# Supplementary Material to: Weakly Supervised Dense Event Captioning in Videos

This Supplementary material provides proof of proposition 1, and presents more details about the model, training, and testing.

#### A **Proof of proposition 1(Fixed-point iteration)**

Proposition 1 (Fixed-Point-Iteration). We define the iteration as

$$\boldsymbol{S}(t+1) = l_{\theta_1}(\boldsymbol{V}, g_{\theta_2}(\boldsymbol{V}, \boldsymbol{S}(t))), \qquad (18)$$

where S(t) will converge to the fixed-point solution, i.e.  $S^* = l_{\theta_1}(V, g_{\theta_2}(V, S^*))$ , if the initial start S(0) surrounds the fixed point  $S^*$  sufficiently and the function  $l_{\theta_1}(V, g_{\theta_2}(V, S))$  is locally Lipschitz continuous around  $S^*$  with Lipschitz constant L < 1.

*Proof.* Since the function  $f(S) = l_{\theta_1}(V, g_{\theta_2}(V, S))$  is Lipschitz continuous with Lipschitz constant  $L \leq 1$ , we have,

$$|\boldsymbol{S}(t+1) - \boldsymbol{S}(t)| = |f(\boldsymbol{S}(t)) - f(\boldsymbol{S}(t-1))|$$
  

$$\leq L|\boldsymbol{S}(t) - \boldsymbol{S}(t-1)|,$$
  

$$\vdots$$
  

$$\leq L^{t}|\boldsymbol{S}(1) - \boldsymbol{S}(0)|.$$
(19)

Since L < 1,  $|S(t+1) - S(t)| \to 0$  converges to zeros as  $t \to \infty$ . Thus, S(t) will converge to the fixed point that is the solution of Eq.(2) in the paper.

# **B** Video Encoder & Sentence Encoder

In § 3, we have introduced the details of our model except the RNN-based video and sentence encoder used in feature extracting, which are considered not crucial in the paper. But for reproduction of our work, we will introduce them detailedly in this section.

#### B.1 Video Encoder.

Reviewing § 3, the video encoder aims to encode the video frames into a set of vectors. For this purpose, a C3D[1] network is constructed to extract frame-block features and resolve short-term dependency, while a followed GRU[2] is leveraged for long-term dependency.

We denote the input video as  $V = \{f_t\}_{t=0}^{T_v}$  where  $T_v$  is the temporal length of the video and  $f_t$  specifies the *t*-th frame. Following the implementation by[3], we first extract the non-overlapping C3D features from the original image frames with a interval  $\delta$ , i.e., computing  $\{x_t^{(v)} = C3D(f_{(t-1)*\delta+1} : f_{t*\delta})\}_{t=0}^{T_v/\delta}$ , where  $\delta = 16$ .

The extracted C3D features are fed into a GRU network to learn the long-term dependence. Specially, we set the initial hidden state as zero, i.e.  $h_0^{(v)} = \mathbf{0}$ , and compute the following outputs and hidden states recursively by

32nd Conference on Neural Information Processing Systems (NeurIPS 2018), Montréal, Canada.

$$\mathbf{v}_t, \mathbf{h}_t^{(v)} = GRU_{video}(\mathbf{x}_t^{(v)}, \mathbf{h}_{t-1}^{(v)}).$$
 (20)

The output sequence  $\{v_t\}_{t=0}^{T_v/\delta}$  and hidden states  $\{h_t^{(v)}\}_{t=0}^{T_v/\delta}$  are used in the followed processes of sentence localizing and caption generating as video representation (In the original paper, we omit  $\delta$  and denote  $T_v/\delta$  as  $T_v$  and the extract feature as  $\{v_t\}_{t=0}^{T_v}$  for simplify).

