# Hierarchical LSTMs with Adaptive Attention for Visual Captioning

Jingkuan Song, *Member, IEEE,* Xiangpeng Li, Lianli Gao, and Heng Tao Shen

**Abstract**—Recent progress has been made in using attention based encoder-decoder framework for image and video captioning. Most existing decoders apply the attention mechanism to every generated word including both visual words (e.g., "gun" and "shooting") and non-visual words (e.g. "the", "a"). However, these non-visual words can be easily predicted using natural language model without considering visual signals or attention. Imposing attention mechanism on non-visual words could mislead and decrease the overall performance of visual captioning. Furthermore, the hierarchy of LSTMs enables more complex representation of visual data, capturing information at different scales. To address these issues, we propose a hierarchical LSTM with adaptive attention (hLSTMat) approach for image and video captioning. Specifically, the proposed framework utilizes the spatial or temporal attention for selecting specific regions or frames to predict the related words, while the adaptive attention is for deciding whether to depend on the visual information or the language context information. Also, a hierarchical LSTMs is designed to simultaneously consider both low-level visual information and high-level language context information to support the caption generation. We initially design our hLSTMat for video captioning task. Then, we further refine it and apply it to image captioning task. To demonstrate the effectiveness of our proposed framework, we test our method on both video and image captioning tasks. Experimental results show that our approach achieves the state-of-the-art performance for most of the evaluation metrics on both tasks. The effect of important components is also well exploited in the ablation study.

**Index Terms**—Video Captioning, Image Captioning, Adaptive Temporal Attention, Hierarchical Structure.

✦

## 1 INTRODUCTION

P REVIOUSLY, visual content understanding [1], [2] and natural language processing (NLP) are not correlative with each other. Integrating visual content with natural language learning to generate descriptions for images, especially for videos, has been regarded as a challenging task. Video captioning is a critical step towards machine intelligence and many applications in daily scenarios, such as video retrieval [3], [4], video understanding, blind navigation and automatic video subtitling.

Thanks to the rapid development of deep Convolutional Neural Network (CNN), recent works have made significant progress for image captioning [5], [6], [7], [8], [9], [10], [11]. However, compared with image captioning, video captioning is more difficult due to the diverse sets of objects, scenes, actions, attributes and salient contents. Despite the difficulty there have been a few attempts for video description generation [12], [13], [14], [15], [16], which are mainly inspired by recent advances in translating with Long Short-Term Memory (LSTM). The LSTM is proposed to overcome the vanishing gradients problem by enabling the network to learn when to forget previous hidden states and when to update hidden states by integrating memory units. LSTM has been successfully adopted to several tasks, e.g., speech recognition, language translation and image captioning [12], [17]. Thus, we follow this elegant recipe and choose to extend LSTM to generate the video sentence with semantic content.

Early attempts were proposed [12], [13], [14], [15] to directly connect a visual convolution model to a deep LSTM networks. For example, Venugopalan *et al.* [12] translate videos to sentences by

directly concatenating a deep neural network with a recurrent neural network. More recently, attention mechanism [18] is a standard part of the deep learning toolkit, contributing to impressive results in neural machine translation [19], visual captioning [6], [14] and question answering [20]. Visual attention models for video captioning make use of video frames at every time step, without explicitly considering the semantic attributes of the predicted words. For example, some words (i.e., "man", "shooting" and "gun") belong to visual words which have corresponding canonical visual signals, while other words (i.e., "the", "a" and "is") are non-visual words, which require no visual information but language context information [7]. In other words, current visual attention models make use of visual information for generating each work, which is unnecessary or even misleading. Ideally, video description not only requires modeling and integrating their sequence dynamic temporal attention information into a natural language but also needs to take into account the relationship between sentence semantics and visual content [21], which to our knowledge has not been simultaneously considered.

Although the encoder-decoder framework with attention mechanisms has been widely applied to solve other sequence generation tasks, e.g., machine translation [22] and image annotation [23]. In standard sequence generation tasks, an input is encoded to a vector embedding, and then decoded to an output string of words using RNNs, e.g., LSTM. However, such framework is in essential an one-pass forward process. When a model predicts the next word, it can only leverage the generated words but not the future unknown words. To humans, deliberation action is a common behavior in their daily, e.g., reading, writing or understanding an image. During the process, global information of both the past and future are leveraged. Xie *et. al.* [24] designed a deliberation network which included two levels of decoders, and it is proved to be effective for neural machine translation. The first one generates

---

- *J. Song, X. Li, L. Gao and H. Shen are with the Future Media Center and School of Computer Science and Engineering, The University of Electronic Science and Technology of China, Chengdu, China, 611731. E-mail: lianli.gao@uestc.edu.cn*

a coarse sentence and corresponding hidden states. The second decoder refines the sentence with deliberation. In the network, the second decoder could leverage the global information of both the past and future parts. Also, the success of deep neural networks is commonly attributed to the hierarchy that is introduced due to the several layers. Each layer processes some part of the task we wish to solve, and passes it on to the next. However, most of the current visual captioning tasks utilize LSTM with a single layer.

To tackle these issues, in this paper we propose a unified encoder-decoder framework (see Fig. 1), named hLSTMat, a Hierarchical LSTMs with adaptive temporal attention model for visual captioning. Specifically, first, in order to extract more meaningful visual features, we adopt a deep neural network to extract region-level or frame-level 2D CNN features for each video or image. Next, we integrate a hierarchical LSTMs consisting of two layers of LSTMs and adaptive attention to decode visual information and language context information to support the generation of sentences for video and image description. Moreover, the proposed novel adaptive attention mechanism automatically decides whether to rely on visual information or not.

It is worthwhile to highlight the main contributions of this proposed approach: 1) We introduce a novel hLSTMat framework which automatically decides when and where to use visual information, and when and how to adopt the language model to generate the next word for visual captioning. Specifically, the spatial and temporal attention is used to decide where to look at visual information and the adaptive attention decides when to rely on language context information. 2) A hierarchical LSTMs is designed to enrich the representation ability of the LSTM. Specifically, when the two LSTMs are connected at each time step, the hierarchical LSTMs can obtain both low-level visual information and high-level language context information. On the other hand, when the two LSTMs are connected sequentially, the second LSTM can be considered as a deliberation process to refine the first LSTM because it is based on the rough global features captured by the first LSTM. 3) Extensive experiments are conducted on mainstream benchmark datasets for both video and image captioning tasks. Experimental results show that our approach achieves the state-of-the-art performance for most of the evaluation metrics on both tasks. The effect of important components is also well exploited in the ablation study.

The rest of our paper is organized as follows: Firstly, works about visual captioning are introduced in Sec. 2. Then we describe the details of our proposed model hLSTMat in Sec. 3. We first instantiate our hLSTMat for the task of video captioning in Sec. 4, and show the experimental results for video captioning in Sec. 5. Then we further instantiate our hLSTMat for the task of image captioning in Sec. 6, and show the experimental results for image captioning in Sec. 7. Finally, we conclude our method in Sec. 8.

## 2 RELATED WORK

Image and video captioning, which automatically translates images and videos into natural language sentences to describe their content, is a very important task. Recent years, we have witnessed the significant growth of captioning and previous work is mainly focusing on two aspects: feature extraction and model innovation.

### 2.1 Feature Extraction

Local feature or local image descriptor is at the core of many computer vision tasks, and classical methods such as HOF, HOG

and SIFT have been extensively used in various computer vision applications such as image classification, object detection, image retrieval and segmentation. Alone with the booming of handcrafted descriptors in the past decade and the large collected datasets like ImageNet, more and more learning based descriptors appears, such as AlexNet [25], GoogLenet [26], VGG [27] and ResNet [28]. Different from handcrafted descriptors which are mostly driven by intuitions or domain expert knowledge, learning based methods are driven by the large scale annotated datasets and the rapid development of GPU [29]. To date, deep learning has revolutionized almost the whole computer vision fields and even other areas such as natural language processing and medical image research. Inspired by the success of recurrent neural network (RNN) in neural machine translation and deep convolutional neural network (CNN) in various visual analysis tasks, image captioning is proposed since associating image captioning with machine translation makes sense. They can be placed in the same framework, called Encoder-Decoder framework. As a result, numerous discriminative CNN features have been proposed an adopted in image captioning task to extract appearance features, such as AlexNet [25], GoogLenet [26], VGG [27] and ResNet [28]. Compared with learning based image descriptor, video descriptor has not yet seen the substantial gains in performance that have been achieved in other areas by CNNs, e.g., image classification and object detection. C3D [30] model is proposed to explore temporal information and good at extracting motion information. More video dataset is required to improve the performance of learning based video descriptor. Therefore, in this paper, we adopt CNNs to extract image/frame appearance feature and C3D model to capture video motion information.

### 2.2 Model Innovation

Except improving the image/video captioning with better features, researchers are also focusing on captioning model innovation, which plays an irreplaceable role in improving the performance.

**Stage 1: Basic Encoder-Decoder.** At the earlier stage of visual captioning, several models such as [5], [13], [31] have been proposed by directly bring together previous advances in natural language processing and computer vision. More specifically, semantic representation of an image is captured by a CNN network and then decoded into a caption using various architectures, such as recurrent neural networks. For example, Venugopalan *et al.* [13] proposed a S2VT approach, which incorporates a stacked LSTM by firstly reading the visual sequence, comprised of RGB and/or optical flow CNN outputs, and then generating a sequence of words. Oriol *et al.* [5] presented a generative model based on a deep recurrent neural network. This model consists of a vision CNN followed by a language generating RNN, and is trained to maximize the likelihood of the target description sentence given the training image. In [31], the proposed framework consists of three parts: a compositional semantics language model, a deep video model and a joint embedding model. In the joint embedding model, the distance of the outputs of the deep video model and compositional language model is minimized in the joint space.

**Stage 2: Visual Attention.** Later on, researchers found that different regions in images and frames in videos have different weights, and thus various attention mechanisms are introduced to guide captioning models by telling where to look at for sentence generation, such as [6], [14], [32], [33]. Yao *et al.* [14] proposed to incorporate both the local dynamic of videos as well as their global
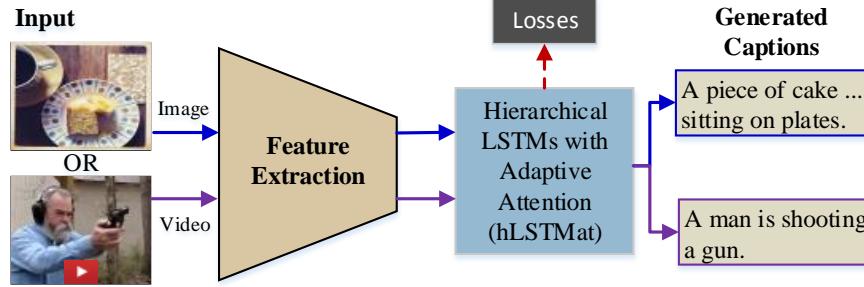
Fig. 1. The framework of our proposed hLSTMat for visual captioning. Given an input image or video, an encoder is firstly applied to extract the features. Then hierarchical LSTM with adaptive attention component plays the role of an decoder, by using the hierarchical LSTM to extract different level of information, and an adaptive attention to decide whether to depend on the visual information or the language context information. The losses are defined on the generated captions and the groundtruth to guide the learning of network parameters. Note that when this framework is applied to different visual captioning tasks, i.e., image and video captioning, there are differences in terms of feature extractor, network structure and losses.

temporal structure for describing videos. For simplicity, they were focusing on highlighting only the region having the maximum attention. A Hierarchical Recurrent Neural Encoder (HRNE) [32], was introduced to generate video representation with emphasis on temporal modeling by applying a LSTM along with attention mechanism to each temporal time step. In [34], it combines multiple forms of attention for video captioning. Temporal, motion and semantic features are weighted via an attention mechanism.

**Stage 3: Semantic Attention.** Semantic attention has been proposed in previous work [37], [38], [39] by adopting attributes or concepts generated by other pre-trained models to enhance captioning performance. Basically, semantic attention is able to attend to a semantically important concepts or attributes or region of interest in an image, and able to weight the relative strength of attention paid on multiple concepts [37]. In addition, Yao *et al.* presented Long Short-Term Memory with Attributes (LSTM-A) to integrates attributes into the successful CNNs and RNNs for image captioning. Variant architectures were constructed to feed image features and attributes into RNNs in different ways to explore the mutual but also fuzzy relationship between them. Moreover, Pan *et al.* [39] proposed Long Short-Term Memory with Transferred Semantic Attributes (LSTM-TSA), which takes advantages of incorporating transferred semantic attributes learnt from images and videos into sequence learning for video captioning.

**stage 4: Deep Reinforcement Learning based Approaches.** Recently, several researchers have started to utilize reinforcement learning to optimize image captioning [40], [41], [42], [43]. To collaboratively generate caption, Ren *et al.* [43] incorporated a "policy network" for generating sentence and a "value network" for evaluating predicting sentence globally, and they improved the policy and value networks with deep reinforcement learning. Furthermore, Liu *et al.* [42] proposed to improve image captioning via policy gradient optimization of a linear combination of SPICE and CIDEr, while Rennis *et al.* [41] utilized the output of its own test-time inference algorithm to normalize the rewards it experiences. All the previous methods show that reinforcement learning has the potential to boost image captioning.

# 3 HIERARCHICAL LSTM WITH ADAPTIVE ATTENTION FOR VISUAL CAPTIONING

In this section, we first briefly describe how to directly use the basic Long Short-Term Memory (LSTM) as the decoder for visual captioning task. Then we introduce our novel encoder-decoder framework, named Hierarchical LSTM with Adaptive Attention (hLSTMat) (See Fig. 1) for both video and image captioning.

## 3.1 A Basic LSTM for Visual Captioning

To date, modeling sequence data with Recurrent Neural Networks (RNNs) has shown great success in the process of machine translation, speech recognition, image and video captioning [9], [10], [12], [13] etc. Long Short-Term Memory (LSTM) is a variant of RNN to avoid the vanishing gradient problem [44].

