

Hierarchical Recurrent Neural Network for Video Summarization

Bin Zhao

School of Computer Science and
Center for OPTical IMagery Analysis
and Learning (OPTIMAL),
Northwestern Polytechnical
University
Xi'an, Shaanxi, P. R. China 710072
binzhao111@gmail.com

Xuelong Li

Xi'an Institute of Optics and Precision
Mechanics, Chinese Academy of
Sciences
Xi'an, Shaanxi, P. R. China 710019
xuelong_li@opt.ac.cn

Xiaoqiang Lu

Xi'an Institute of Optics and Precision
Mechanics, Chinese Academy of
Sciences
Xi'an, Shaanxi, P. R. China 710019
luxq666666@gmail.com

ABSTRACT

Exploiting the temporal dependency among video frames or subshots is very important for the task of video summarization. Practically, RNN is good at temporal dependency modeling, and has achieved overwhelming performance in many video-based tasks, such as video captioning and classification. However, RNN is not capable enough to handle the video summarization task, since traditional RNNs, including LSTM, can only deal with short videos, while the videos in the summarization task are usually in longer duration. To address this problem, we propose a hierarchical recurrent neural network for video summarization, called H-RNN in this paper. Specifically, it has two layers, where the first layer is utilized to encode short video subshots cut from the original video, and the final hidden state of each subshot is input to the second layer for calculating its confidence to be a key subshot. Compared to traditional RNNs, H-RNN is more suitable to video summarization, since it can exploit long temporal dependency among frames, meanwhile, the computation operations are significantly lessened. The results on two popular datasets, including the Combined dataset and VTW dataset, have demonstrated that the proposed H-RNN outperforms the state-of-the-arts.

KEYWORDS

deep learning, video summarization, hierarchical recurrent neural network

1 INTRODUCTION

Nowadays, video data are increasing explosively with the popularity of camera devices. There is a great demand for automatic techniques to handle these videos efficiently. Particularly, video summarization is one of the techniques that provide a viewer-friendly way to browse the huge amount of video data [36]. In general, it generates the video summary by shortening the video content into a compact

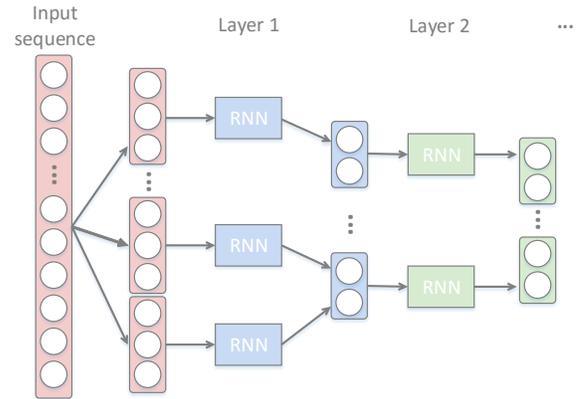


Figure 1: A compact architecture of hierarchical recurrent neural network. It is able to capture the long-range temporal dependency by processing the long input sequence hierarchically.

version [32]. Practically, there are several ways for video summarization to shorten the video. In this paper, we focus on the most popular one, i.e., key subshot selection.

There is a stable growing interest in video summarization. Earlier works are mostly based on clustering or dictionary learning [1, 5, 6, 37], where the cluster centers or dictionary elements are taken as the most representative subshots in the video, and then selected into the summary as key subshots. Later, to emphasize the salient visual attributes in frames (people, objects, etc.), several property models are designed to capture the importance, representativeness and interestingness of the summary, and utilized to score the subshots [8, 9, 18, 24]. Naturally, those subshots with higher scores are selected into the summary. These approaches demonstrate promising results, and usually exceed the clustering and dictionary learning ones.

However, all these approaches are not capable enough to model the rich information in video content. Recently, inspired by the great success of deep learning, *Convolutional Neural Network* (CNN) and *Recurrent Neural Network* (RNN) are introduced to video summarization, where CNN is utilized to extract deep visual features and RNN is employed to predict the probability of one subshot to be selected into the summary [32, 36]. This architecture has achieved the state-of-the-art results in video summarization. Apart from the

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '17, Mountain View, CA, USA

© 2017 ACM. 978-1-4503-4906-2/17/10...\$15.00

DOI: 10.1145/3123266.3123328

great ability of CNN in visual feature extraction, this mainly benefits from the capability of RNN in exploiting the temporal dependency among frames [36].

Unfortunately, RNN only works well for short frame sequence [21, 30]. Even for LSTM, one kind of RNN that is the most excellent in long frame sequence modeling, the favorable video length is less than 80 frames [21]. While to video summarization, most of the videos contain thousands of frames. In this case, it is difficult for RNN to capture this long-range temporal dependency of videos. Thus, current approaches that apply RNN directly to video summarization may restrict the quality of video summary. To address this problem, we propose a hierarchical structure of RNN. As depicted in Figure 1, the hierarchical RNN is composed of multi-layers, and each layer is with one or more short RNNs, by which the long input sequence is processed hierarchically. Actually, the hierarchical RNN is a general architecture which varies according to specific tasks.

In this paper, a specialized hierarchical RNN is designed for the task of video summarization, called as H-RNN. Detailedly, it contains two layers. The first layer is a LSTM, which is utilized to process video subshots generated by cutting the video evenly, and the intra-subshot temporal dependency is encoded in the final hidden. Then, the final hidden of each subshot is input to the second layer. Specifically, the second layer is a bi-directional LSTM, which is composed of a forward and a backward LSTM. It is employed to exploit the inter-subshot temporal dependency and determine whether a certain subshot is valuable to be a key subshot.