#### **B.2** Sentence Encoder

Similar to the video encoder, sentence encoder aims to encode natural sentences into vector representations. Given a sentence  $C = \{w_t\}_{t=0}^{T_c}, w_t \in \{0, 1\}^V$  denote the one-hot encoding of the *t*-th word, where V is the vocabulary size. Then,  $w_t$  is embedded into a k-dimensional vector by  $x_t^{(c)} = W_e w_t$ , with  $W_e \in \mathcal{R}^{k \times V}$  being the trainable embedding matrix. Similar to the video encoder, the embedded features  $\{x_t^{(c)}\}_{t=0}^{T_c}$  are fed into a GRU model. Specially, the initial hidden state  $h_0^{(c)}$  is set to be zero, and the outputs and following hidden states are computed as:

$$\boldsymbol{c}_t, \boldsymbol{h}_t^{(c)} = GRU_{text}(\boldsymbol{x}_t^{(c)}, \boldsymbol{h}_{t-1}^{(c)}).$$
(21)

After the encoding stage, the final outputs  $\{c_t\}_{t=0}^{T_c}$  and hidden states  $\{h_t^{(c)}\}_{t=0}^{T_c}$  are used in the followed networks.

#### **B.3** More Details

**I.** As both **Sentence Localizer**  $l_{\theta_1}$  and **Caption Generator**  $g_{\theta_2}$  use the video features for further processing, they can choose whether to share the same video encoder parameters or not. In our previous experiments, we found that both strategies reach similar performances. The reported results are obtained using different video encoder parameters.

II. Also, as training a C3D network from scratch is very time-consuming and error-prone, we directly adopt the public-available C3D features<sup>1</sup>. The public-available features are 500-way features reduced from the original 4096-way output of C3D's fc7 layer using PCA. In our experiments, the features are denoted as  $\{x_t^{(v)}\}_{t=0}^{T_v/\delta}$  and fed into the  $GRU_{video}$  for further processing.

# C Training & Testing Details

#### C.1 Training

Because our model consists of two submodels and several term losses, there may be several ways to train the whole model. In this section, we introduce the strategy we used in our experiments. Firstly, the model is pretrained with our predefined *Fake Proposal*, i.e.  $S^{(f)} = (0.5, 1)$ . After pretraining for several rounds, we train the proposed three term losses alternatively. Details are shown in Algorithm1.

### C.2 Testing

As our model is not directly trained on the dense event captioning problem, we provide an extra explanation on the strategy we used in testing. In short, our training strategy and losses force the model meets the requirement of fixed-point iteration, so we use a random bunch of segments  $\{S_i^{(r)}\}_i^{N_r}$  as initial temporal segments, feed them into the cycle system  $S'_i = l_{\theta_1}(V, g_{\theta_2}(V, S_i^{(r)}))$ . Considering that those random segments will converge to the true event temporal segments or diverge to some unknown random segment, the distance between  $S'_i$  and  $S_i^{(r)}$  should be small if they are in the neighborhood of a certain event. Also, we measure the distance between  $S'_i$  and  $S'_i$  to further reduce redundancy(This step is not critical, but can reduce the repetition of temporal segments.). Specifically:

• if  $dist(S'_i, S^{(r)}_i) > \Theta_1$ , remove  $S'_i$  from the predicted temporal segment set

<sup>&</sup>lt;sup>1</sup>http://activity-net.org/challenges/2016/download.html#c3d

• if  $dist(\boldsymbol{S}_i',\boldsymbol{S}_j') < \Theta_2,$  merge  $\boldsymbol{S}_i',\boldsymbol{S}_j'$  as:  $\boldsymbol{S}_i'' = union(\boldsymbol{S}_i',\boldsymbol{S}_j')$ 

where  $dist(\cdot)$  computes the IoU between  $S_1, S_2$ :

$$dist(\mathbf{S}_1, \mathbf{S}_2) = \frac{intersection(\mathbf{S}_1, \mathbf{S}_2)}{union(\mathbf{S}_1, \mathbf{S}_2)}$$
(22)

Details about testing are shown in Algorithm2.