**LSTM Unit.** A basic LSTM unit consists of three gates (input $\mathbf{i}_t$, forget $\mathbf{f}_t$ and output $\mathbf{o}_t$), a single memory cell $\mathbf{m}_t$. Specifically, $\mathbf{i}_t$ allows incoming signals to alter the state of the memory cell or block it. $\mathbf{f}_t$ controls what to be remembered or be forgotten by the cell, and somehow can avoid the gradient from vanishing or exploding when back propagating through time. Finally, $\mathbf{o}_t$ allows the state of the memory cell to have an effect on other neurons or prevent it. Basically, the memory cell and gates in a LSTM block are defined as follows:

$$\begin{aligned} \mathbf{i}_t &= \sigma(\mathbf{W}_i\mathbf{y}_t + \mathbf{U}_i h_{t-1} + \mathbf{b}_i) \\ \mathbf{f}_t &= \sigma(\mathbf{W}_f\mathbf{y}_t + \mathbf{U}_f h_{t-1} + \mathbf{b}_f) \\ \mathbf{o}_t &= \sigma(\mathbf{W}_o\mathbf{y}_t + \mathbf{U}_o h_{t-1} + \mathbf{b}_o) \\ \mathbf{g}_t &= \phi(\mathbf{W}_g\mathbf{y}_t + \mathbf{U}_g h_{t-1} + \mathbf{b}_g) \\ \mathbf{m}_t &= \mathbf{f}_t \odot \mathbf{m}_{t-1} + \mathbf{i}_t \odot \mathbf{g}_t \\ \mathbf{h}_t &= \mathbf{o}_t \odot \phi(\mathbf{m}_t) \end{aligned} \tag{1}$$

where the weight matrices $\mathbf{W}$, $\mathbf{U}$, and $\mathbf{b}$ are parameters to be learned. $\mathbf{y}_t$ represents the input vector for the LSTM unit at each time $t$. $\sigma$ represents the logistic sigmoid non-linear activation function mapping real numbers to $(0, 1)$, and it can be thought as knobs that LSTM learns to selectively forget its memory or accept the current input. $\phi$ denotes the hyperbolic tangent function tanh. $\odot$ is the element-wise product with the gate value. For convenience, we denote $\mathbf{h}_t, \mathbf{m}_t = \text{LSTM}(\mathbf{y}_t, \mathbf{h}_{t-1}, \mathbf{m}_{t-1})$ as the computation function for updating the LSTM internal state.

**Visual Captioning.** Given an input video or image $\mathbf{x}$, an encoder network $\phi_E$ encodes it into a representation space:

$$\mathbf{V} = \phi_E(\mathbf{x}). \tag{2}$$

where $\phi_E$ usually denotes a CNN neural network. Here, LSTM is chosen as a decoder network $\phi_D$ to model $\mathbf{V}$ to generate a description $\mathbf{z} = \{z_1, \cdots, z_T\}$ for $\mathbf{x}$, where $T$ is the description length. In addition, the LSTM unit updates its internal state $\mathbf{h}_t$

and the $t$-th word $z_t$ based on its previous internal state $\mathbf{h}_{t-1}$, the previous output $y_t$ and the representation $\mathbf{V}$:

$$(\mathbf{h}_t, z_t) = \phi_D(\mathbf{h}_{t-1}, y_t, \mathbf{V})$$

In test stage, the previous output $\mathbf{z}_{t-1}$ is the current input vector $\mathbf{y}_t$. And the initial $\mathbf{y}_0$ represents the begin of the sentence $(BOS)$.

In addition, the LSTM updates its internal state recursively until the end-of-sentence tag is generated. For simplicity, we named this simple method as basic-LSTM.

### 3.2  Hierarchical LSTMs with Adaptive Temporal Attention for Visual Captioning

In this subsection, we introduce our hLSTMat framework, which consists of three major components: 1) an encoder; 2) an attention based hierarchical LSTM decoder; and 3) the losses.

#### 3.2.1  CNN Encoders

The goal of an encoder is to compute features that are compact and representative and can capture the most related visual information for the decoder. Thanks to the rapid development of CNNs, a great success has been made in large scale image recognition task [28], object detection [45] and visual captioning [12]. High-level features can been extracted from upper or intermediate layers of a deep CNN network. Therefore, a set of well-tested CNN networks, such as the ResNet-152 model [28] which has achieved the best performance in ImageNet Large Scale Visual Recognition Challenge, can be used as a candidate encoder for our framework. Depending on the type of input, different types of features can be extracted. For example, we can extract frame-level features for videos and we can extract region-level features for images. The details of feature extraction are task-specific, and will be given in each specific task in Section 4 and Section 6.

#### 3.2.2  hLSTMat Decoder

In the basic model, a LSTM is used as the decoder to predict the captions. The vanilla LSTM model is comprised of a single hidden LSTM layer followed by a standard feedforward output layer. However, the success of deep neural networks is commonly attributed to the hierarchy that is introduced due to the several layers [46]. Therefore, we proposed a hierarchical LSTM based decoder. It can be seen as a processing pipeline, in which each layer solves a part of the task before passing it on to the next, until finally the last layer provides the output.

Also, we define an adaptive attention layer based on the hierarchical LSTM for deciding whether to depend on the visual information or the language context information.

#### 3.2.3  Losses

After the top MLP layer predicts the probability distribution of each word in the vocabulary, maximum likelihood estimation (MLE) loss is usually utilized to guide the learning of parameters:

$$\ell_{MLE} = -\sum_{t=1}^{T} log P(z_t | z_{<t}, \mathbf{V}, \Theta) \qquad (3)$$

where $z_t$ is the $t$-th word in the ground truth caption, $z_{<t}$ are the first $t$-1 words in the ground truth caption, $T$ denotes the total number of words in sentence, $\mathbf{V}$ denotes the features of the corresponding input, and $\Theta$ are model parameters. Therefore, Eq. 3 is regarded as the MLE loss function to optimize the model.

There are also other types of losses for visual captioning. We consider our method introduced above as "agent" to interact with external environment (i.e., words, visual features), and $\Theta$ as the policy to conduct an action to predict a word. After the whole caption is generated, the agent observers a reward. Since CIDEr is proposed to evaluate the quality of visual captioning model. We can design our reward functions using CIDEr, and the loss based on deep reinforcement learning can be defined.

## 4  HLSTMAT FOR VIDEO CAPTIONING

In this subsection, we instantiate our hLSTMat and apply it to the task of video captioning. As shown in Fig. 2, there are three major components: 1) a CNN Encoder, 2) an attention based hierarchical LSTM decoder and 3) the losses. We give the details for each component below.

### 4.1  CNN Encoders

Inspired by [14], we preprocess each video clip by selecting equally-spaced 28 frames out of the first 360 frames and then feeding them into a CNN network proposed in [28]. Thus, for each selected frame we obtain a 2,048-D feature vector, which are extracted from the $pool5$ layer. For simplicity, given a video $\mathbf{x}$, we extract $L$=28 frame-level features $\mathbf{V}$, represented as $\mathbf{V} = \{\mathbf{v}_1, \cdots, \mathbf{v}_L\}$.

Motion features catches the long range dependencies of video, which is quite essential for video information. The C3D motion net reads 16 adjacent frames as a segment and catches a 4096-D motion feature vector as output. For each video, we attain 10 segment features to represent the whole video.

### 4.2  Attention based Hierarchical Decoder

Our decoder (see Fig. 3) integrates two LSTMs. The bottom LSTM layer is used to efficiently decode visual features, and the top LSTM is focusing on mining deep language context information for video captioning. We also incorporate two attention mechanisms into our framework. A temporal attention is to guide which frame to look, while the adaptive temporal attention is proposed to decide when to use visual information and when to use sentence context information. The top MLP layer is to predict the probability distribution of each word in the vocabulary.

Unlike vanilla LSTM decoder, which performs mean pooling over 2D features across each video to form a fixed-dimension representation, attention based LSTM decoder is focusing on a subset of consecutive frames to form a fixed-dimensional representation at each time $t$.

- Bottom LSTM Layer. For the bottom LSTM layer, the updated internal hidden state depends on the current word $y_t$, previous hidden state $\mathbf{h}_{t-1}$ and memory state $\mathbf{m}_{t-1}$:

$$\begin{aligned} \mathbf{h}_0, \mathbf{m}_0 &= \left[ \mathbf{W^{ih}}; \mathbf{W^{ic}} \right] Mean(\{\mathbf{v_i}\}) \\ \mathbf{h}_t, \mathbf{m}_t &= LSTM(\mathbf{y}_t, \mathbf{h}_{t-1}, \mathbf{m}_{t-1}) \end{aligned} \qquad (4)$$

  where $\mathbf{y}_t = \mathbf{E}[y_t]$ denotes a word feature of a single word $y_t$. $Mean(\cdot)$ denotes a mean pooling of the given feature set $\mathbf{v_i}$. $\mathbf{W^{ih}}$ and $\mathbf{W^{ic}}$ are parameters that need to be learned.
- Top LSTM Layer. For the top LSTM, it takes the output of the bottom LSTM unit output $\mathbf{h}_t$, previous hidden state
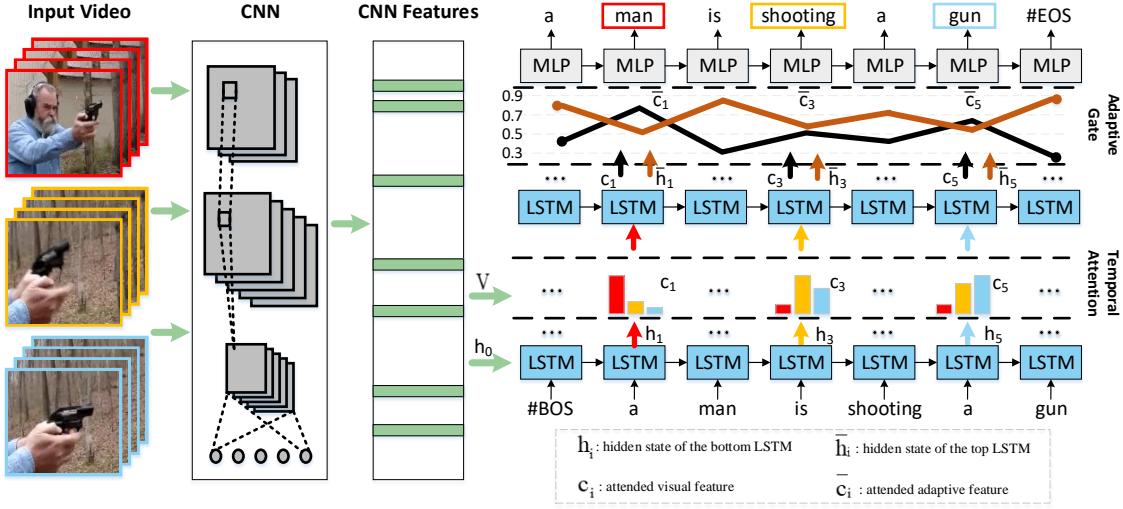
Fig. 2. The framework of our proposed method hLSTMat. To illustrate the effectiveness of our hLSTMat, each generated visual words (i.e., "man", "shooting" or "gun") is generated with visual information extracting from a set of specific frames. For instance, "man" is marked with "red", this indicates it is generated by using the frames marked with red bounding boxes, "shooting" is generated replying on the frames marked with "orange". Other non-visual words such as "a" and "is" are relying on the language model.
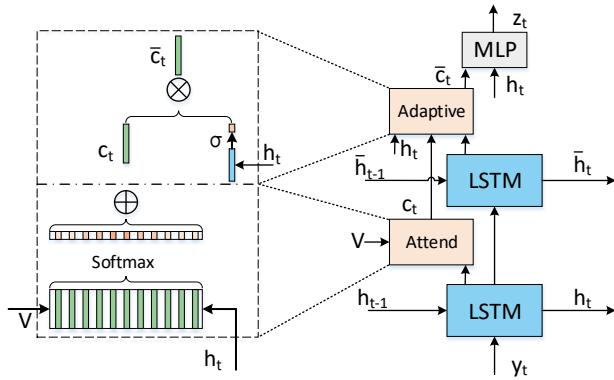


Fig. 3. An illustration of the proposed method generating the t-th target word $z_t$ given a video.

$\bar{\mathbf{h}}_{t-1}$ and the memory state $\bar{\mathbf{m}}_{t-1}$ as inputs to obtain the hidden state $\bar{\mathbf{h}}_t$ at time $t$, and it can be defined as below:

$$\bar{\mathbf{h}}_0 = 0, \bar{\mathbf{m}}_0 = 0$$
$$\bar{\mathbf{h}}_t, \bar{\mathbf{m}}_t = LSTM(\mathbf{h}_t, \bar{\mathbf{h}}_{t-1}, \bar{\mathbf{m}}_{t-1}) \quad (5)$$

- Attention Layers. In addition, for attention based LSTM, context vector is in general an important factor, since it provides meaningful visual evidence for caption generation [14]. In order to efficiently adjust the choose of visual information or sentence context information for caption generation, we defined an adaptive temporal context vector $\bar{\mathbf{c}}_t$ and a temporal context vector $\mathbf{c}_t$ at time $t$. See below:

$$\bar{\mathbf{c}}_t = \psi(\mathbf{h}_t, \bar{\mathbf{h}}_t, \mathbf{c}_t), \quad \mathbf{c}_t = \varphi(\mathbf{h}_t, \mathbf{V}) \quad (6)$$

where $\psi$ denotes the function of our adaptive gate, while $\varphi$ denotes the function of our temporal attention model. Moreover, $\bar{\mathbf{c}}_t$ denotes the final context vector through our adaptive gate, and $\mathbf{c}_t$ represents intermediate vectors

calculated by our temporal attention model. These two attention layers will be described in details in Sec. 4.4 and Sec. 4.6.

- MLP layer. To output a symbol $z_t$, a probability distribution over a set of possible words is obtained using $\mathbf{h}_t$ and our adaptive temporal attention vector $\bar{\mathbf{c}}_t$:

$$\mathbf{p}_t = softmax\left(\mathbf{U}_p\phi(\mathbf{W}_p[\mathbf{h}_t; \bar{\mathbf{c}}_t] + \mathbf{b}_p) + \mathbf{d}\right) \quad (7)$$

where $\mathbf{U}_p$, $\mathbf{W}_p$, $\mathbf{b}_p$ and $\mathbf{d}$ are parameters to be learned. Next, we can interpret the output of the softmax layer $\mathbf{p}_t$ as a probability distribution over words:

$$P(z_t|z_{<t}, \mathbf{V}, \Theta) \quad (8)$$

where $\mathbf{V}$ denotes the features of the corresponding input video, and $\Theta$ are model parameters.

### 4.3 Losses

To learn $\Theta$ in our modal, we minimize the MLE loss defined in Eq. 3. After the parameters are learned, we choose BeamSearch [5] method to generate sentences for videos, which iteratively considers the set of the $k$ best sentences up to time $t$ as candidates to generate sentence of time $t + 1$, and keeps only best k results of them. Finally, we approximates $D = argmax_{D'} Pr(D'|X)$ as our best generated description. In our entire experiment, we set the beam size of BeamSearch as 5.