Generally, compared to current RNN-based approaches in video summarization, H-RNN has the following advantages:

- 1) H-RNN can model long-range temporal dependency with a short time step. As a result, it reduces the information loss in frame sequence modeling meanwhile the computation operations are reduced significantly.
- 2) The hierarchical structure of H-RNN increases the nonlinear fitting ability of traditional RNN, which has been demonstrated extremely helpful for visual tasks [25, 27].
- 3) H-RNN exploits the intra-subshot (i.e., among frames in the subshot) and inter-subshot temporal dependency in the two layers, respectively. This hierarchical structure is more suitable for video data, since video temporal structure is intrinsically layered as frames and subshots [23].

2 RELATED WORKS

There have been a variety of video summarization approaches proposed in the literature. Generally, existing approaches can be classified into unsupervised ones and supervised ones.

Unsupervised approaches select key subshots according to manually designed criteria [16, 20, 22, 37], such as representative to the video content and diverse with each other, etc. Clustering is one of the most popular unsupervised summarization approaches [1, 5, 20]. Practically, with hand-crafted features, similar frames are grouped into the same cluster, and the cluster centers are taken as the most representative elements and selected into the summary. Earlier works apply clustering algorithms to video summarization directly [10, 38]. Later, more works integrate the domain knowledge of video data into clustering algorithms [5, 20]. As in [5], the frames are initially clustered in sequential order, as consecutive frames

are similar and more probably to be allocated to the same cluster. Other works construct more comprehensive models based on the idea of clustering [3, 22]. As in [22], the video is transformed into an undirected graph, and the summary is generated by partitioning this graph into clusters. More recently, a co-clustering approach is proposed to simultaneously summarize several videos with the same topic by their co-occurrence, i.e., similar subshots shared by these videos are selected into the summary [3].

Dictionary learning is another important unsupervised summarization approach [4, 6, 19, 37]. This kind of approaches seek to select a few key subshots to compose a compact dictionary so as to represent the video content. [6] supposes that the original video can be reconstructed by its summary sparsely. Based on this point, the summary is generated by sparse coding. Furthermore, the Locality-constrained Linear Coding (LLC) is introduced to [17] to preserve the local similarity of video subshots when reconstructing the original video. Besides, to improve the efficiency, [37] propose a quasi real-time dictionary learning approach to summarize the video, which updates the video summary on-the-fly by adding those elements that cannot be well reconstructed by current video summary.

Supervised approaches learn the hidden summarization patterns from human generated summaries, which have been drawing increasing attention and getting more promising results than unsupervised ones [7, 35]. In supervised approaches, property models are usually taken as the decision criteria to select key subshots [8, 14, 18, 24]. For example, [14] and [8] build importance model and interesting model to score the subshots, and those subshots with higher scores are selected into the summary. [18] designs a smoothness model to constrain that the selected subshots have a smooth story line. Moreover, [9] employs three property models, i.e., interestingness, representativeness, uniformity, to build a comprehensive score function. Some other works even utilize auxiliary information to summarize videos, such as web image priors [12], video titles [26], and video category labels [24], etc.

More recently, deep learning is introduced to video summarization, including CNN and RNN. [32] builds a deep rank model relying on two CNNs i.e., AlexNet [13] and C3D [29] stitched with two Multi-Layer Perceptron (MLP) behind their final pooling layers. Given frames or subshots as input, the deep rank model outputs a ranking score. Naturally, a higher score indicates higher probability of that frame or subshot to be selected into the summary. Besides, LSTM is employed in [31, 36] to model the video sequence and rank the video subshots, which has achieved the state-of-the-art results in video summarization. However, due to the weakness in long temporal dependency exploitation, the input sequence to the LSTM is generated by the mean pooling or uniform sampling of frame features, which causes inevitable information loss. Actually, the proposed approach in this paper is essentially developed to solve this problem.

3 OUR APPROACH

In this section, we first provide a brief review to RNN, especially LSTM, since it is the build block of the proposed approach. Then, we present the hierarchical structure of RNN and our video summarization specialized hierarchical RNN, i.e., H-RNN.

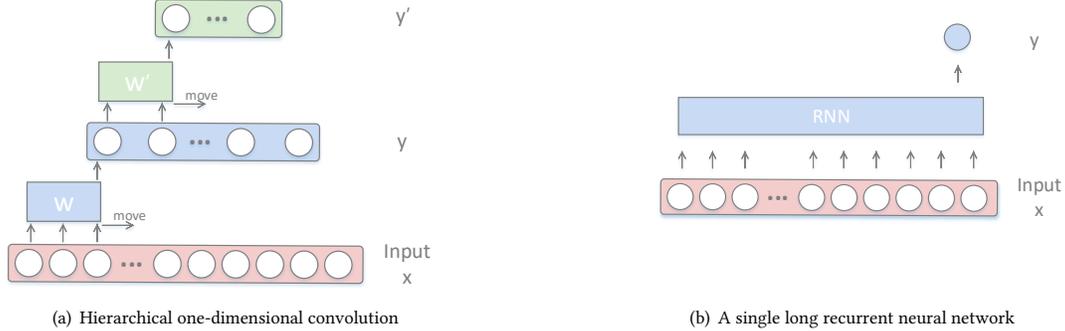


Figure 2: The structures of hierarchical convolution and RNN. Actually, it becomes a comparison between hierarchical RNN and single long RNN, if we replace the filters in (a) (i.e., w and w') with short RNNs. Compared to (b), hierarchical RNN is more efficient in long temporal dependency exploitation, meanwhile, the computation operations is significantly reduced.

