# Supplementary Materials: Towards Bridging Event Captioner and Sentence Localizer for Weakly Supervised Dense Event Captioning

Shaoxiang Chen and Yu-Gang Jiang\* Shanghai Key Lab of Intelligent Information Processing, School of Computer Science, Fudan University

{sxchen13, ygj}@fudan.edu.cn

### Algorithm 1: Steps In One Training Iteration

Input: Video V, Annotated Caption C (can be batched). Inherit the Event Captioner  $\mathcal{C}$ , Sentence Localizer  $\mathcal{L}$ , and Induced Set Attention Block (ISAB)  $\mathcal{I}$  from previous iteration: Load a predefined set of event proposals  $G, Q \leftarrow \emptyset$ ; for  $G_i \in \boldsymbol{G}$  do  $\mathcal{C}(\mathcal{I}(V), G_i) \to (C_i, S_i);$ // Output captions and features. // ISAB is abstracted as a feature transformer for clarity. Q.append( $(G_i, C_i, S_i)$ ); Select the best  $(\hat{G}, \hat{C}, \hat{S})$  from Q based on caption confidences;  $\mathcal{L}(V, \hat{S}) \rightarrow (W^s, W^v);$ // Localize the sentence. // Losses that connect  ${\mathcal C}$  and  ${\mathcal L}$ . For  $\mathcal{C}$ , compute captioning loss with  $(\hat{C}, C)$  guided by  $W^s$ : For  $\mathcal{L}$ , compute localization loss between  $\hat{G}$  and  $W^{v}$ ; // Losses within  $\mathcal{L}$ .

For  $\mathcal{L}$ , compute contrastive loss and MIL loss within  $\mathcal{L}$ ; Perform gradient-based parameter updates for  $\{\mathcal{C}, \mathcal{L}, \mathcal{I}\}$ ;

# **1. Algorithm Description**

We show the simplified sequence of steps during training and inference in Algorithms 1 and 2.

# 2. More Model Details

**Masked Temporal Attention.** The masked temporal attention (adapted from [1]) in our event captioner (Eq. (16))

$$\begin{split} & \boldsymbol{h}_{t}^{(1)} = \texttt{LSTM}^{(1)}(\texttt{att}(\widetilde{\boldsymbol{V}}, M(:, G_{i}), \boldsymbol{h}_{t-1}^{(1)}), \boldsymbol{h}_{t-1}^{(1)}), \\ & \boldsymbol{h}_{t}^{(2)} = \texttt{LSTM}^{(2)}([\texttt{embed}(w_{t-1}), \boldsymbol{h}_{t}^{(1)}], \boldsymbol{h}_{t-1}^{(2)}), \end{split}$$

\*Corresponding author.

#### Algorithm 2: The Inference Process for One Video

Input: Video V. Output: Segment predictions with captions. Load trained parameters for  $\mathcal{C}, \mathcal{L}, \mathcal{I}$ ; Initialize the proposals  $G, Q \leftarrow \emptyset$ ; for *iter*  $\leftarrow 1$  to *max\_iter* do // max\_iter can be 1. for  $G_i \in \boldsymbol{G}$  do  $\mathcal{C}(\mathcal{I}(m{V}),G_i) 
ightarrow (m{C}_i,m{S}_i);$  // Proposal to sentence. Q.append( $(G_i, C_i, S_i)$ ); Select the best  $(\hat{G}, \hat{C}, \hat{S})$  from Q; // Similar to training.  $\mathcal{L}(V, \hat{S}) \rightarrow (W^s, W^v);$ // Localize the sentence.  $G \leftarrow TAG(W^v);$ // Generate new proposals with TAG, and refine Gfor the next iteration. return top K of G, captioned and sorted by caption confidences.