### 4.4 Temporal Attention Model

As mentioned above, temporal context vector $\mathbf{c}_t$ is an important factor in encoder-decoder framework. To deal with the variability of the length of videos, a simple strategy [12] is used to compute the average of features across a video, and this generated feature is used as input to the model at each time step:

$$\mathbf{c}_t = \frac{1}{L}\sum_{l=1}^{L}\mathbf{v}_l \quad (9)$$

However, this strategy effectively collapses all frame-level features into a single vector, neglecting the inherent temporal structure and leading to the loss of information. Instead of using a simple average strategy (see Eq. 9), we wish to take the dynamic weight sum of the temporal feature vectors according to attention weights $\alpha_t^i$, which are calculated by a soft attention. For each $\mathbf{v}_i$ at $t$ time step, we use the follow function to calculate $\mathbf{c}_t$:

$$\mathbf{c}_t = \frac{1}{L} \sum_{l=1}^{L} \alpha_t^l \mathbf{v}_l \qquad (10)$$

where at $t$ time step $\sum_{l=1}^{L} \alpha_t^l = 1$.

In this paper, we integrate two LSTM layers, a novel temporal attention model for computing the context vector $\mathbf{c}_t$ in Eq. 6 proposed in our framework. Given a set of video features $\mathbf{V}$ and the current hidden state of the bottom layer LSTM $\mathbf{h}_t$, we feed them into a single neural network layer, and it returns an unnormalized relevant scores $\varepsilon_t$. Finally, a softmax function is applied to generate the attention distribute over the $L$ frames of the video:

$$
\begin{aligned}
\varepsilon_t &= \mathbf{W}^{\mathrm{T}} \left( \mathbf{W}_a \mathbf{h}_t + \mathbf{U}_a \mathbf{V} + \mathbf{b}_a \right) \\
\alpha_t &= softmax(\varepsilon_t)
\end{aligned}
\qquad (11)
$$

where $\mathbf{w}^T$, $\mathbf{W}_a$, $\mathbf{U}_a$ and $\mathbf{b}_a$ are parameters to be learned. $\alpha_t \in \mathbb{R}^n$ is the attention weight which quantifies the relevance of features in $\mathbf{V}$. Different from [14], we utilize the current hidden state instead of previous hidden state $\mathbf{h}_t$ generated by the first LSTM layer to obtain the context vector $\mathbf{c}_t$, which focuses on salient feature in the video.

## 4.5 Spatial Attention Model

As temporal attention mentions above, different frames play different roles when generating caption. It is intuitive to think that different regions also have different weights in video captioning. Frame-level image feature extract the global feature of the whole frame. Compared with global feature, spatial features of the image contains more details in the image and we use attention mechanism to select important regions in captioning. For each frame of video, we extract $res5c$ layer from pre-trained ResNet model which denotes the region feature of video frame. For each step, we can calculate the spatial context vector using current hidden state $\mathbf{h_t}$ as:

$$\mathbf{c}_t = \frac{1}{N} \sum_{n=1}^{N} \alpha_t^n \mathbf{r}_n \qquad (12)$$

where $N$ is the number of regions, $\mathbf{r}_i$ is the spatial feature of video. Similar to temporal attention, a simple softmax function is used to generate the spatial attention map as Eq. 11. However, spatial attention model use region feature $R$ instead of frame level feature.

$$
\begin{aligned}
\varepsilon_t &= \mathbf{W}^{\mathrm{T}} tanh \left( \mathbf{W}_a \mathbf{h}_t + \mathbf{U}_a \mathbf{R} + \mathbf{b}_a \right) \\
\alpha_t &= softmax(\varepsilon_t)
\end{aligned}
\qquad (13)
$$

R denotes spatial feature extracted from convolution neural network, which represents region feature of the image.

TABLE 1
The operation modules and their number used in our three proposed models.

| model | Bot LSTM | Attention | Top LSTM | Adap ATT |
|---|---|---|---|---|
| hLSTMat | 1 | 1 | 1 | 1 |
| ParA | 1 | 2 | 1 | 1 |
| Two Stream | 2 | 2 | 2 | 1 |

## 4.6 Adaptive Temporal Attention Model

In this paper we propose an adaptive temporal attention model to compute a context vector $\bar{\mathbf{c}}_t$ in Eq. 6, shown in Fig. 3, to encourage that a decoder uses nearly no visual information from video frames to predict the non-visual words, and use the most related visual information to predict visual words. In our hierarchical LSTM network, the hidden state in the bottom LSTM layer is a latent representation of what the decoder already knows. With the hidden state $\mathbf{h}_t$, we extend our temporal attention model, and propose an adaptive model that is able to determine whether it needs to attend the video to predict the next word. In addition, a sigmoid function is applied to the hidden state $\mathbf{h}_t$ to further filter visual information.

$$
\begin{aligned}
\bar{\mathbf{c}}_t &= \beta_t \mathbf{c}_t + (1 - \beta_t)\bar{\mathbf{h}}_t \\
\beta_t &= sigmoid(\mathbf{W}_s \mathbf{h}_t)
\end{aligned}
\qquad (14)
$$

where $\mathbf{W}_s$ denotes the parameters to be learned and $\beta_t$ is adaptive gate at each time $t$. In our adaptive temporal attention model, $\beta_t$ is projected into the range of $[0, 1]$. When $\beta_t = 1$, it indicates that full visual information is considered, while when $\beta_t = 0$ it indicates that no visual information is considered to generate the next word.

## 4.7 Multiple Features based hLSTMat

Some visual words (e.g., "gun" and "dog") and verbs (e.g., "shooting" and "eating") can be identified from a still image from the appearance alone, while others, (e.g., distinguishing "yawing" from "singing" or "walking" from "running") require motion cues as well. The hLSTMat using single appearance feature cannot accurately generate video captions. In this section, we consider three architectures for fusing appearance and motion features.

**hLSTMat with concatenation fusion (ConF).** In order to make use of two types of features (i.e., appearance and motion cues), a concatenation fusion is introduced by firstly concatenating appearance and motion features to form a new video feature, see Fig. 4. In other words, the hLSTMat takes an integrated video feature as input to generate video descriptions.

**Two-stream hLSTMat Networks (Two-stream).** Two-stream architecture has been widely used in video recognition, and it can achieve very good performance [47]. Motivated by this, we utilize two different pre-trained hLSTMat models. Each hLSTMat model is trained by taking one type of video feature as the input. As a result, each model separately generates a distribution of the words. For two models, their generated distributions are represented as $\mathbf{p}_{t1}$ and $\mathbf{p}_{t2}$, respectively. To predict the word at time $t$, it uses the following function Eq. 15 to calculate the final distribution $\mathbf{p}_t$:

$$\mathbf{p}_t = \frac{1}{2}(\mathbf{p}_{t1} + \mathbf{p}_{t2}) \qquad (15)$$

Here, we devise our video captioning architecture accordingly, dividing it into two streams, as shown in Fig. 4. Each stream is implemented using a hLSTMat network.
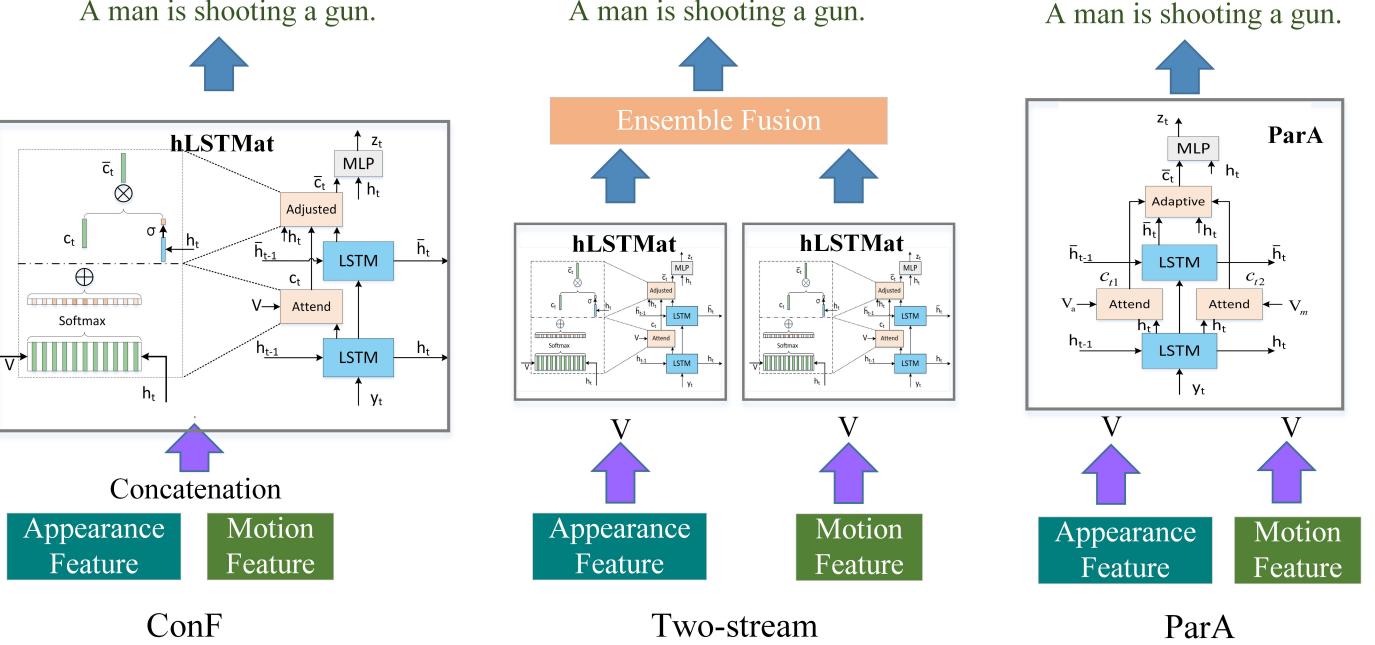
Fig. 4. We construct variants of architectures of spatial-temporal hLSTMat. The left indicates the hLSTMat with concatenation fusion (ConF). The middle one represents the Two-stream hLSTMat Networks (Two-stream). The right one illustrates the architectures of parallel adaptive temporal attention model (ParA).

**Parallel Adaptive Temporal Attention Model (ParA).** Two-stream hLSTMat Networks fuse the distribution scores using averaging, and the experimental results show that two-stream hLSTMat significantly outperforms hLSTMat. However, a limitation of two-stream hLSTMat is that it doubles the training time and the number of learning parameters. In order to solve this problem, we propose a parallel adaptive temporal attention model, which enables hLSTMat to make use of complementary information for word prediction. The opreation modules, which we use in our hLSTMat, ParA and Two-stream model, are listed as Tab. 1. It is obvious that Two-stream model uses double training parameters compared with hLSTMat and there are two attention modules in ParA to deal with two kinds of features. From this table, both of Two-stream and ParA use two attention modules. But for other modules, Two-stream model uses double parameter compared with ConF. The model is shown in Fig. 4.

As illustrated in Fig. 4, appearance ($\mathbf{V}_{static}$) and motion ($\mathbf{V}_{motion}$) features are sent to two parallel *attend* components to extract meaningful appearance and motion evidences. Next, we propose an adaptive temporal model to compute the context vector, which is defined as $\bar{\mathbf{c}}_t$. It is a mixture of the video spatial and temporal features as well as the language context information. Mathematically, the ParA model is defined as follows:

$$\bar{\mathbf{c}}_t = \beta_{t1}\mathbf{c}_{t1} + \beta_{t2}\mathbf{c}_{t2} + \beta_{t3}\bar{\mathbf{h}}_t$$
$$[\beta_{t1}, \beta_{t2}, \beta_{t3}] = softmax(\mathbf{W}_s\mathbf{h}_t) \qquad (16)$$

where $\beta_{t1}$, $\beta_{t2}$, $\beta_{t3}$ are the new sentinel gates at time $t$. $\mathbf{c}_{t1}$ and $\mathbf{c}_{t2}$ represents the "appearance sentinel" vector and the "motion sentinel" vector, respectively. If $\beta_{t1} + \beta_{t2} = 1$, then $\beta_{t3} = 0$. This implies that no visual information is considered, and only the evidence of appearance and motion are considered. In other words, the model can predict visual words and non-visual words on the basis of weight of appearance, motion and language context

information. In addition, all the sentinel gates are computed by a softmax, where $\mathbf{W}_s$ are parameters to be learned.

In terms of initialization of $\mathbf{h}_0$, we define it as follows:

$$\mathbf{h}_0, \mathbf{m}_0 = \left[\mathbf{W}^{ih}; \mathbf{W}^{ic}\right] Mean(\{\mathbf{v_{static}}, \mathbf{v_{motion}}\}) \qquad (17)$$

where $\mathbf{W}^{ih}$ and $\mathbf{W}^{ic}$ are parameters that need to be learned. Note that this model can be easily extended to support multiple features.

## 5 EXPERIMENTS FOR VIDEO CAPTIONING

We evaluate our algorithm on the task of video captioning. Specifically, we firstly study the influence of CNN encoders. Secondly, we explore the effectiveness of the proposed components. Next, we compare our results with the state-of-the-art methods on three benchmark datasets.

### 5.1 Datasets

We consider three publicly available datasets that have been widely used in previous work.

**The Microsoft Video Description Corpus (MSVD).** This video corpus consists of 1,970 short video clips, approximately 80,000 description pairs and about 16,000 vocabulary words [48]. Following [13], [14], we split the dataset into training, validation and testing set with 1,200, 100 and 670 videos, respectively.

**MSR Video to Text (MSR-VTT).** In 2016, Xu *et al.* [49] proposed the largest video benchmark for video understanding, and especially for video captioning. This dataset contains 10,000 web video clips, and each clip is annotated with approximately 20 natural language sentences. In addition, it covers the most comprehensive categorizes (i.e., 20 categories) and a wide variety of visual content, and contains 200,000 clip-sentence pairs.

**Large Scale Movie Description Challenge (LSMDC).** It is a dataset of movies with aligned descriptions sourced from movie

scripts. Also, it is based on the two previous datasets: the MPII Movie Description dataset (MPII-MD [51]) and the Montreal Video Annotation Dataset (M-VAD [52]). In total, this dataset contains a parallel corpus of 118,114 sentences and 118,081 video clips from 202 movies. In our experiments, it is divided into training, validation, public test, and blind test set, and each contains 10,1079, 7,408, 10,053, and 9,578 videos, respectively.