3.1 Recurrent Neural Network

A standard RNN is constructed by extending a feedforward network with an extra feedback connection, so that it can model sequence. Practically, it can interpret the input sequence (x_1, x_2, \dots, x_n) into another sequence (y_1, y_2, \dots, y_n) iteratively by the following equations:

$$h_t = \phi(W_h x_t + U_h h_{t-1} + b_h), \quad (1)$$

$$y_t = \phi(U_y h_t + b_y), \quad (2)$$

where h_t is the hidden state, t denotes the t -th time step, ϕ stands for the activation function, and W , U and b are the training weights and biases.

Principally, the standard RNN should work efficiently in sequence modeling. However, it is really hard to train for the gradient vanishing problem [2]. Then, LSTM is designed to address this issue, which is the most popular variant of standard RNN [11]. Specifically, it is extended from standard RNN with an extra memory cell, which is utilized to selectively memorize the previous inputs. In fact, there are several variants of LSTM, and they are similar with each other. In this paper, the one proposed in [33] is employed, which is most widely used in video-based tasks. Detailedly, the calculation of hidden state h_t and memory cell c_t is formulated as:

$$i_t = \sigma(W_{ix} x_t + U_{ih} h_{t-1} + b_i), \quad (3)$$

$$f_t = \sigma(W_{fx} x_t + U_{fh} h_{t-1} + b_f), \quad (4)$$

$$o_t = \sigma(W_{ox} x_t + U_{oh} h_{t-1} + b_o), \quad (5)$$

$$g_t = \phi(W_{gx} x_t + U_{gh} h_{t-1} + b_g), \quad (6)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t, \quad (7)$$

$$h_t = o_t \odot \phi(c_t), \quad (8)$$

where σ denotes the sigmoid function, and all the W s, U s, b s are the training weights and bias. Besides, i_t , f_t and o_t are three gates, which are most important to LSTM. Concretely, the input gate i_t controls whether to record current input x_t , the forget gate f_t decides whether to drop previous memory cell c_{t-1} , and the

output gate o_t determines the information in current memory cell c_t transferred to the hidden state h_t .

3.2 Hierarchical Recurrent Neural Network

The motivation for designing hierarchical RNN is to improve its capability to exploit long-range temporal dependency of the videos. Actually, it is originally inspired by the operation of one-dimensional convolution. As depicted in the first layer of Figure 2(a), a one-dimensional filter w is utilized to exploit the sequential information by performing convolutional operations on the input sequence x :

$$y = w * x, \quad (9)$$

where y denotes the output sequence, and $*$ stands for the convolutional operation. It can be observed that although the filter w is much short than x , at each time step, it operates appropriately on a subsequence of x and outputs a much shorter sequence y . Particularly, if the convolution stride is set as n , $|y|$ is just $1/n$ of $|x|$, where $|\cdot|$ denotes the length of the sequence. Furthermore, in the second layer, another filter w' is applied to y , and outputs a shorter sequence y' . Naturally, more filters can also be applied to higher layers, until the final output is generated. Then, the hierarchical structure is formed, and the long sequence x is processed by several short filters hierarchically.

Inspired by this, we construct a similar hierarchical structure for RNN. Actually, the filters in different layers of Figure 2(a) are replaced by short RNNs, and the convolutional operation is just like successively processing several short subsequences which are cut from the long input sequence with or without overlap (the length of the subsequence is equal to the length of the filter RNN). Specifically, the RNN in the first layer exploits short-range temporal dependency, and longer ones are captured by higher layer RNNs. Intuitively, the long-range temporal dependency is captured by the hierarchical structure of several short RNNs. Moreover, compared to single RNN that operates on the long sequence directly in Figure 2(b), our hierarchical RNN can not only reduce the information loss in long sequence modeling, but also the computation operations