is formulated as

$$\begin{split} \mathtt{att}(\widetilde{\boldsymbol{V}}, \boldsymbol{M}(:, G_i), \boldsymbol{h}_{t-1}^{(1)}) &= \sum_{l=1}^{L} \beta_l \widetilde{\boldsymbol{V}}_l, \\ \mathtt{where} \quad \beta_l &= \frac{\boldsymbol{M}(l, G_i) \exp(r_l)}{\sum_l \boldsymbol{M}(l, G_i) \exp(r_l)}, \\ r_l &= W_r( \tanh(W_{r_V} \widetilde{\boldsymbol{V}}_l + W_{r_h} \boldsymbol{h}_{t-1}^{(1)} + b_r)), \end{split}$$

where  $W_r$ ,  $W_{r_V}$ ,  $W_{r_h}$  and  $b_r$  are trainable parameters. Masked Temporal Attention is mainly for temporally aggregating video features within the segment indicated by the mask. The value of  $M(l, G_i)$  is 1 if  $s_i \leq l \leq e_i$ , and otherwise is 0.

**More Implementation Details.** The video and sentence features are linearly transformed with a ReLU activation in

	Value	Μ	С	R	B@4
Concept number	1024	7.37	19.86	12.82	1.25
	2186	7.49	21.21	13.02	1.33
	All (5305)	7.45	20.73	12.96	1.29
Concept type	Verb (615)	7.18	19.61	12.79	1.23
	Noun (1571)	7.34	19.92	12.85	1.26
	Verb+Noun (2186)	7.49	21.21	13.02	1.33

Table 1: Study on the number and type of concepts.

Eq. (6). The number of heads in the MAB( $\cdot$ ) function is generally set to 4. The FFN( $\cdot$ ) has two fully connected layers with ReLU activation and a residual connection from its inputs to the outputs. The  $\gamma$  in Eq. (18) is set to 0.1. A dropout layer with ratio 0.5 is applied to each of the LSTM outputs (Eq. (16)). As in most video captioning methods, we generate captions using beam search and the beam size is set to 5.

# 3. More Ablation Experiments

Table 1 above shows the effects of the number and type of concepts. We can conclude that (1) as the number increases, there is a point (2186, our final choice) where we can get more information from more concepts to enhance the features without having too much noise and making the concepts too long-tailed; (2) it is better to include verbs and nouns since they are both important aspects of an event. But noun concepts are the majority and yield better performances than verb concepts.

# 4. Visualization of the Learned Space

To show how our method models the video-sentence interaction and learn their joint space, we visualize the learned features with t-SNE<sup>1</sup> in Figure 1. We can see that the video and sentence subspace are initially separated, and the localization loss we optimize will pull them close and make them overlap as the training proceeds. More specifically, we plot the features of multiple clips from a video, and we can clearly observe that the sentence and video features are gradually pulled close. Note that although we did not explicitly enforce the similar videos' (or sentences') features to be close to each other, the features from the same video's clips and sentences appear to be clustered, which may be because the underlying concepts of the video are captured and can make it easier to cluster similar clips.

# **5. More Qualitative Results**

In Figures 2 and 3, we provide more qualitative results generated by our method on the testing set. Each figure contains three examples. Figure 2 shows high-quality examples, where our method can correctly localize and describe the key events, e.g. 'a game of cricket' and 'volleyball'. From the examples in Figure 3, we can see that one weakness of the method is that it can generate repetitive captions for multiple segments. This is due to that inter-segment relation is not modeled and the captioning process of each segment is independent.

## References

 Li Yao, Atousa Torabi, Kyunghyun Cho, Nicolas Ballas, Christopher J. Pal, Hugo Larochelle, and Aaron C. Courville. Describing videos by exploiting temporal structure. In *ICCV*, pages 4507–4515, 2015. 1

<sup>&</sup>lt;sup>1</sup>We also note that t-SNE is sensitive to its hyperparameter settings.



(c) Feature space at 14k iterations.

(d) Feature space at 20k iterations.

Figure 1: t-SNE visualization of the joint feature space (around 2,000 sentence-clip pairs shown here) learned at different training stages. The '+' markers denote samples from a video bjtjeUcoxkg.mp4.



Figure 2: More qualitative examples.



Figure 3: More qualitative examples.