## 5.2 Implementation Details

### 5.2.1 Preprocessing

For MSVD dataset, we convert all descriptions to lower cases, and then use wordpunct_tokenizer method from the NLTK toolbox to tokenize sentences and remove punctuations. Therefore, it yields a vocabulary of 13,010 in size for the training split. For MSR-VTT dataset, captions have been tokenized, thus we directly split descriptions using blank space, thus it yields a vocabulary of 23,662 in size for training split. For LSMDC, vocabulary is composed of 25,610 words, which are extracted from captions by the NLTK toolbox. Inspired by [14], we preprocess each video clip by selecting equally-spaced 28 frames out of the first 360 frames and then feeding them into a CNN network proposed in [28]. Thus, for each selected frame we obtain a 2,048-D feature vector, which are extracted from the $pool5$ layer. As for spatial feature, we extract region features from $res5c$ layer.

### 5.2.2 Training details

In the training phase, in order to deal with sentences with arbitrary length, we add a begin-of-sentence tag <BOS> to start each sentence and an end-of-sentence tag <EOS> to end each sentence. In the testing phase, we input <BOS> tag into our attention-based hierarchical LSTM to trigger video description generation process. For each word generation, we choose the word with the maximum probability and stop until we reach <EOS>.

In addition, all the LSTM unit sizes are set as 512 and the word embedding size is set as 512, empirically. Our objective function Eq. 3 is optimized over the whole training video-sentence pairs with mini-batch 64 in size of MSVD 256 in size of MSR-VTT and LSMDC. We adopt adadelta [53], which is an adaptive learning rate approach, to optimize our loss function. In addition, 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 epochs are reached, or until the evaluation metric does not improve on the validation set at the patience of 20.

### 5.2.3 Evaluation metrics

To evaluate the performance, we employ three different standard evaluation metrics: BLUE [54], METEOR [55], and CIDEr [56].

## 5.3 The Effect of Different CNN Encoders

To date, there are 6 widely used CNN encoders including C3D, GoogleNet, Inception-V3, ResNet-50, ResNet-101 and ResNet-152 to extract visual features. In this sub-experiment, we study the influence of different CNN encoders on our framework. The experiments are conducted on the MSVD and MSR-VTT datasets, and the results are shown in Tab. 2 and Tab. 3.

From Tab. 2, we find that by taking ResNet-152 as the visual decoder, our method performs the best with 82.9% B@1, 72.2% B@2, 63.0% B@3, 53.0% B@4 and 33.6% METEOR, while Inception-v3 is a strong competitor, with 82.7% B@1,

TABLE 2
Experiment results on the MSVD dataset. We use different features to verify our hLSTMat method.

| Model | B@1 | B@2 | B@3 | B@4 | METEOR |
|---|---|---|---|---|---|
| C3D | 79.9 | 68.2 | 58.3 | 47.5 | 30.5 |
| GoogleNet | 80.8 | 68.6 | 58.9 | 48.5 | 31.9 |
| Inception-v3 | 82.7 | 72.0 | 62.5 | 51.9 | 33.5 |
| ResNet-50 | 80.9 | 69.1 | 59.5 | 49.0 | 32.3 |
| ResNet-101 | 82.2 | 70.9 | 61.4 | 50.8 | 32.7 |
| ResNet-152 | **82.9** | **72.2** | **63.0** | **53.0** | **33.6** |

TABLE 3
Experimental Results on MSR-VTT dataset. We use different features (i.e., C3D, Inception-v3 and ResNet-152) to further verify our hLSTMat.

| Model | B@1 | B@2 | B@3 | B@4 | METEOR | CIDEr |
|---|---|---|---|---|---|---|
| C3D | 75.4 | 61.0 | 48.3 | 36.7 | 25.8 | 37.8 |
| Inception-v3 | 72.7 | 60.3 | 48.7 | 38.2 | 26.1 | 43.1 |
| ResNet-152 | **73.9** | **61.6** | **49.5** | **38.4** | **26.3** | **43.4** |

72.0% B@2, 62.5% B@3, 51.9% B@4 and 33.5% METEOR. However, the gap between ResNet-152 and Inception-v3 is very small. Compared with other appearance features, C3D obtains the lowest scores for both BLUE and METEOR. Unlike GoogleNet, Inception and ResNet, which can be pre-trained on a large still image classification dataset (e.g., ImageNet), the C3D is trained on video dataset [30], and the available datasets for video action classification (e.g., Sport1M and UCF101) are still rather small. Therefore, we can conclude that in general, GoogleNet is worse than ResNet, and for ResNet, deeper networks result in a better performance. Also, the appearance features outperform the motion features for video captioning.

To further evaluate the performance of different features, we conducted more experiments on the MSR-VTT dataset by comparing C3D, Inception-v3 and ResNet-152, and show the results in Tab. 3. We can have similar observation to Tab. 2 that ResNet-152 performs the best and C3D obtains the worst results.

## 5.4 Architecture Exploration and Comparison

In this sub-experiment, we explore the impact of three proposed components, the hierarchical LSTMs, temporal attention, and adaptive attention. Specifically, we compare the basic LSTM proposed in Sec.3.1 "basic LSTM", "basic+adaptive attention" which adds adaptive attention to the basic LSTM, "hLSTMt" which removes the adaptive mechanism from the hLSTMat, and "hLSTMat", as well as comparing them with the state of the art methods: MP-LSTM [12] and SA [14]. In order to conduct a fair comparison, all the methods take ResNet-152 as the encoder. We conduct the experiments on the MSVD dataset and the results are shown in Tab. 4.

Tab. 4 shows that in general, our hLSTMat achieves the best results in all metrics with 82.9% B@1, 72.2% B@2, 63.0% B@3, 53.0% B@4, 33.6% METEOR and 73.8% CIDEr. By comparing "basic LSTM" with "SA" which adds temporal attention to the basic LSTM, we can observe that the temporal attention can improve the performance of video captioning. Specifically, the performance is improved by 1% B@1, 1% B@2, 1.9% B@3, 1.7% B@4, 0.6% METEOR and 2.1% CIDEr. Moreover, by comparing "basic LSTM" with "basic+adaptive attention", and "hLSTMt" with "hLSTMt", we find that adaptive attention plays an important role for video captioning. This indicates that the

TABLE 4
The effect of different components and the comparison with the state-of-the-art methods on the MSVD dataset. The default encoder for all methods is ResNet-152.

| Model | B@1 | B@2 | B@3 | B@4 | METEOR | CIDEr |
|---|---|---|---|---|---|---|
| basic LSTM | 80.6 | 69.3 | 59.7 | 49.6 | 32.7 | 69.9 |
| MP-LSTM [12] | 81.1 | 70.2 | 61.0 | 50.4 | 32.5 | 71.0 |
| SA [14] | 81.6 | 70.3 | 61.6 | 51.3 | 33.3 | 72.0 |
| basic+adaptive attention | 80.9 | 69.7 | 61.1 | 50.2 | 31.6 | 71.5 |
| hLSTMt | 82.5 | 71.9 | 62.0 | 52.1 | 33.3 | 73.5 |
| hLSTMat | **82.9** | **72.2** | **63.0** | **53.0** | **33.6** | **73.8** |

TABLE 5
The Performance of three variants architectures of temporal based hLSTMat and spatial based hLSTMat on the MSVD dataset. C, R indicate C3D and ResNet-152, respectively. T, S indicates Temporal attention and Spatial attention.

| Model | B@1 | B@2 | B@3 | B@4 | METEOR | CIDEr |
|---|---|---|---|---|---|---|
| Motion(C) | 79.9 | 68.2 | 58.3 | 47.5 | 30.5 | 59.5 |
| Appearance(R) | 82.9 | 72.2 | 63.0 | 53.0 | 33.6 | **73.8** |
| ConF (C+R) | 81.9 | 71.5 | 62.3 | 52.5 | 32.2 | 65.1 |
| ParA (C+R) | 82.5 | 72.4 | 63.7 | 53.6 | 32.3 | 66.8 |
| Two-stream (C+R) | **83.4** | **73.4** | **64.1** | **54.0** | **33.2** | 71.3 |
| Temporal-Resnet | 82.9 | 72.2 | 63.0 | 53.0 | 33.6 | **73.8** |
| Spatial-Resnet | 82.7 | 72.7 | 63.2 | 53.5 | 33.2 | 72.3 |
| ParA (S+T) | 83.1 | 73.1 | 63.7 | 53.7 | **33.9** | 69.8 |
| Two-stream (S+T) | **83.3** | **73.6** | **64.6** | **54.3** | 33.5 | 72.8 |

adaptive attention mechanism deciding whether to look at the visual information or language context information is beneficial for more accurate caption prediction. However, the performance improvement of adaptive attention is not as significant as temporal attention. By comparing "hLSTMt" with "SA", it is observed that the a significant improvement is achieved by the hierarchical LSTMs. The performance is improved by 0.9% B@1, 1.6% B@2, 0.4% B@3, 0.8% B@4, 0% METEOR and 1.5% CIDEr.

## 5.5 The Effect of Spatial-temporal hLSTMat

In this section, we aim to evaluate the performance of our proposed multiple features based hLSTMat: Two-stream and ParA, both introduced in Sec. 4.7. We compare them with baselines: hLSTMat and ConF which is the simplest extension of hLSTMat. C3D (C) ResNet152 (R) are utilized in ours experiments. By default, we use temporal attention and ResNet152 for video captioning. In the first group of experiments, we use C3D (C) ResNet152 (R) as multiple features, and in the second group, we utilize temporal (T) and spatial attention (S) on ResNet152 features as multiple features. The experiments are conducted on two datasets: MSVD and MSR-VTT, and the results are shown in Tab. 5 and Tab. 6 respectively.

In the first part of Tab. 5, we observe that in general, combining multiple features can achieve better performance than using single feature for most of the evaluation metrics of video captioning. Interestingly, using the single appearance feature achieves the best performance for CIDEr. One possible reason is that.... Another observation is that Two-stream consistently outperforms ParA and ConF for all evaluation metrics.

The experimental results show that Two-stream performs best, ParA goes to second, and the ConF performs the worst. However, in terms of the number of training parameters, which needs to be learned, and the training time, Two-stream doubles ParA. In Section 4.3, we find out that C3D performs worse than all the appearance features, due to lacking of enough training datasets.

TABLE 6
The Performance of three variants architectures of temporal based hLSTMat and spatial based hLSTMat on the MSR-VTT dataset. C, R indicate C3D and ResNet-152, respectively. T, S indicates Temporal attention and Spatial attention.

| Model | B@1 | B@2 | B@3 | B@4 | METEOR | CIDEr |
|---|---|---|---|---|---|---|
| Motion(C) | 75.4 | 61.0 | 48.3 | 36.7 | 25.8 | 37.8 |
| Appearance(R) | 73.9 | 61.6 | 49.5 | 38.4 | 26.3 | **43.4** |
| ConF (C+R) | 74.3 | 61.7 | 49.1 | 37.4 | 25.9 | 41.6 |
| ParA (C+R) | 74.5 | 61.1 | 48.8 | 38.1 | 26.0 | 42.3 |
| Two-stream (C+R) | **75.2** | **62.7** | **50.3** | **38.7** | **26.8** | 41.9 |
| Temporal-Resnet | 73.9 | 61.6 | 49.5 | 38.4 | 26.3 | **43.4** |
| Spatial-Resnet | 74.2 | 60.9 | 48.8 | 38.3 | 26.3 | 40.8 |
| ParA (S+T) | 75.2 | 61.4 | 48.9 | 38.0 | 26.4 | 39.7 |
| Two-stream (S+T) | **76.2** | **62.9** | **50.6** | **39.7** | **27.0** | 42.1 |

TABLE 7
The performance comparison with the state-of-the-art methods on MSVD dataset. (V) denotes VGGNet, (O) denotes optical flow, (G) denotes GoogleNet, (C) denotes C3D and (R) denotes ResNet-152.

| Model | B@1 | B@2 | B@3 | B@4 | METEOR | CIDEr |
|---|---|---|---|---|---|---|
| Random Choice | 43.3 | 21.0 | 11.6 | 5.1 | 12.8 | 1.5 |
| S2VT(V) [13] | - | - | - | - | 29.2 | - |
| S2VT(V+O) | - | - | - | - | 29.8 | - |
| HRNE(G) [32] | 78.4 | 66.1 | 55.1 | 43.6 | 32.1 | - |
| HRNE-SA (G) | 79.2 | 66.3 | 55.1 | 43.8 | 33.1 | - |
| LSTM-E(V) [57] | 74.9 | 60.9 | 50.6 | 40.2 | 29.5 | - |
| LSTM-E(C) | 75.7 | 62.3 | 52.0 | 41.7 | 29.9 | - |
| LSTM-E(V+C) | 78.8 | 66.0 | 55.4 | 45.3 | 31.0 | - |
| p-RNN(V) [33] | 77.3 | 64.5 | 54.6 | 44.3 | 31.1 | 62.1 |
| p-RNN(C) | 79.7 | 67.9 | 57.9 | 47.4 | 30.3 | 53.6 |
| p-RNN(V+C) | 81.5 | 70.4 | 60.4 | 49.9 | 32.6 | 65.8 |
| LSTM-TSA$_I$ [39] | 81.0 | 69.6 | 60.2 | 50.2 | 32.4 | 71.5 |
| LSTM-TSA$_V$ | 82.1 | 70.7 | 61.1 | 50.5 | 32.6 | 71.7 |
| LSTM-TSA$_{I,V}$ | 82.8 | 72.0 | 62.8 | 52.8 | 33.5 | **74.0** |
| hLSTMat (T) | 82.9 | 72.2 | 63.0 | 53.0 | **33.6** | 73.8 |
| ParA (S+T) | 83.1 | 73.1 | 63.7 | 53.7 | 33.9 | 69.8 |
| **Two-stream (S+T)** | **83.3** | **73.6** | **64.6** | **54.3** | 33.5 | 72.8 |

In other words, the potential of spatial-temporal hLSTMat networks is not fully expressed. Therefore, we conduct the third experiment by replacing C3D with Inception-v3 and the results are demonstrated in Tab. 5. By observing the experimental results, we find out that Two-stream (I+R) performs best with all the evaluation metrics (i.e., 84.1% B@1, 74.9% B@2, 66.2% B@3, 56.2% B@4, 34.7% METEOR and 81.1% CIDEr). In addition, compared with ConF (C+R) and ParA (C+R), ConF (I+R) and ParA (I+R) increases the performance, respectively. Following this thought, we run spatial-temporal hLSTMat on the MSR-VTT and the results are shown in Tab. 6. From which, we can see that Two-stream (I+R) achieves the best results with 77.1% B@1, 64.0% B@2, 51.3% B@3, 40.0% B@4, 27.3% METEOR and 45.7% CIDEr.