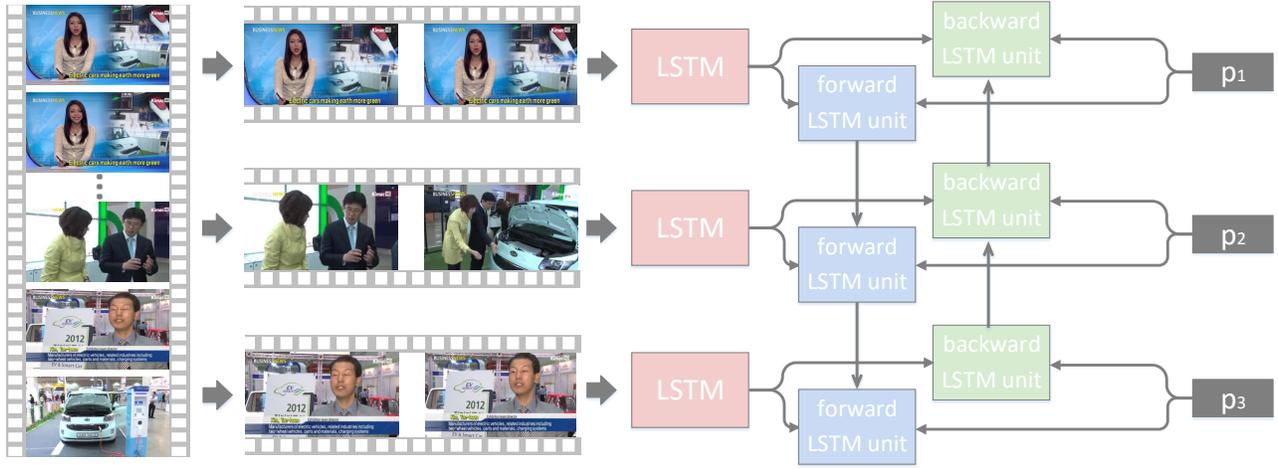


Figure 3: The architecture of the proposed approach for video summarization, i.e., H-RNN. It contains two layers, where the first layer is a LSTM and the second layer is a bi-directional LSTM (forward and backward). The two layers exploit the intra-subshot and inter-subshot temporal dependency, respectively, and the output of the second layer is utilized to predict the confidence of each subshot to be selected into the summary.

are significantly lessened, which is detailedly discussed for the task of video summarization in the next subsection.

Generally, the hierarchical RNN can be composed of multi-layers and each layer with several RNNs. In other words, it is a general architecture that varies according to specific tasks.

3.3 Video Summarization with H-RNN

The videos for summarization are usually with long durations, i.e., about thousands of frames. Moreover, according to [23], the video structure is obviously layered that frames form the subshots and subshots form the video. Therefore, the hierarchical RNN is quite applicable to the task of video summarization.

In this part, H-RNN is developed for the task of video summarization. As described in Figure 3, it contains two layers, where the first layer is a LSTM and the second layer is a bi-directional LSTM. The details about summarizing the video with H-RNN is presented as follows.

Firstly, the frame sequence (f_1, f_2, \dots, f_T) is separated into several subsequences, denoted as subshots $(f_1, f_2, \dots, f_s), (f_{s+1}, f_{s+2}, \dots, f_{2s}), \dots, (f_{m-s+1}, f_{m-s+2}, \dots, f_T)$, where f_i stands for the feature of frame i , T denotes the total frames in the video, m is the number of subshots, and s is the length of each subshots. Practically, if the final subshot is shorter than s , it is padded with zeros.

Then, the subshots are input to the first layer LSTM, which is formulated as follows:

$$\tau_i = LSTM(f_{i \cdot s + 1}, f_{i \cdot s + 2}, \dots, f_{(i+1) \cdot s}), \quad (10)$$

where $LSTM(\cdot)$ is short for Equ. (3)-(8), τ_i denotes the final hidden state of the i -th subshot. Actually, the short temporal dependency in the subshot is captured by τ_i . Thus, it is taken as the representation of the i -th subshot.

Next, the sequence $(\tau_1, \tau_2, \dots, \tau_m)$ is input to the second layer. As aforementioned, the second layer is a bi-directional LSTM. Actually, bi-directional LSTM is composed of a forward LSTM and a backward

LSTM. The main difference between them is that the backward LSTM operates reversely. Therefore, the calculation in the second layer is formulated as:

$$h_t^f = LSTM(\tau_t, h_{t-1}^f), \quad (11)$$

$$h_t^b = LSTM(\tau_t, h_{t+1}^b), \quad (12)$$

where h_t^f and h_t^b are the t -th output hidden state of forward LSTM and backward LSTM, respectively.

Finally, the output of the second layer is employed to predict the confidence of a certain subshot to be selected into the video summary. It is formulated as:

$$p_t = softmax\left(\tanh\left(W_p \left[h_t^f, h_t^b, \tau_t\right] + b_p\right)\right), \quad (13)$$

where W_p and b_p are the parameters to be learned. The softmax function is utilized to constrain the sum of the elements in p_t to be 1. Actually, p_t is a two-dimensional vector, each element of which indicates the possibility of the t -th subshot is key or non-key. It can be observed from Equ. (13) that p_t is determined jointly by the hidden state of forward and backward LSTM, i.e., h_t^f and h_t^b , together with the representation of the t -th subshot τ_t . It is because that for subshot t , h_t^f and h_t^b capture the front and behind temporal dependency, respectively, and τ_t contains the intra-subshot dependency. All these information are quite important for the determination of whether to select subshot t into the summary.

In this paper, the proposed H-RNN is trained end-to-end. Given the reference summaries generated manually, the parameters in H-RNN are learned by:

$$\Theta = \arg \min_{\Theta} \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^{m^{(i)}} L(p_t^{(i)}, g_t^{(i)}), \quad (14)$$

where N is the number of videos in the training set. $m^{(i)}$ denotes the number of shots in video i . $L(\cdot)$ stands for the loss function, which measures the cross-entropy between the generated probability distribution $p_t^{(i)}$ and the ground truth $g_t^{(i)}$. Practically, $g_t^{(i)}$ is a binary vector (indicates whether the subshot is key or not) or decimal vector (indicates the confidence of the subshot to be a key subshot).

Equ. (14) is optimized with *Stochastic Gradient Descent* method based on *Back Propagation Through Time* algorithm. Actually, H-RNN is easier to train than traditional LSTM, because the computation operations are much lessened [15]. For example, for a video with 2000 frames, traditional LSTM needs 2000 operations to handle the whole frame sequence. While in H-RNN, the video is handled if the length of the first layer LSTM is 40 and the second layer bi-directional LSTM is 50, i.e., just 140 computation operations are needed totally. That is to say, more than ninety percent computation operations are reduced by H-RNN.

