, F., Sebe, N.: Joint graph learning and video segmentation via multiple cues and topology calibration. In: Proceedings of the 2016 ACM Conference on Multimedia Conference, MM 2016, Amsterdam, The Netherlands, October 15-19, 2016, pp. 831–840 (2016)
42. Song, J., He, T., Fan, H., Gao, L.: Deep discrete hashing with self-supervised pairwise labels. *CoRR* [abs/1707.02112](https://arxiv.org/abs/1707.02112) (2017)
43. Sutskever, I., Vinyals, O., Le, Q.V.: Sequence to sequence learning with neural networks. In: NIPS, pp. 3104–3112 (2014)
44. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S.E., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: CVPR, pp. 1–9 (2015)
45. Thonnat, M., Rota, N.: Image understanding for visual surveillance applications (2000)
46. Tran, D., Bourdev, L.D., Fergus, R., Torresani, L., Paluri, M.: C3D: generic features for video analysis. *ICCV*
47. Vedantam, R., Zitnick, C.L., Parikh, D.: Cider: Consensus-based image description evaluation. In: CVPR, pp. 4566–4575 (2015)
48. Venugopalan, S., Rohrbach, M., Donahue, J., Mooney, R.J., Darrell, T., Saenko, K.: Sequence to sequence - video to text. In: ICCV, pp. 4534–4542 (2015)
49. Venugopalan, S., Xu, H., Donahue, J., Rohrbach, M., Mooney, R.J., Saenko, K.: Translating videos to natural language using deep recurrent neural networks. In: NAACL HLT, pp. 1494–1504 (2015)
50. Wang, J., Wang, W., Huang, Y., Wang, L., Tan, T.: Multimodal memory modelling for video captioning. *CoRR* [abs/1611.05592](https://arxiv.org/abs/1611.05592) (2016)
51. Wang, L., Xiong, Y., Wang, Z., Qiao, Y., Lin, D., Tang, X., Gool, L.V.: Temporal segment networks: Towards good practices for deep action recognition. In: Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part VIII, pp. 20–36 (2016)
52. Xu, J., Mei, T., Yao, T., Rui, Y.: MSR-VTT: A large video description dataset for bridging video and language. In: CVPR, pp. 5288–5296 (2016)
53. Yao, L., Torabi, A., Cho, K., Ballas, N., Pal, C.J., Larochelle, H., Courville, A.C.: Describing videos by exploiting temporal structure. In: ICCV, pp. 4507–4515 (2015)

54. Yao, T., Pan, Y., Li, Y., Qiu, Z., Mei, T.: Boosting image captioning with attributes. CoRR **abs/1611.01646** (2016)
55. Yu, H., Wang, J., Huang, Z., Yang, Y., Xu, W.: Video paragraph captioning using hierarchical recurrent neural networks. In: CVPR, pp. 4584–4593 (2016)
56. Yu, Y., Choi, J., Kim, Y., Yoo, K., Lee, S.H., Kim, G.: Supervising neural attention models for video captioning by human gaze data. In: CVPR (2017)
57. Zeiler, M.D.: ADADELTA: an adaptive learning rate method. CoRR **abs/1212.5701** (2012)
58. Zhu, X., Huang, Z., Shen, H.T., Zhao, X.: Linear cross-modal hashing for efficient multimedia search. In: ACM MM, pp. 143–152 (2013)
59. Zhu, X., Li, X., Zhang, S.: Block-row sparse multiview multilabel learning for image classification. IEEE transactions on cybernetics **46**(2), 450–461 (2016)
60. Zhu, X., Zhang, L., Huang, Z.: A sparse embedding and least variance encoding approach to hashing. IEEE transactions on image processing **23**(9), 3737–3750 (2014)