## 5.6 Compare with the-state-of-the-art Methods

### 5.6.1 Results on MSVD dataset

In this subsection, we show the comparison of our approach with the baselines on the MSVD dataset. Some of the above baselines only utilize video features generated by a single deep network, while others (i.e., S2VT, LSTM-E and p-RNN) make uses of both single network and multiple network generated features. Therefore, we first compare our method with approaches using static frame-level features extracted by a single network. In addition, we compare our spatial-temporal approaches with methods utilized
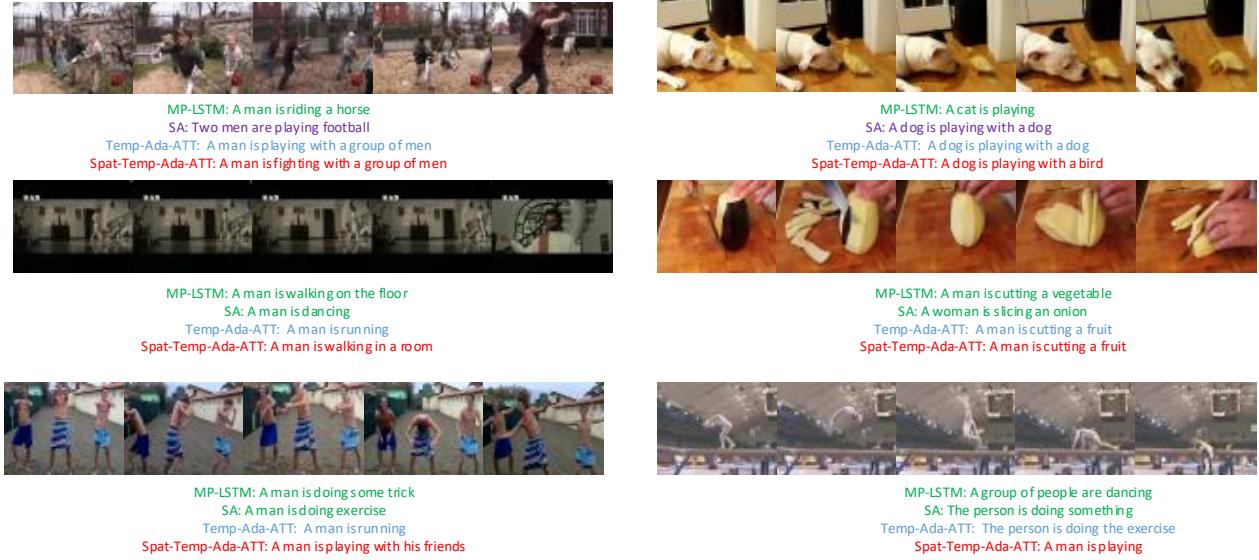
Fig. 5. Examples of sentences produced by MP-LSTM and the proposed method: Temporal Adaptive attention model and Spatial Temporal Adaptive attention. Each video is represented by five frames.

TABLE 8
The performance comparison with the state-of-the-art methods on MSR-VTT dataset.

| Model | B@4 | METEOR |
|---|---|---|
| Random Choice | 3.2 | 10.5 |
| MP-LSTM (V) | 34.8 | 24.8 |
| MP-LSTM (C) | 35.4 | 24.8 |
| MP-LSTM (V+C) | 35.8 | 25.3 |
| SA (V) | 35.6 | 25.4 |
| SA (C) | 36.1 | 25.7 |
| SA (V+C) | 36.6 | 25.9 |
| Rank-1 v2t_navigator | 40.8 | 28.2 |
| Rank-2 Aalto | 39.8 | 28.2 |
| Rank-3 VideoLAB | 39.1 | 27.7 |
| Rank-4 ruc-uva | 38.7 | 26.9 |
| Rank-5 Fudan-ILC | 38.7 | 26.8 |
| hLSTMat (R) | 38.3 | 26.3 |
| ParA(S+T) | 38.5 | 25.7 |
| Two-stream (S+T) | **39.7** | **27.0** |

TABLE 9
The performance comparison with the state-of-the-art methods on LSMDC dataset.

| Methods | CIDEr | B@4 | METEOR | ROUGE |
|---|---|---|---|---|
| BUPT CIST AI lab | 0.072 | 0.005 | **0.075** | 0.152 |
| IIT Kanpur | 0.042 | 0.004 | 0.070 | 0.138 |
| Aalto University | 0.037 | 0.002 | 0.033 | 0.069 |
| Shetty and Laaksonen | 0.044 | 0.003 | 0.046 | 0.108 |
| S2VT | 0.082 | 0.007 | 0.070 | 0.149 |
| Base-SAN | 0.090 | 0.005 | 0.066 | 0.150 |
| CT-SAN | 0.100 | **0.008** | 0.071 | **0.159** |
| hLSTMat(I) | 0.097 | 0.007 | 0.057 | 0.150 |
| hLSTMat(R) | 0.102 | 0.007 | 0.058 | 0.147 |
| Two-stream (R+I) | **0.104** | 0.007 | 0.056 | 0.146 |

two deep features. The results are shown in Tab. 7 and we have the following observations:

- Compared with the best counterpart (i.e., p-RNN) which

only takes spatial information, our method has 8.7% improvement on B@4 and 2.5% on METEOR.

- The hierarchical structure in HRNE reduces the length of input flow and composites multiple consecutive input at a higher level, which increases the learning capability and enables the model encode richer temporal information of multiple granularities. Our approach (53.0% B@4, 33.6% METEOR) performs better than HRNE (43.6% B@4, 32.1% METEOR) and HRNE-SA (43.8% B@4, 33.1% METEOR). This shows the effectiveness of our model.

- Our hLSTMat (53.0% B@4, 33.6% METEOR) can achieve better results than our hLSTMt (52.1% B@4, 33.3% METEOR). This indicates that it is beneficial to incorporate the adaptive temporal attention into our framework.

- Our hLSTMat achieves a good performance (53.0% B@4 and 33.6% METEOR) using static frame-level features compared with approaches combining multiple deep features. For S2VT (V+O), LSTM-E (V+C) and p-RNN (V+C), they use two networks VGGNet/GoogleNet and optical flow/C3D to capture videos' spatial and temporal information, respectively. Compared with them, our hLSTMat only utilizes ResNet-152 to capture frame-level features, which proves the effectiveness of our hierarchical LSTM with adaptive temporal attention model.

In addition, simultaneously considering both spatial and temporal video information can enhance the video caption performance. VGGNet and GoogleNet are used to extract spatial information, while optical flow and C3D are utilized for capturing temporal information. For example, compared with LSTM-E (V) and LSTM-E (C), LSTM-E (V+C) achieves higher B@4 (45.3%) and METEOR (31.0%). For p-RNN, p-RNN (V+C) (49.9% B@4 and 32.6% METEOR ) performs better than both p-RNN(V) (44.3% B@4 and 31.3% METEOR) and p-RNN(C) (47.4% B@4 and 30.3% METEOR). Compared with all two features based
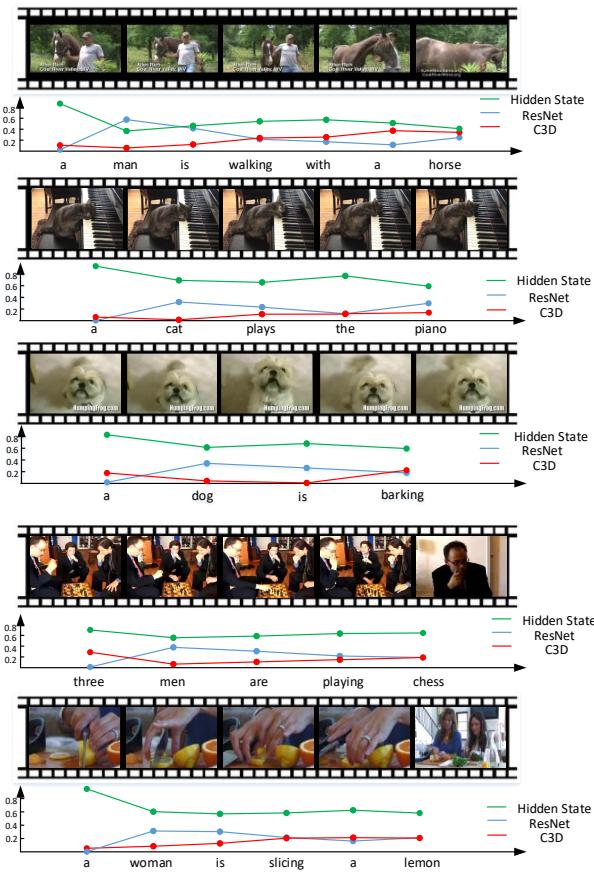
Fig. 6. Examples of attention weights changes along with the generation of captions.

approaches LSTM-E (V+C), p-RNN (V+C) and LSTM-TSA$_{I,V}$, our Two-stream (R+I) achieves the best performance with 84.1% B@1, 74.9% B@2, 66.2% B@3, 56.2% B@4, 34.7% METEOR and 81.1% CIDEr.

We adopt questionnaires collected from ten users with different academic backgrounds. Given a video caption, users are asked to score the following aspects: 1) Caption Accuracy, 2) Caption Information Coverage, 3) Overall Quality. Results show that our method outperforms others at 'Overall Quality', and 'Caption Accuracy' with small margin. But it has lower value for 'Information coverage' than p-RNN.

### 5.6.2 Results on MSR-VTT dataset

We compare our model with the state-of-the-art methods on the MSR-VTT dataset, and the results are shown in Tab. 8. Our Two-stream model performs the best on all metrics with 40.0% @B4 and 27.3% METEOR, while our hLSTMat (R) goes to second with 38.3% @B4 and 26.3% METEOR. Compared with our method hLSTMat using only temporal attention, the performance is improved by 1.7% for @B4, and 1.0% for METEOR. This verifies the effectiveness of our hLSTMat approach. In addition, our two-stream (R+I) increase the state-of-the-art method SA (V+C) by 3.4% on B@4 and 1.4% METEOR. This proofs the importance of two-stream networks. Besides, We compared our method with the MSR-VTT 2016 challenge results and we can see that our method can attain comparable performance. The

leaderboard of MSR-VTT 2016 challenge can be easily found on the Internet. [1]

### 5.6.3 Results on LSMDC dataset

In this section, we report the performance of different methods on the LSMDC dataset, see Tab. 9. All the previous methods are LSMDC participants, which have no obligation to disclose their identities or used technique. Among comparable methods, our approach Two-stream obtains the highest CIDEr (10.4%). In terms of B@4 and ROUGE, CT-SAN ranks first with 0.8% and 15.9%, respectively. And BUPT CIST AI lab achieves best performance on METEOR. In addition, compared with MSR-VTT (i.e., each video contains approximately 20 natural language sentences), the LSMDC dataset does not contain sufficient descriptions for each video. For 118,081 video clips, it only has 118,114 sentences. In average, each video only has one language description. However, in reality, human can provides multiple natural language sentences for describing a video. In other words, a video can be described from different aspects, but the the ground truth only provides one description. This might be the major reason why all the methods gain low performance scores on this dataset in terms of all evaluation metrics.

### 5.6.4 Qualitative analysis

In this section, we display several examples of actual videos in different methods. Here, we use three model MP-LSTM, our proposed Temporal Adaptive attention and Spatial-Temporal Adaptive attention to generate sentence and their sentences in green, blue and read respectively. Compared with

## 5.7 Feature Weights Analysis

In this section, we discuss the feature weights learned by the spatial-temporal framework: ParA by demonstrating some examples in Fig. 6. More specifically, we aim to answer how ParA leverages the three important factors including the motion feature (C3D), appearance feature (ResNet152) and language context information (i.e., hidden state value), and what are the their weight trends in a captioning generation process?

Fig. 6 provides five video captioning examples and each shows the weights change of ResNet152 (R), C3D (C) and hidden state (H). The generated natural language sentences are shown under the horizontal time axis, and each word is positioned at the time step when it is generated. The vertical axis stands for the weight scores and their sum weights equals to 1. From the five examples, we can see that the syntactic information is quite essential to captioning and the weight of language content information is always larger than other two feature weights. This is because the H not only contains the syntactic information but also contains the previous appearance and motion information.

If we look at the trend over time, we can see that 1) when generating non-visual words such as "a", "the","are" and "is", the over all trend of H is upwards. This indicates that context information plays a more important role in non-visual words generation. In the second example, weights of H is quite high (i.e., both more than 0.7) for generating "a" and "the". 2) when generating the subjective in a sentence, such as "man", "cat","dog", "men" or "women", the trend of R grows dramatically, while the trend of C goes down. 3) when generating actions such as "plays the

---

1. http:http://ms-multimedia-challenge.com/2016/leaderboard

piano", "playing chess" and "slicing a lemon", the trend of C is upwards. In general, by comparing appearance and motion features, an interesting phenomenon is found: appearance feature's weight tends to be larger in generating nouns such as "man", "cat", "dog" and "horse", while motion feature's weight tends to be higher in generating verbs such as "walking", "plays", "barking" and "slicing.

## 6  HLSTMAT FOR IMAGE CAPTIONING

In this subsection, we instantiate our hLSTMat and apply it to the task of image captioning. As shown in Fig. 7, there are three major components: 1) a CNN Encoder, 2) an attention based hierarchical LSTM decoder and 3) the losses. We give the details for each component below.

### 6.1  CNN Encoder

In this paper, we adopt a pretrained ResNet-101 [59] (pool5) to extract global visual feature, and use Faster R-CNN [60] to produce bounding boxes and then apply them on the pretrained ResNet-101 to extract $L$ region features. For simplicity, given an image $\mathbf{x}$, we extract $L+1$ image features, represented as $\{\mathbf{v}_g, \mathbf{V}\}$, where $\mathbf{v}_g$ is the global visual feature and $\mathbf{V} = \{\mathbf{v}_1, \cdots, \mathbf{v}_L\}$ indicates the $L$ region visual features.

### 6.2  Hierarchical LSTM with Adaptive Attention

Our hLSTMat consists of two LSTM blocks. The first LSTM aims to prepare the hidden states and visual attention for generating a preliminary version of the captions, while the second LSTM is applied as a proofreading process to refine them. Both LSTM are built upon the basic LSTM (Eq. 1).