4 EXPERIMENTS

To verify the effectiveness of the proposed approach, it is tested on two popular datasets, i.e., Combined and VTW, and compared with several state-of-the-art approaches on video summarization.

4.1 Setup

4.1.1 Dataset. The first dataset is combined with three popular datasets, i.e., SumMe [8], TVsum [26] and MED [24]. In this paper, it is called the Combined dataset. The intuition lying behind the combination is that the videos in the three datasets are similar both in their visual contents and styles. Moreover, the combination of these datasets can address the problem of lacking of training data, which is widely used in video summarization. Detailedly, the Combined dataset consists of 235 videos, the average duration is 2 minutes, about 3000 frames for each video. In this paper, the Combined dataset is split into a training set of 180 videos, and a testing set of 55 videos.

The second dataset, i.e., VTW, is originally proposed for the task of video captioning, which totally contains 18100 videos [34]. Fortunately, 2000 of them are labeled with subshot-level highlight scores that indicate the confidence of each subshot to be selected into the summary, so they are employed in this paper. Specifically, these videos are open-domain that crawled from YouTube, and the average duration is 1.5 minutes, about 2000 frames for each video. In this paper, the selected 2000 videos in the VTW dataset is divided into two parts, 1500 for training and 500 for testing.

4.1.2 Feature. Both the shallow features and deep features are considered in this paper. Similar to prior works, for shallow features, color histogram, optical flow and SIFT features are extracted for each frame, they exploit the appearance, motion and local information, respectively [26]. While for deep features, GoogLeNet is employed to extract the frame features, which is widely used in computer vision tasks [28].

4.1.3 Evaluation. As in prior works [8, 9, 36], the quality of the generated summary is evaluated by comparing to the reference summary created by human. There are three most frequently used evaluation metrics, i.e., precision (the correct subshots to subshots

in the generated summary), recall (the correct subshots to subshots in the reference summary), F-measure (the harmonic mean of precision and recall), which are also employed in this paper.

4.1.4 Parameters. We have tested several H-RNNs under different lengths of LSTMs in each layer, and they get stable performance when the LSTM length varies from 25-60 in each layer (40 is the most favorable). Therefore, in our H-RNN, the LSTM length of the first and second layer are both fixed as 40, so that the input flow at most goes through 80 steps. It can reduce the information loss caused by longer input flow. As a result, our H-RNN can handle the frame sequence less than 1600. For videos contains more than 1600 frames, they are sampled to meet the constraint, while for videos with fewer than 1600 frames, they are padded with zeros.

4.2 Results on the Combined dataset

Table 1 presents the performance of various approaches on the Combined dataset. It provides the results of different approaches on both shallow feature and deep feature. In this paper, for a fair comparison, these approaches are equipped with the same feature. Specifically, for all approaches, the shallow feature is generated by the combination of color histogram, optical flow and SIFT, and the deep feature is extracted by the pool5 layer of GoogLeNet. The two kinds of features are widely used in video summarization. From Table 1, it can be observed that most of the approaches perform better with deep feature rather than shallow feature, it benefits from the great capability of CNN to extract the visual information in the video.

In Table 1, the compared approaches are from different types. The first five are non-RNN-based approaches. VSUMM and LiveLight summarize the video based on clustering and dictionary learning, respectively. CSUV and LSMO exceed them. Specifically, CSUV builds an *Interestingness* model to predict the importance of each subshot. Moreover, LSMO is an extension of CSUV that combines several property models, including the *Interestingness* model proposed in CSUV, together with *Representativeness* model to constrain the key subshots to be representative to the video content, and *Uniformity* model to demand key subshots distributing uniformly. Actually, the following two models are utilized to exploit the relationships between key subshots, i.e., temporal dependency. The better performance of LSMO than CSUV indicates that the temporal dependency is necessary for the task of video summarization. Summary Transfer also exploits the temporal dependency, meanwhile, it utilizes the category label of videos, which have achieved state-of-the-art results in non-RNN-based approaches. While the rest of the approaches are all RNN-based (i.e., LSTM). They perform better than non-RNN-based approaches even without auxiliary information, like the video category label in Summary Transfer. The better performance of RNN-based approaches has verified the superiority of LSTM in temporal dependency modeling.

Actually, vsLSTM and dppLSTM are two approaches most related to H-RNN. Concretely, vsLSTM is constructed by integrating a bi-directional LSTM with a *Multi-Layer Perception* (MLP). Practically, it is reported that MLP is helpful in improving the summary quality [32, 36]. However, it also increases the computation burden and the training parameters in the network. The proposed H-RNN outperforms vsLSTM even without the MLP, it benefits from the

Table 1: The results (F-measure) of various approaches on the Combined dataset. (The scores in bold indicate the best values.)