For the $t$-th time step, $\mathbf{y}_t$ is the input vector of the LSTM unit, $\mathbf{h}_t$ is the output vector of the LSTM unit, and $\mathbf{h}_{t-1}$ is the output of the LSTM unit at the $t$-1 time step. For simplicity, we refer to the operation procedure of basic LSTM with the following notation:

$$\mathbf{h}_t = LSTM(\mathbf{y}_t, \mathbf{h}_{t-1}) \tag{18}$$

#### 6.2.1  The First LSTM

We modify the basic LSTM to generate an initial text sequence feature as depicted in Fig. 7. We define the input of the LSTM unit as:

$$\mathbf{y}_t^1 = [\mathbf{v}_g, \mathbf{h}_{t-1}^2, \mathbf{w}_t] \tag{19}$$

where $\mathbf{v}_g$ is the global image visual feature, $\mathbf{h}_{t-1}^2$ is the output of the second LSTM layer, at the $t$-1 time step, and $\mathbf{w}_t$ is the feature of the current word derived by embedding an one-hot word vector. It is easy to understand that the current hidden states are based on the global visual feature, the previous hidden states and the $t$-th word. We further use $\mathbf{h}_{t-1}^2$ from a higher-level LSTM in order to utilize the more precise information to guide the learning of $\mathbf{h}_t^1$. The we get $\mathbf{h}_t^1 = LSTM(\mathbf{y}_t^1, \mathbf{h}_{t-1}^1)$.

Traditionally, the hidden state of the LSTM is directly applied to guide which region should be focused on. The LSTM provides a temporal shortcut path to avoid vanishing gradients. Here, we provide an additional word shortcut from the $t$-th high frequency word for efficient training. As a result, we use a residual shortcut connection to further reduce vanishing gradients:

$$\tilde{\mathbf{h}}_t^1 = \mathbf{W}_{rd}[\mathbf{w}_t; \mathbf{h}_t^1] \tag{20}$$

where $\mathbf{h}_t^1$ is the hidden states at the $t$-th step, $\mathbf{W}_{rd}$ is the parameter to be learned, and [;] is a concatenation operation.

Given $L$ image region visual features $\{\mathbf{v}_1, \cdots, \mathbf{v}_L\}$ and the hidden state $\tilde{\mathbf{h}}_t^1$, we aim to selectively utilize certain region visual features by defining a visual attention mechanism below:

$$\mathbf{e}_t^1 = \mathbf{w}_{e1}^T \tanh(\mathbf{W}_{v1}\mathbf{V}, \mathbf{W}_{h1}\tilde{h}_t^1) \tag{21}$$

$$\alpha_t^1 = soft\max(\mathbf{e}_t^1) \tag{22}$$

where $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2 ..., \mathbf{v}_L]$ indicates the $L$ image region features. And $\mathbf{w}_{e1}^T, \mathbf{W}_{v1}, \mathbf{W}_{h1}$ are parameters to be learned. $\alpha_t^1 \in \mathbb{R}^L$ is the attention weight, which is then applied on region features to locate the important visual information:

$$\hat{\mathbf{v}}_t^1 = \sum_{l=1}^{L} \alpha_{l,t}^1 \mathbf{v}_l \tag{23}$$

where $\hat{\mathbf{v}}_t^1$ is the attended visual information which can be used together with $\tilde{\mathbf{h}}_t^1$ to generate the primary $t$-th word.

#### 6.2.2  The Second LSTM with Adaptive Attention

By integrating the softmax layer and the loss functions to the first LSTM layer, we can generate a preliminary word at each step. In this subsection, we design a second LSTM layer with adaptive attention as a deliberate process to further purify the captioning results. Similar to Eq. 14, we first define the language context information $s_t$ as:

$$\mathbf{g}_t = \sigma(\mathbf{W}_x\mathbf{y}_t^2 + \mathbf{W}_h\mathbf{h}_{t-1}^2) \tag{24}$$

$$\mathbf{s}_t = \mathbf{g}_t \odot \tanh(\mathbf{m}_t^2) \tag{25}$$

where $\mathbf{W}_x, \mathbf{W}_h$ are parameters to be learned, $\mathbf{y}_t^2$ is the input of the LSTM unit. $\odot$ is an element-wise product and $\sigma$ represents the logistic sigmoid activation. More specifically,

$$\mathbf{y}_t^2 = [\mathbf{v}_g, \tilde{\mathbf{h}}_t^1, \hat{\mathbf{v}}_t^1] \tag{26}$$

After we get $\mathbf{h}_t^2 = LSTM(\mathbf{y}_t^2, \mathbf{h}_{t-1}^2)$ and $\mathbf{s}_t$ from the second LSTM, we compute an attention vector to determine when and where to look at the visual and context information. To compute the attention vector, we firstly obtain $\mathbf{e}_t^2$ by:

$$\mathbf{e}_t^2 = \mathbf{w}_{z2}^T \tanh(\mathbf{W}_{v2}V, \mathbf{W}_{h2}\mathbf{h}_t^2) \tag{27}$$

where $\mathbf{w}_{z2}^T, \mathbf{W}_{v2}, \mathbf{W}_{h2}$ are parameters to be learned. $\mathbf{h}_t^2$ is the LSTM output at the $t$-th time step. Next, we use the following function to calculate the attention weights $\alpha_t^2$.

$$\alpha_t^2 = soft\max([\mathbf{e}_t^2; \mathbf{w}_a^T \tanh(\mathbf{W}_s\mathbf{s}_t + \mathbf{W}_{h3}\mathbf{h}_t^2)]) \tag{28}$$

where $\mathbf{w}_a^T, \mathbf{W}_s, \mathbf{W}_{h3}$ are parameters to be learned. $\alpha_t^2 \in \mathbb{R}^{L+1}$ contains weights for image region features and the sequential context information $\mathbf{s}_t$. Finally, the attended results can be obtained by:

$$\hat{\mathbf{v}}_t^2 = \sum_{i=1}^{L+1} \alpha_{i,t}^2 \mathbf{v}_i \tag{29}$$

where $\mathbf{v}_{L+1}$ equals to $\mathbf{s}_t$.

To generate the $t$-th word, we combine the output of the first layer $\tilde{\mathbf{h}}_t^1$, the output of the second LSTM $\mathbf{h}_t^2$ and the attended visual features $\mathbf{v}_t^2$ together by the function below:

$$\tilde{\mathbf{h}}_t^2 = \mathbf{W}_{sd}[\tilde{\mathbf{h}}_t^1; \mathbf{h}_t^2; \hat{\mathbf{h}}_t^2] \tag{30}$$

where $\mathbf{W}_{sd}$ is the parameter to be learned. We use softmax function to calculate the probability of the $t$-th word:

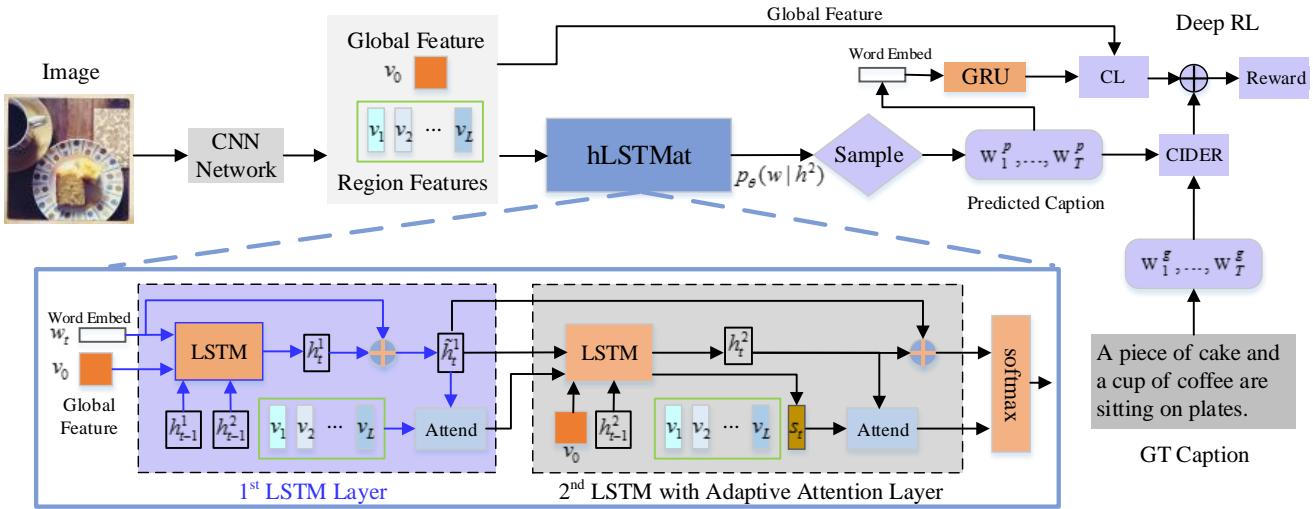$$\mathbf{p}_t = softmax(\tilde{\mathbf{h}}_t^2) \tag{31}$$

Fig. 7. The proposed framework for DA, our model uses residual shortcut connection to improve information flow through two LSTMs. And adaptive attention is applied to calculate weights of features when predicting new word.
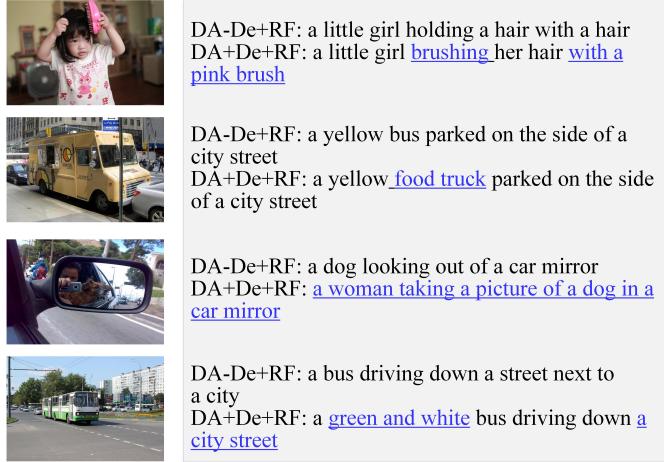


Fig. 8. Examples of image captions generated by our DA. For each image, the first and second sentences are generated by the first-pass (DA-De+RF) and the second-pass respectively (DA+De+RF). This demonstrates that our deliberate attention can generate more precise concepts (e.g., "food truck" vs "bus") and more reasonable descriptions (e.g., "a women taking a picture of a dog in a car mirror" vs "a dog looking out of a car mirror") by using a deliberate process.

## 6.3 Two-step Training

In essential, our training process can be divided into two stages. Firstly, the parameters $\Theta$ in image captioning model are pre-trained by minimizing MLE loss (Eq. 3), and then we fine-tune the model with reinforcement learning.

### 6.3.1 Training with Reinforcement Learning

Inspired by the previous work [61], [62], we consider our model introduced above as "agent" to interact with external environment (i.e., words, global and region visual features), and $(\Theta)$ as the policy to conduct an action to predict a word. After the whole caption is generated, the agent observers a reward. Since CIDEr is proposed to evaluate the quality of image captioning model.

We design our reward functions by combing contrastive loss (CL) with CIDEr. Next, we introduce the definition of CL.

**Contrastive Loss (CL).** Given an image $\mathbf{x}$ and caption $\mathbf{c}$, we obtain caption and image features by RNN and CNN, respectively. We take global image feature $\mathbf{v}_0$ as image features. Each word in $\mathbf{c}$ is embedded and then input into a RNN network to derive a caption feature. We define caption feature as $\mathbf{c}_0$. Next, we map two features into a common space by $\mathbf{W}_v^T$ and $\mathbf{W}_c^T$, respectively:

$$f(\mathbf{x}) = \mathbf{W}_v^T \mathbf{v}_0 \tag{32}$$

$$f(\mathbf{c}) = \mathbf{W}_c^T \mathbf{c}_0 \tag{33}$$

Furthermore, cosine similarity is used to compute similarity between an image and caption:

$$s(\mathbf{x}, \mathbf{c}) = \frac{f(\mathbf{x}) \cdot f(\mathbf{c})}{\|f(\mathbf{x})\| \, \|f(\mathbf{c})\|} \tag{34}$$

Parameters in the such model are learned by minimizing contrastive loss (CL), which is a sum of two hinge losses:

$$\ell_{CON}(\mathbf{x}, \mathbf{c}) = \max_{\mathbf{c}'} \left[ \alpha + s(\mathbf{x}, \mathbf{c}') - s(\mathbf{x}, \mathbf{c}) \right]_+ \\ + \max_{\mathbf{x}'} \left[ \alpha + s(\mathbf{x}', \mathbf{c}) - s(\mathbf{x}, \mathbf{c}) \right]_+ \tag{35}$$

where $[\mathbf{x}]_+ \equiv \max(\mathbf{x}, 0)$. In Eq. 35, $(\mathbf{c}, \mathbf{x})$ indicates that a pair of caption and image is matched. $\mathbf{c}$ correctly describes image $I$. Both $(\mathbf{x}, \mathbf{c}')$ and $(\mathbf{x}', \mathbf{c})$ are mismatched. $((\mathbf{x}, \mathbf{c}'))$ suggests that $\mathbf{c}'$ is the incorrect description of $\mathbf{x}$, while $(\mathbf{x}', \mathbf{c})$ suggests that $\mathbf{c}$ is the incorrect description of $\mathbf{x}'$. $\alpha$ is used to ensure minimum gap between scores of $(\mathbf{c}, \mathbf{x})$ with $(\mathbf{x}, \mathbf{c}')$ and $(\mathbf{x}', \mathbf{c})$.

**CIDEr+CL** In order to optimize the parameters in our model, the objective becomes to maximize the reward obtained from reward function by learning parameters. An update is implemented by computing the gradient of the expected reward:

$$\nabla_\Theta E[R(\hat{\mathbf{c}}, \mathbf{x})] \approx (R(\hat{\mathbf{c}}, \mathbf{x}) - R(\mathbf{c}^*, \mathbf{x})) \nabla_\Theta \log p(\hat{\mathbf{c}}|\mathbf{x}; \Theta) \tag{36}$$

where $R(\hat{\mathbf{c}}, \mathbf{x}) = CIDEr(\hat{\mathbf{c}}) - \ell_{CON}(\hat{\mathbf{c}}, \mathbf{x})$ is reward function. $\hat{\mathbf{c}}$ is caption generated from (31). The baseline is computed by $\mathbf{c}^* = (BOS, \mathbf{w}_1^*, ..., \mathbf{w}_T^*)$, which is obtained as:

$$\mathbf{w}_t^p = \arg\max_{\mathbf{w}} (\mathbf{w}|\mathbf{w}_{0...t-1}^p, \mathbf{x}) \tag{37}$$

# 7 EXPERIMENTS FOR IMAGE CAPTIONING

## 7.1 Datasets

In this paper, we utilize two datasets including COCO [63] and Flickr30K [64] to evaluate the performance of our proposed DA network.

**COCO.** It is the largest dataset for image captioning, which consists of $82,783$, $40504$ and $40,775$ images for training, validation and testing, respectively. In COCO dataset, each image has 5 captions annotated by human beings. Following previous work [61], [65], [66], [67], we use the "Karpathy" splits proposed in [68]. In this split, $113,287$, $5,000$ and $5,000$ images are used for training, validation and testing, respectively.

**Flickr30K.** It consists of $31,783$ images collected from Flickr. In particular, each image is associated with 5 crowd-scoured descriptions. Most of the images are about human beings performing activities. Following previous work [66], we use 29k images for training, 1k for validation and 1k for testing.

For both COCO and Flickr30K dataset, we conduct a pre-processing procedure by firstly truncating captions longer than 16 words. Next, all modified captions are converted to lower case. For each dataset, we build a vocabulary. For COCO, the vocabulary contains $9,487$ words, while for Flickr30K the vocabulary has 7k words.

## 7.2 Evaluation Metric

Five generally used evaluation metrics are adopted to evaluate the performance of image captioning, including BLUE [69], ROUGEL [70], METEOR [22], CIDEr [71] and SPICE [72]. More specifically, for the COCO dataset, we also report the SPICE subclass scores on 5k validation sets, including Color, Attribute, Cardinality, Object, Relation and Size. All the SPICE subclass scores are scaled up by 100.

## 7.3 Implementation Details

To extract the global visual feature, we use pre-trained ResNet-101 [59] to process the input image and the output of "pool5" (2048-D) is used as global appearance feature. In terms of region visual features, we utilize the same region features used in [67]. As a result, for each image, we obtain 36 region features and the dimension for each region feature is also 2048-D. In addition, the hidden state size of two LSTM in our DA network is set to be 512.

To compute the contrastive loss, we use RNN as the text encoder to extract caption features and use pool5 of ResNet-101 to extract image features. The dimension of word embedding is 512, the RNN hidden state size is set as 1024 and the dimension of embedded image feature is 1024. To train the model for calculate contrastive loss, we set epochs as 27 for both COCO and Flicker30K.

For training process, MLE is firstly used to pre-train the DA model. Next, we train it with reinforcement learning by using CIDEr and CL as reward value. For MLE training, the epoch is set as 150 for both COCO and Flickr30K, while for reinforcement learning the epoch is set as 200 for COCO and 150 for Flickr30K. All models are trained by using Adam and the batch size is set as 128. We initialize the learning rate with 5e-4 and update it by a decreasing factor 0.8 in every 15 epochs. When conduct testing, beam search is applied to predict captions, with beam size setting as 5.

## 7.4 Ablation Study

In order to figure out the contribution of each component, we conduct the following ablation studies on the COCO dataset. Specifically, we remove the deliberation process (De) and reinforcement learning (RF) respectively from our DA model, and have four experiments: DA-De-RF, DA-De+RF, DA+De-RF and DA+De+RF. The experimental results are shown in Tab. 10.

**The Influence of Deliberation.** From Tab. 10, we can see that with or without reinforcement learning, our DA+De models, including DA+De+RF and DA+De-RF, perform better than DA-De models (DA-De-RF and DA-De+DF). More specifically, DA+De-RF outperforms DA-De-RF on all 12 metrics, in particular with an increase of 2% on BLUE4, 1% on METEOR, 1.6% ROUGE, 7% CIDEr and 1.1% SPICE. In addition, compared with DA-De+RF, DA+De+RF obtains higher performance on all 12 metrics. This verifies the importance of our deliberation mechanism.

In order to further demonstrate the role of deliberation mechanism, we show four visual examples in Fig. 8. In Fig. 8, each image contains two descriptions. The first sentence is generated by the DA-De+RF and the second sentence is generated by the DA+De+RF model. We can see that without deliberation process, the model can generate a sentence which contains error information (e.g., "a hair with a hair") or inconsistent semantic information, e.g., "a dog looking out of a car mirror". With the deliberation component, our DA model can provide more precise description, such as "a yellow food truck" instead of "a yellow bus"; and a more reasonable description: "a women taking a picture of a dog in a car mirror" instead of "a dog looking out of a car mirror". These examples also show that the first residual attention based layer can detect primary objects (e.g., girl, hair) and activities (e.g., holding), while the second residual attention based layer can refine the activities (e.g., brushing) and detect the relationship between objects (e.g., "a women taking a picture of a dog" and "a dog in a car mirror").

In addition, Fig. 9 shows the attended image regions of the first residual based attention of the DA+De+RF. For each generated word, we visualize the attention weights on individual pixels, outlining the region with the maximum attention weight in orange. Moreover, for each word, we display the corresponding visual weight of second residual based attention layer. From Fig. 9, we find that our first layer attention is able to locate the right objects, which enables the DA+De+RF to accurately describe objects occurred in the input image. On the other hand, the visual weights in the second layer are obviously higher when our model predicts words related to objects (e.g., car and pizza).

**The Influence of Reinforcement Learning.** From Tab. 10, we can see that the results clearly show the advantage of our reinforcement learning. From the first block of Tab. 10, we can see that DA-De+RF achieves 4.2% BLUE1, 1.5% BLUE4, 0.9% METEOR, 1.5% ROUGE, 12.4% CIDEr and 1.7% SPICE performance gain compared with DA-De-DF. Moreover, compared with DA+De-RF, DA+De+DF performs better with an increase of 4.1% BLUE1, 1.8% BLUE4, 1.1% METEOR, 2.0% ROUGE, 13.7% CIDEr and 1.8% SPICE. The increases of performance clearly demonstrate the effectiveness of the proposed reinforcement learning component.

In order to further demonstrate the discriminability, we show some qualitative results in Fig. 10. By observing the two images, a human can aware that two men are performing the same activity (i.e., doing a tricks on a skateboard), but they are playing at

TABLE 10
Ablation study results obtained from the COCO dataset.

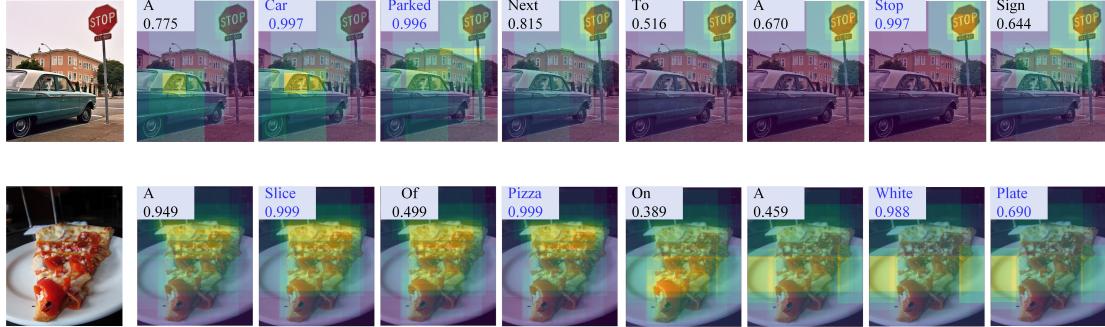| model | BLUE1 | BLUE4 | METEOR | ROUGE | CIDEr | SPICE | Att | Card | Col | Obj | Rel | Size |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| DA-De-RF | 74.2 | 33.7 | 26.4 | 54.6 | 104.9 | 19.4 | 9.3 | 2.1 | 11.0 | 35.5 | 5.2 | 3.8 |
| DA-De+RF | 78.4 | 35.2 | 27.3 | 56.5 | 117.3 | 21.1 | 9.2 | 13.0 | 10.5 | 39.2 | 5.6 | 2.8 |
| DA+De-RF | 75.8 | 35.7 | 27.4 | 56.2 | 111.9 | 20.5 | 10.8 | 6.1 | 14.6 | 36.8 | 5.5 | **5.6** |
| DA+De+RF | **79.9** | **37.5** | **28.5** | **58.2** | **125.6** | **22.3** | **11.2** | **14.4** | **15.2** | **40.3** | **6.4** | 3.7 |



Fig. 9. Visualization of first residual attention map of the DA+De+RF model. The sentence is generated by the DA+De+RF. The region with the maximum attention weight is in orange. We also show each word with the corresponding visual weight in the second residual attention block.

TABLE 11
Performance on Flicker30k test split. DA refers to the DA+De+RF mode in Tab. 10.

| model | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 | METEOR | CIDEr |
|---|---|---|---|---|---|---|
| DeepVS [73] | 57.3 | 36.9 | 24.0 | 15.7 | 15.3 | 24.7 |
| Hard-Attention [74] | 66.9 | 43.9 | 29.6 | 19.9 | 18.5 | - |
| ATT-FCN [65] | 64.7 | 46.0 | 32.4 | 23.0 | 18.9 | - |
| Adaptive-Attention [66] | 67.7 | 49.4 | 35.4 | 25.1 | 20.4 | 53.1 |
| DA | **73.8** | **55.1** | **40.3** | **29.4** | **23.0** | **66.6** |
| Relative Improvement | 6.1 | 5.7 | 4.9 | 4.3 | 2.6 | 13.5 |

TABLE 12
Single-model image captioning performance on the COCO Karpathy test split. B-4, M, R, C and S are BLUE4, METEOR, ROUGE, CIDEr and SPICE scores, respectively. All methods are based on the reinforcement learning.

| model | B4 | MR | R | C | S | Att | Card | Col | Obj | Rel | Size |
|---|---|---|---|---|---|---|---|---|---|---|---|
| SCST:Att2in [62] | 31.3 | 26.0 | 54.3 | 101.3 | - | - | - | - | - | - | - |
| SCST:Att2all [62] | 30.0 | 25.9 | 53.4 | 99.4 | - | - | - | - | - | - | - |
| ATTN+C+D(1) [61] | 36.3 | 27.3 | 57.1 | 114.1 | 21.1 | 9.5 | 10.5 | 9.3 | 39.0 | 5.9 | 2.6 |
| Up-Down [67] | 36.3 | 27.7 | 56.9 | 120.1 | 21.4 | 10.0 | **18.4** | 11.4 | 39.1 | **6.5** | 3.2 |
| DA | **37.5** | **28.5** | **58.2** | **125.6** | **22.3** | **11.2** | 14.4 | **15.2** | **40.3** | 6.4 | **3.7** |



DA+De-RF : a man doing a trick on a skateboard in the air
DA+De+RF : a man doing a trick on a skateboard on a street



DA+De-RF : a man doing a trick on a skateboard in the air
DA+De+RF : a man doing a trick on a skateboard on a park

Fig. 10. The role of RF. Examples of image captions generated by DA+De-RF and DA+De+RF.

different places (i.e., on a street and on a park). The first man plays skateboard on the street, while the second man plays skateboard

on a park. From Fig. 10, we can see that DA+De-RF generates the same description for both images, while DA+De+RF is able to generate precise caption to describe their difference (i.e., on a street or on a park). This indicates the reinforcement learning with contrastive loss improves the discriminability of our DA model.

## 7.5 Comparing with State-of-the-Art

In this section, we use DA to represent our model DA+De+RF for convenience.

**Flickr30K.** For Flickr30K dataset, we compare our DA with DeepVS [73], Hard-Attention [74], ATT-FCN [65] and Adaptive-Attention [66] and the comparison results are shown in Tab. 11. From Tab. 11, we can see that our DA outperforms the counterparts by a large margin. Specifically, compared with Adaptive [66], DA has 6.1% BLUE-1, 5.7% BLUE-2, 4.9% BLUE-3, 4.3% BLUE-4, 2.6% METEOR and 13.5% CIDEr increases. The improvement is significant, especially for BLUE-n and CIDEr.

TABLE 13
Results on the online MSCOCO test sever. DA is a single model, both SCST:Att2all and Up-Down are an ensemble of 4 models. LSTM-A3 utilizes Resnet-152 based visual feature while DA uses RestNet101 based visual feature.

| | BLEU-1 | | BLEU-2 | | BLUE-3 | | BLUE-4 | | METEOR | | ROUGE-L | | CIDEr | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | c5 | c40 | c5 | c40 | c5 | c40 | c5 | c40 | c5 | c40 | c5 | c40 | c5 | c40 |
| Review Net | 72.0 | 90.0 | 55.0 | 81.2 | 41.4 | 70.5 | 31.1 | 59.7 | 25.6 | 34.7 | 53.5 | 68.6 | 96.5 | 96.9 |
| Adaptive | 74.8 | 92.0 | 58.4 | 84.5 | 44.4 | 74.4 | 33.6 | 63.7 | 26.4 | 35.9 | 55.5 | 70.5 | 104.2 | 105.9 |
| PG-BCMR | 75.4 | - | 59.1 | - | 44.5 | - | 33.2 | - | 25.7 | - | 55 | - | 101.3 | - |
| SCST:Att2all | 78.1 | 93.7 | 61.9 | 86.0 | 47.0 | 75.9 | 35.2 | 64.5 | 27.0 | 35.5 | 56.3 | 70.7 | 114.7 | 116.7 |
| LSTM-A3 | 78.7 | 93.7 | 62.7 | 86.7 | 47.6 | 76.5 | 35.6 | 65.2 | 27.0 | 35.4 | 56.4 | 70.5 | 116.0 | 118.0 |
| Up-Down | **80.2** | **95.2** | **64.1** | **88.8** | **49.1** | **79.4** | **36.9** | **68.5** | 27.6 | 36.7 | 57.1 | **72.4** | 117.9 | **120.5** |
| DA | 79.4 | 94.4 | 63.5 | 88.0 | 48.7 | 78.4 | 36.8 | 67.4 | **28.2** | **37.0** | **57.7** | 72.2 | **120.5** | 122.0 |



Fig. 11. Examples of image captions generated by Up-Down, DA and human beings. The first column, both Up-Down and our DA provide accurate descriptions. The middle column shows that in some cases our DA is able to provide better descriptions, while the third column indicates that in complex situations both Up-Down and DA fail.

**COCO.** We conduct two types of evaluations on the COCO dataset. The first is conducted offline by using the "Karpathy" split that have been widely used in prior work. The second one is conducted online and the captioning results are obtained on the online MSCOCO test server.

For offline evaluation, all the image captioning models are single-model. Here, we compare DA with SCST:Att2in [62], SCST:Att2all [62], ATTN+C+D(1) [61] and Up-Down [67]. The offline evaluation results are shown in Tab. 12. It is clear that our DA performs the best on the widely used evaluation metrics, e.g., BLUE4, METEOR, ROUGE, CIDEr and SPICE scores. Up-Down [67] is a strong competitor, and it performs the best for "Card" and "Rel".

For online evaluation, we compare with previous published works, including Review Net [75], Adaptive [66], PG-BCMR [76], SCST:Att2all [62], LSTM-A3 [77], Up-Down [67]. The experimental results are shown in Tab. 13. From Tab. 13, we can see that beside METEOR, Up-down in general performs the best. In fact, both Up-Down and SCST:Att2All are an ensemble of 4 models, while our DA uses a single-model. Although LSTM-A3 utilizes better visual features extracted from the ResNet-152, our DA with ResNet-101 visual features obtains higher performances, especially the METEOR scores reaching 28.2% on c5 and 37.0% on c40. We believe that the performance of DA could be further boosted via an ensemble of multiple DA-based models.