Feature	shallow feature			deep feature		
	SumMe	Tvsum	MED	SumMe	Tvsum	MED
VSUMM [5]	0.328	0.390	0.260	0.335	0.391	0.263
LiveLight [37]	0.357	0.460	0.258	0.384	0.477	0.262
CSUV [8]	0.393	0.532	0.277	–	–	–
LSMO [9]	0.397	0.548	0.283	0.403	0.568	0.285
Summary Transfer [35]	0.397	0.543	0.292	0.409	0.541	0.297
vsLSTM [36]	0.406	0.571	0.288	0.421	0.580	0.293
dppLSTM [36]	0.407	0.579	0.294	0.429	0.597	0.296
H-RNN	0.421	0.602	0.312	0.443	0.621	0.311

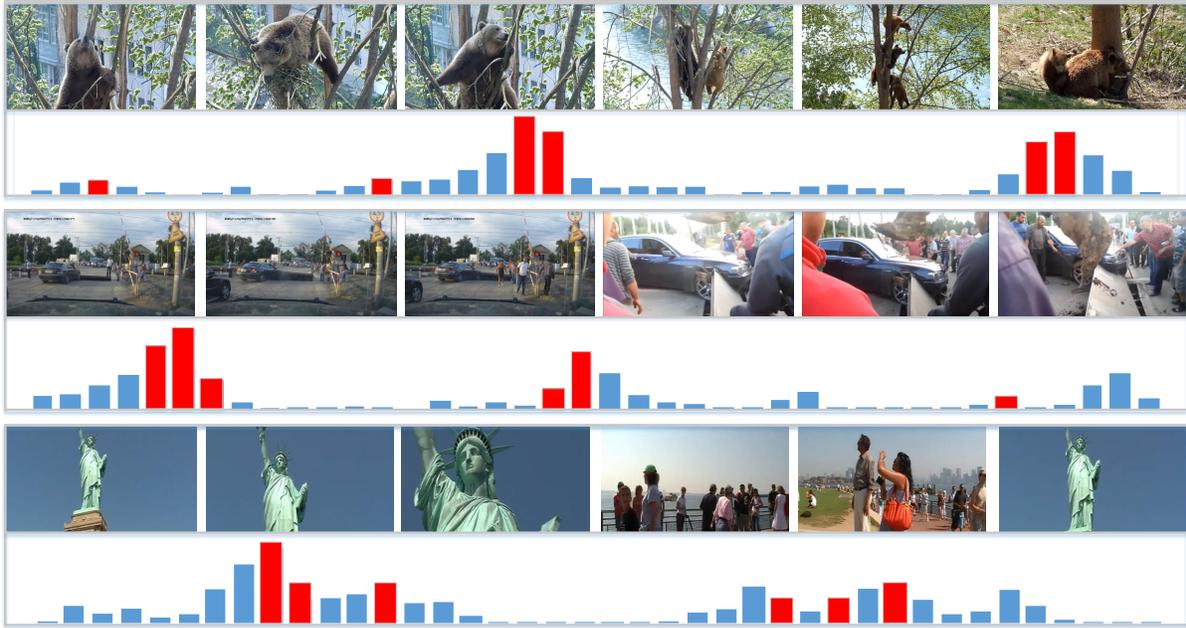


Figure 4: Three example summaries generated by H-RNN, where the key subshots in the summary are represented with several frames (one frame for each shot), and the histograms below the frames denote the distribution of human-annotated shot scores. The red histograms indicates the indexes of selected key subshots. Specifically, six key shots are displayed for each video, and they are corresponding to the upper frames sequentially.

Table 2: The results of various approaches on the VTW dataset.

Feature	shallow feature			deep feature		
	Precision	Recall	F-measure	Precision	Recall	F-measure
CSUV [8]	0.367	0.423	0.399	–	–	–
HD-VS [32]	–	–	–	0.392	0.483	0.427
vsLSTM [36]	0.388	0.490	0.433	0.397	0.495	0.446
H-RNN	0.421	0.522	0.467	0.432	0.528	0.478

nonlinear fitting ability enhanced by the hierarchical structure of

LSTM. Besides, dppLSTM is an extension of vsLSTM by adding a

Table 3: The results of various LSTM-based approaches on the VTW dataset.

Metrics	Precision	Recall	F-measure
single LSTM (mean pool)	0.366	0.442	0.387
single LSTM (uniform sampling)	0.336	0.468	0.391
bi-directional LSTM (mean pool)	0.383	0.502	0.437
bi-directional LSTM (uniform sampling)	0.387	0.495	0.431
H-RNN (single)	0.392	0.523	0.455
H-RNN (bi-directional)	0.441	0.542	0.480

Determinantal Point Process (DPP) model, which is proved effective in representative and diversity subset selection [7, 36]. But the dpplSTM is much more complex than vsLSTM, also is hard to train. Generally, our H-RNN performs better than dpplSTM with a more compact architecture and with less computation. In conclusion, the better performance of H-RNN than vsLSTM and dpplSTM have verified the effectiveness and efficiency of H-RNN in the task of video summarization.

In Figure 4, we present exemplar summaries generated by H-RNN. From the displayed key subshots and the human-annotated scores below them, it can be observed that most high score subshots are selected into our summaries, and the generated summaries can represent the original video content well. It indicates that, in most occasions, the summaries generated by H-RNN basically meet the human demand.

4.3 Results on the VTW dataset

Table 2 shows the results of various approaches on the VTW dataset. Actually, many existing summarization approaches are based on manually designed criteria, so that they are quite dataset dependent. As a result, some of them are not suitable for the VTW dataset, and get very poor performance. For simplicity, they are not listed here.