### 7.6 Qualitative Analysis

Here, we show some qualitative results in Fig. 11. From the above Tables (i.e., Tab. 11, Tab. 12 and Tab. 13), we can see that the previously Adaptive [66] performs the best on the Flicker30k while Up-Down [67] performs the best on the COCO dataset. Due

to the reason that Up-Down [67] releases the code while Adaptive [66] does not, we show some captioning examples obtained by our DA and Up-Down. The first column with two examples show that both DA and Up-Down are able to provide accurate description. For the middle column, we can see that our DA provides more accurate descriptions, especially for describing the relationships among objects. In this case, Up-Down fails to detect all objects and the relationships among them. The third column shows two images and both of them have complex background and their corresponding descriptions contain rich semantic information. For those two images, both DA and Up-Down fails. One possible reason is that a human being can conduct reasoning based on his or her background knowledge while at this stage both DA and Up-Down cannot. This proposes a potential research direction for image captioning.

## 8 CONCLUSION AND FUTURE WORK

In this paper, we firstly introduce a novel hLSTMat encoder-decoder framework, which integrates a hierarchical LSTMs, temporal attention and adaptive temporal attention to automatically decide when to make good use of visual information or when to utilize sentence context information, as well as to simultaneously considering both low-level video visual features and language context information. Experiments show that hLSTMat achieves state-of-the-art performances on both MSVD and MSR-VTT datasets. Secondly, we propose two spatial-temporal networks, namely ParA and Two-stream, and they further improve the performance of video captioning on all three datasets.

## ACKNOWLEDGMENTS

## REFERENCES

[1] J. Song, L. Gao, F. Nie, H. T. Shen, Y. Yan, and N. Sebe, "Optimized graph learning using partial tags and multiple features for image and video annotation," *IEEE Transactions on Image Processing*, vol. 25, no. 11, pp. 4999–5011, 2016.

[2] L. Gao, P. Wang, J. Song, Z. Huang, J. Shao, and H. T. Shen, "Event video mashup: From hundreds of videos to minutes of skeleton," in *AAAI*, 2017, pp. 1323–1330.

[3] J. Wang, T. Zhang, J. Song, N. Sebe, H. T. Shen *et al.*, "A survey on learning to hash," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.

[4] J. Song, L. Gao, L. Liu, X. Zhu, and N. Sebe, "Quantization-based hashing: a general framework for scalable image and video retrieval," *Pattern Recognition*, 2017.

[5] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and tell: A neural image caption generator," in *CVPR*, 2015, pp. 3156–3164.

[6] K. Xu, J. Ba, R. Kiros, A. Courville, R. Salakhutdinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," *arXiv preprint arXiv:1502.03044*, 2015.

[7] J. Lu, C. Xiong, D. Parikh, and R. Socher, "Knowing when to look: Adaptive attention via a visual sentinel for image captioning," *arXiv preprint arXiv:1612.01887*, 2016.

[8] A. Karpathy, A. Joulin, and F. F. F. Li, "Deep fragment embeddings for bidirectional image sentence mapping," in *NIPS*, 2014, pp. 1889–1897.

[9] H. Fang, S. Gupta, F. Iandola, R. K. Srivastava, L. Deng, P. Dollár, J. Gao, X. He, M. Mitchell, J. C. Platt *et al.*, "From captions to visual concepts and back," in *CVPR*, 2015, pp. 1473–1482.

[10] X. Chen and C. L. Zitnick, "Learning a recurrent visual representation for image caption generation," *arXiv preprint arXiv:1411.5654*, 2014.

[11] L. Chen, H. Zhang, J. Xiao, L. Nie, J. Shao, and T.-S. Chua, "Sca-cnn: Spatial and channel-wise attention in convolutional networks for image captioning," *arXiv preprint arXiv:1611.05594*, 2016.

[12] S. Venugopalan, H. Xu, J. Donahue, M. Rohrbach, R. Mooney, and K. Saenko, "Translating videos to natural language using deep recurrent neural networks," *arXiv preprint arXiv:1412.4729*, 2014.

[13] S. Venugopalan, M. Rohrbach, J. Donahue, R. Mooney, T. Darrell, and K. Saenko, "Sequence to sequence-video to text," in *ICCV*, 2015, pp. 4534–4542.

[14] L. Yao, A. Torabi, K. Cho, N. Ballas, C. Pal, H. Larochelle, and A. Courville, "Describing videos by exploiting temporal structure," in *ICCV*, 2015.

[15] G. Li, S. Ma, and Y. Han, "Summarization-based video caption via deep neural networks," in *ACM Multimedia*, 2015, pp. 1191–1194.

[16] C. Gan, C. Sun, L. Duan, and B. Gong, "Webly-supervised video recognition by mutually voting for relevant web images and web video frames," in *ECCV*, 2016.

[17] K. Cho, A. Courville, and Y. Bengio, "Describing multimedia content using attention-based encoder-decoder networks," *Multimedia, IEEE Transactions on*, vol. 17, no. 11, pp. 1875–1886, 2015.

[18] J. Gu, Z. Lu, H. Li, and V. O. K. Li, "Incorporating copying mechanism in sequence-to-sequence learning," in *ACL*, 2016.

[19] M.-T. Luong, H. Pham, and C. D. Manning, "Effective approaches to attention-based neural machine translation," *arXiv preprint arXiv:1508.04025*, 2015.

[20] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola, "Stacked attention networks for image question answering," in *CVPR*, 2016.

[21] Z. Gan, C. Gan, X. He, Y. Pu, K. Tran, J. Gao, L. Carin, and L. Deng, "Semantic compositional networks for visual captioning," *CoRR*, vol. abs/1611.08002, 2016.

[22] M. Denkowski and A. Lavie, "Meteor universal: Language specific translation evaluation for any target language," in *Proceedings of the ninth workshop on statistical machine translation*, 2014, pp. 376–380.

[23] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma *et al.*, "Visual genome: Connecting language and vision using crowdsourced dense image annotations," *International Journal of Computer Vision*, vol. 123, pp. 32–73, 2017.

[24] Y. Xia, F. Tian, L. Wu, J. Lin, T. Qin, N. Yu, and T. Liu, "Deliberation networks: Sequence generation beyond one-pass decoding," in *NIPS*, 2017, pp. 1782–1792.

[25] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NIPS*, 2012, pp. 1097–1105.

[26] 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, pp. 1–9.

[27] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[28] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.

[29] Y. Tian, B. Fan, and F. Wu, "L2-net: Deep learning of discriminative patch descriptor in euclidean space," in *CVPR*, July 2017.

[30] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3d convolutional networks," in *ICCV*. IEEE, 2015, pp. 4489–4497.

[31] R. Xu, C. Xiong, W. Chen, and J. J. Corso, "Jointly modeling deep video and compositional text to bridge vision and language in a unified framework." in *AAAI*. Citeseer, 2015, pp. 2346–2352.

[32] P. Pan, Z. Xu, Y. Yang, F. Wu, and Y. Zhuang, "Hierarchical recurrent neural encoder for video representation with application to captioning," in *CVPR*, 2016, pp. 1029–1038.

[33] H. Yu, J. Wang, Z. Huang, Y. Yang, and W. Xu, "Video paragraph captioning using hierarchical recurrent neural networks," in *CVPR*, 2016.

[34] X. Long, C. Gan, and G. de Melo, "Video captioning with multi-faceted attention," *arXiv preprint arXiv:1612.00234*, 2016.

[35] C. Hori, T. Hori, T. Lee, Z. Zhang, B. Harsham, J. R. Hershey, T. K. Marks, and K. Sumi, "Attention-based multimodal fusion for video description," in *ICCV*, 2017, pp. 4203–4212.

[36] V. Ramanishka, A. Das, D. H. Park, S. Venugopalan, L. A. Hendricks, M. Rohrbach, and K. Saenko, "Multimodal video description," in *ACM Multimedia*, 2016, pp. 1092–1096.

[37] Q. You, H. Jin, Z. Wang, C. Fang, and J. Luo, "Image captioning with semantic attention," in *CVPR*, 2016.

[38] T. Yao, Y. Pan, Y. Li, Z. Qiu, and T. Mei, "Boosting image captioning with attributes," *arXiv preprint arXiv:1611.01646*, 2016.

[39] Y. Pan, T. Yao, H. Li, and T. Mei, "Video captioning with transferred semantic attributes," *arXiv preprint arXiv:1611.07675*, 2016.

[40] M. Ranzato, S. Chopra, M. Auli, and W. Zaremba, "Sequence level training with recurrent neural networks," *arXiv preprint arXiv:1511.06732*, 2015.

[41] S. J. Rennie, E. Marcheret, Y. Mroueh, J. Ross, and V. Goel, "Self-critical sequence training for image captioning," *arXiv preprint arXiv:1612.00563*, 2016.

[42] S. Liu, Z. Zhu, N. Ye, S. Guadarrama, and K. Murphy, "Improved image captioning via policy gradient optimization of spider," *arXiv preprint arXiv:1612.00370*, 2016.

[43] Z. Ren, X. Wang, N. Zhang, X. Lv, and L.-J. Li, "Deep reinforcement learning-based image captioning with embedding reward," *arXiv preprint arXiv:1704.03899*, 2017.

[44] Y. Bengio, P. Simard, and P. Frasconi, "Learning long-term dependencies with gradient descent is difficult," *IEEE transactions on neural networks*, vol. 5, no. 2, pp. 157–166, 1994.

[45] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *NIPS*, 2015, pp. 91–99.

[46] M. Hermans and B. Schrauwen, "Training and analysing deep recurrent neural networks," in *NIPS*, 2013, pp. 190–198.

[47] K. Simonyan and A. Zisserman, "Two-stream convolutional networks for action recognition in videos," in *NIPS*, 2014, pp. 568–576.

[48] D. L. Chen and W. B. Dolan, "Collecting highly parallel data for paraphrase evaluation," in *ACL*, 2011, pp. 190–200.

[49] J. Xu, T. Mei, T. Yao, and Y. Rui, "Msr-vtt: A large video description dataset for bridging video and language," in *CVPR*, 2016.

[50] A. Rohrbach, A. Torabi, M. Rohrbach, N. Tandon, C. Pal, H. Larochelle, A. Courville, and B. Schiele, "Movie description," *International Journal of Computer Vision*, vol. 123, no. 1, pp. 94–120, 2017.

[51] A. Rohrbach, M. Rohrbach, N. Tandon, and B. Schiele, "A dataset for movie description," in *CVPR*, 2015, pp. 3202–3212.

[52] A. Torabi, C. Pal, H. Larochelle, and A. Courville, "Using descriptive video services to create a large data source for video annotation research," *arXiv preprint arXiv:1503.01070*, 2015.

[53] M. D. Zeiler, "Adadelta: an adaptive learning rate method," *arXiv preprint arXiv:1212.5701*, 2012.

[54] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *ACL*, 2002, pp. 311–318.

[55] S. Banerjee and A. Lavie, "Meteor: An automatic metric for mt evaluation with improved correlation with human judgments," in *Proceedings*

*of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, 2005, pp. 65–72.

[56] R. Vedantam, C. Lawrence Zitnick, and D. Parikh, "Cider: Consensus-based image description evaluation," in *CVPR*, 2015, pp. 4566–4575.

[57] Y. Pan, T. Mei, T. Yao, H. Li, and Y. Rui, "Jointly modeling embedding and translation to bridge video and language," in *CVPR*, 2016, pp. 4594–4602.

[58] Y. Yu, H. Ko, J. Choi, and G. Kim, "End-to-end concept word detection for video captioning, retrieval, and question answering," in *CVPR*, 2017, pp. 3261–3269.

[59] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.

[60] S. Ren, K. He, R. B. Girshick, and J. Sun, "Faster R-CNN: towards real-time object detection with region proposal networks," in *NIPS*, 2015, pp. 91–99.

[61] R. Luo, B. L. Price, S. Cohen, and G. Shakhnarovich, "Discriminability objective for training descriptive captions," vol. abs/1803.04376, 2018.

[62] S. J. Rennie, E. Marcheret, Y. Mroueh, J. Ross, and V. Goel, "Self-critical sequence training for image captioning," in *CVPR*, 2017, pp. 1179–1195.

[63] T. Y. Lin, M. Maire, S. J. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollar, and C. L. Zitnick, "Microsoft common objects in context," in *ECCV*, 2014, pp. 740–755.

[64] P. Young, A. Lai, M. Hodosh, and J. Hockenmaier, "From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions," *Transactions of the Association for Computational Linguistics*, vol. 2, pp. 67–78, 2014.

[65] Q. You, H. Jin, Z. Wang, C. Fang, and J. Luo, "Image captioning with semantic attention," in *CVPR*, 2016, pp. 4651–4659.

[66] J. Lu, C. Xiong, D. Parikh, and R. Socher, "Knowing when to look: Adaptive attention via a visual sentinel for image captioning," in *CVPR*, 2017, pp. 3242–3250.

[67] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang, "Bottom-up and top-down attention for image captioning and visual question answering," in *CVPR*, 2018.

[68] A. Karpathy and F. Li, "Deep visual-semantic alignments for generating image descriptions," in *CVPR*, 2015, pp. 3128–3137.

[69] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *ACL*, 2002, pp. 311–318.

[70] C.-Y. Lin, "Rouge: A package for automatic evaluation of summaries," *Text Summarization Branches Out*, 2004.

[71] R. Vedantam, C. Lawrence Zitnick, and D. Parikh, "Cider: Consensus-based image description evaluation," in *CVPR*, 2015, pp. 4566–4575.

[72] P. Anderson, B. Fernando, M. Johnson, and S. Gould, "Spice: Semantic propositional image caption evaluation," in *ECCV*.   Springer, 2016, pp. 382–398.

[73] A. Karpathy and F. Li, "Deep visual-semantic alignments for generating image descriptions," in *CVPR*, 2015, pp. 3128–3137.

[74] 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, pp. 2048–2057.

[75] Z. Yang, Y. Yuan, Y. Wu, W. W. Cohen, and R. R. Salakhutdinov, "Review networks for caption generation," in *NIPS*, 2016, pp. 2361–2369.

[76] S. Liu, Z. Zhu, N. Ye, S. Guadarrama, and K. Murphy, "Improved image captioning via policy gradient optimization of spider," in *ICCV*, 2017, pp. 873–881.

[77] T. Yao, Y. Pan, Y. Li, Z. Qiu, and T. Mei, "Boosting image captioning with attributes," in *ICCV*, 2017, pp. 4904–4912.