In Table 2, the results of three compared approaches are provided, where CSUV and vsLSTM have been introduced before. HD-VS summarizes the video by integrating two CNNs, i.e., AlexNet [13] and C3D [29], together with two MLPs after the two CNNs. Specifically, AlexNet is employed to extract the visual information in each frame, and C3D is a 3D convolutional neural network that utilized to exploit the short-range temporal dependency. It can be observed from Table 2 that RNN-based approaches, i.e., vsLSTM and H-RNN, show better results than HD-VS. It is because that LSTM does better in exploiting the temporal dependency than C3D, let alone in long frame sequence. Besides, the even better performance of H-RNN than vsLSTM also shows the superiority of the hierarchical structure of LSTM in video summarization.

To verify the necessity of the structure of H-RNN, the results of several approaches based on LSTM are listed in Table 3. Particularly, considering that LSTM does not work well with long frame sequence, the length of the compared approaches, i.e., single LSTM and bi-directional LSTM, are both fixed as 80. Limited by this, the frame features input to single LSTM and bi-directional LSTM are generated by the mean pooling or uniform sampling of the full frame feature sequence. It can be observed that the two frame feature treatment methods get comparable results. But the significantly

better performance of bi-directional LSTM than single LSTM indicates that both the forward and backward temporal dependency are important for video summarization. It is also the reason that H-RNN (bi-directional) performs better than H-RNN (single), where H-RNN (bi-directional) denotes the second layer of H-RNN is a bi-directional LSTM, and the second layer of H-RNN (single) is a single LSTM. Besides, the better performance of H-RNN than single LSTM and bi-directional LSTM indicates that: 1) The hierarchical structure of H-RNN is more suitable for video summarization since it can deal with long frame sequence, while to single LSTM and bi-directional LSTM, they can only get satisfied results by mean pooling or uniformly sampling the frame feature sequence, which causes inevitable information loss. 2) The hierarchical structure of H-RNN increases the capability of non-linear fitting, which is helpful to the task of video summarization.

5 CONCLUSIONS

In this paper, we propose a hierarchical structure of RNN to enhance the capability of traditional RNN in long-range temporal dependency capturing. Particularly, for the task of video summarization, we design a specialized two-layer RNN according to the layered video structure, called as H-RNN. Particularly, the first layer is a LSTM, which is utilized to exploit the intra-subshot temporal dependency among frames. The second layer is a bi-directional LSTM that can capture both the forward and backward inter-subshot temporal dependency, and the output of the second layer is utilized to predict whether a certain subshot is valuable to be selected into the summary. Compared to current RNN-based approaches, H-RNN is more suitable to the task of video summarization, and the experimental results have verified its superiority.

REFERENCES

- [1] Aya Aner and John R. Kender. 2002. Video Summaries through Mosaic-Based Shot and Scene Clustering. In *Computer Vision - ECCV 2002, 7th European Conference on Computer Vision, Copenhagen, Denmark, May 28-31, 2002, Proceedings, Part IV*. 388–402.
- [2] Yoshua Bengio, Patrice Simard, and Paolo Frasconi. 1994. Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks* 5, 2 (1994), 157–166.
- [3] Wen-Sheng Chu, Yale Song, and Alejandro Jaimes. 2015. Video co-summarization: Video summarization by visual co-occurrence. In *IEEE Conference on Computer Vision and Pattern Recognition*. 3584–3592.
- [4] Yang Cong, Junsong Yuan, and Jiebo Luo. 2012. Towards Scalable Summarization of Consumer Videos Via Sparse Dictionary Selection. *IEEE Trans. Multimedia* 14, 1 (2012), 66–75.
- [5] Sandra Eliza Fontes de Avila, Ana Paula Brandão Lopes, Antonio da Luz Jr., and Arnaldo de Albuquerque Araújo. 2011. VSUMM: A mechanism designed to produce static video summaries and a novel evaluation method. *Pattern Recognition Letters* 32, 1 (2011), 56–68.

- [6] Ehsan Elhamifar, Guillermo Sapiro, and René Vidal. 2012. See all by looking at a few: Sparse modeling for finding representative objects. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*. 1600–1607.
- [7] Boqing Gong, Wei-Lun Chao, Kristen Grauman, and Fei Sha. 2014. Diverse Sequential Subset Selection for Supervised Video Summarization. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*. 2069–2077.
- [8] Michael Gygli, Helmut Grabner, Hayko Riemenschneider, and Luc Van Gool. 2014. Creating Summaries from User Videos. In *European Conference on Computer Vision*. 505–520.
- [9] Michael Gygli, Helmut Grabner, and Luc Van Gool. 2015. Video summarization by learning submodular mixtures of objectives. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3090–3098.
- [10] Youssef Hadi, Fedwa Essannouni, and Rachid Oulad Haj Thami. 2006. Video summarization by k-medoid clustering. In *Proceedings of the 2006 ACM symposium on Applied computing*. ACM, 1400–1401.
- [11] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [12] Aditya Khosla, Raffay Hamid, Chih-Jen Lin, and Neel Sundaresan. 2013. Large-Scale Video Summarization Using Web-Image Priors. In *2013 IEEE Conference on Computer Vision and Pattern Recognition, Portland, OR, USA, June 23-28, 2013*. 2698–2705.
- [13] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. 2012. ImageNet Classification with Deep Convolutional Neural Networks. In *Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States*. 1106–1114.
- [14] Yong Jae Lee, Joydeep Ghosh, and Kristen Grauman. 2012. Discovering important people and objects for egocentric video summarization. In *IEEE Conference on Computer Vision and Pattern Recognition*. 1346–1353.
- [15] David Liu, Gang Hua, and Tsuhan Chen. 2010. A Hierarchical Visual Model for Video Object Summarization. *IEEE Trans. Pattern Analysis and Machine Intelligence* 32, 12 (2010), 2178–2190.
- [16] Tiecheng Liu and John R. Kender. 2002. Optimization Algorithms for the Selection of Key Frame Sequences of Variable Length. In *Computer Vision - ECCV 2002, 7th European Conference on Computer Vision, Copenhagen, Denmark, May 28-31, 2002, Proceedings, Part IV*. 403–417.
- [17] Shiyang Lu, Zhiyong Wang, Tao Mei, Genliang Guan, and David Dagan Feng. 2014. A Bag-of-Importance Model With Locality-Constrained Coding Based Feature Learning for Video Summarization. *IEEE Trans. Multimedia* 16, 6 (2014), 1497–1509.
- [18] Zheng Lu and Kristen Grauman. 2013. Story-driven summarization for egocentric video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2714–2721.
- [19] Qiao Luan, Mingli Song, Chu Yee Liao, Jiajun Bu, Zicheng Liu, and Ming-Ting Sun. 2014. Video Summarization based on Nonnegative Linear Reconstruction. In *IEEE International Conference on Multimedia and Expo*. 1–6.
- [20] Padmavathi Mundur, Yong Rao, and Yelena Yesha. 2006. Keyframe-based video summarization using Delaunay clustering. *Int. J. on Digital Libraries* 6, 2 (2006), 219–232.
- [21] Joe Yue-Hei Ng, Matthew J. Hausknecht, Sudheendra Vijayanarasimhan, Oriol Vinyals, Rajat Monga, and George Toderici. 2015. Beyond short snippets: Deep networks for video classification. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*. 4694–4702.
- [22] Chong-Wah Ngo, Yu-Fei Ma, and HongJiang Zhang. 2003. Automatic Video Summarization by Graph Modeling. In *9th IEEE International Conference on Computer Vision (ICCV 2003), 14-17 October 2003, Nice, France*. 104–109.
- [23] Pingbo Pan, Zhongwen Xu, Yi Yang, Fei Wu, and Yueting Zhuang. 2016. Hierarchical Recurrent Neural Encoder for Video Representation with Application to Captioning. In *2016 IEEE Conference on Computer Vision and Pattern Recognition*. 1029–1038.
- [24] Danila Potapov, Matthijs Douze, Zaïd Harchaoui, and Cordelia Schmid. 2014. Category-Specific Video Summarization. In *European Conference on Computer Vision*. 540–555.
- [25] Karen Simonyan and Andrew Zisserman. 2014. Very Deep Convolutional Networks for Large-Scale Image Recognition. *CoRR* abs/1409.1556 (2014).
- [26] Yale Song, Jordi Vallmitjana, Amanda Stent, and Alejandro Jaimes. 2015. Tvsum: Summarizing web videos using titles. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 5179–5187.
- [27] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to Sequence Learning with Neural Networks. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*. 3104–3112.
- [28] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott E. Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2015. Going deeper with convolutions. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*. 1–9.
- [29] Du Tran, Lubomir D. Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. 2014. C3D: Generic Features for Video Analysis. *CoRR* abs/1412.0767 (2014).
- [30] Subhashini Venugopalan, Marcus Rohrbach, Jeffrey Donahue, Raymond J. Mooney, Trevor Darrell, and Kate Saenko. 2015. Sequence to Sequence - Video to Text. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*. 4534–4542.
- [31] Huan Yang, Baoyuan Wang, Stephen Lin, David P. Wipf, Minyi Guo, and Baining Guo. 2015. Unsupervised Extraction of Video Highlights via Robust Recurrent Auto-Encoders. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*. 4633–4641.
- [32] Ting Yao, Tao Mei, and Yong Rui. 2016. Highlight Detection with Pairwise Deep Ranking for First-Person Video Summarization. In *2016 IEEE Conference on Computer Vision and Pattern Recognition*. 982–990.
- [33] Wojciech Zaremba and Ilya Sutskever. 2014. Learning to Execute. *CoRR* abs/1410.4615 (2014).
- [34] Kuo-Hao Zeng, Tseng-Hung Chen, Juan Carlos Niebles, and Min Sun. 2016. Title Generation for User Generated Videos. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II*. 609–625.
- [35] Ke Zhang, Wei-Lun Chao, Fei Sha, and Kristen Grauman. 2016. Summary Transfer: Exemplar-Based Subset Selection for Video Summarization. In *2016 IEEE Conference on Computer Vision and Pattern Recognition*. 1059–1067.
- [36] Ke Zhang, Wei-Lun Chao, Fei Sha, and Kristen Grauman. 2016. Video Summarization with Long Short-Term Memory. In *Computer Vision - ECCV 2016 - 14th European Conference*. 766–782.
- [37] Bin Zhao and Eric P. Xing. 2014. Quasi Real-Time Summarization for Consumer Videos. In *2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014*. 2513–2520.
- [38] Yueting Zhuang, Yong Rui, Thomas S. Huang, and Sharad Mehrotra. 1998. Adaptive Key Frame Extraction using Unsupervised Clustering. In *Proceedings of the 1998 IEEE International Conference on Image Processing, ICIP-98, Chicago, Illinois, October 4-7, 1998*. 866–870.