Algorithm 1 Training pipeline for the WS-DEC model Input: D// the dataset iterator for training **Input:**  $l_{\theta_1}, g_{\theta_2}$ // the random initialized sentence localizer, caption generator Input:  $\{S_{j}^{(a)}\}_{j=0}^{N_{a}}$ // anchor segments for training  $l_{\theta_1}$ //  $\boldsymbol{S}^{(f)} = (m^{(f)}, w^{(f)})$  is used to pretrain the caption generator Input:  $S^{(f)}$ Output:  $\theta_1, \theta_2$ // trained parameters for sentence localizer and caption generator 1:  $step \leftarrow 0$ 2: while *step < pretrain\_step* do // we pretrain the model with fake proposal for  $(V, \{Ci\}_{i=0}^{N_v}) \in \mathcal{D}$  do 3:  $C \leftarrow \text{RANDOMCHOOSE}(\{Ci\}_{i=0}^{N_v}) // \text{ randomly choose a sentence}$ 4:  $\begin{array}{l} \mathbf{C}' \leftarrow g_{\theta_2}(\mathbf{V}, \mathbf{S}^{(f)}) \, / \text{ obtain the fake generation} \\ \mathcal{L}_c \leftarrow dist(\mathbf{C}, \mathbf{C}') \, / \text{compute the loss} \\ \theta_2 \leftarrow \theta_2 - \eta_2 \frac{\partial \mathcal{L}_c}{\partial \theta_2} \, / \text{ update parameters with SGD} \end{array}$ 5: 6: 7: 8: end for  $step \leftarrow step + 1$ 9: 10: end while 11: while  $step < training\_step$  do // we train the model with three term losses alternatively 12: if training  $L_c$  and  $L_s$  then // train the model with  $\mathcal{L}_c$  and  $\mathcal{L}_s$ for  $(V, \{Ci\}_{i=0}^{N_v}) \in \mathcal{D}$  do 13:  $C \leftarrow \text{RANDOMCHOOSE}(\{Ci\}_{i=0}^{N_v})$ 14:  $\begin{array}{l} \mathbf{C} \leftarrow \mathsf{RANDOM CHOOSE}(\{\mathbf{C}i\}_{i=0}^{*} \\ \mathbf{S}' \leftarrow l_{\theta_1}(\mathbf{V}, \mathbf{C}) \\ \mathbf{C}' \leftarrow g_{\theta_2}(\mathbf{V}, \mathbf{S}' + \delta) \\ \mathbf{S}'' \leftarrow g_{\theta_2}(\mathbf{V}, \mathbf{C}') \\ \mathcal{L} \leftarrow dist(\mathbf{C}, \mathbf{C}') + dist(\mathbf{S}', \mathbf{S}'') \\ \theta_1 \leftarrow \theta_1 - \eta_{c1} \frac{\partial \mathcal{L}}{\partial \theta_1} \\ \theta_2 \leftarrow \theta_2 - \eta_{c2} \frac{\partial \mathcal{L}}{\partial \theta_2} \end{array}$ 15: 16: 17: 18: 19: 20: 21: end for 22: end if if training  $L_a$  then // train the model with  $\mathcal{L}_a$ 23: for  $(V, \{Ci\}_{i=0}^{N_v}) \in \mathcal{D}$  do 24:  $C \leftarrow \text{RANDOMCHOOSE}(\{Ci\}_{i=0}^{N_v})$ 25:  $\{Cj^{(a)}\}_{j=0}^{N_a} \leftarrow \{g_{\theta_2}(V, S_j^{(a)})\}_{j=0}^{N_a}$  //anchor segments  $\rightarrow$  anchor sentences confidence  $\leftarrow l_{\theta_1}(V, C)$  // confidence about each anchor regarding V and C26: 27:  $\{score_{j}^{(a)}\}_{j=0}^{N_{a}} \leftarrow \{\text{ METEOR}(\boldsymbol{C}, \boldsymbol{C}j^{(a)})\}_{j=0}^{N_{a}} // \text{ compute Meteor score}$ 28:  $label \leftarrow \arg \max_{j} \{score_{j}^{(a)}\}_{j=0}^{N_{a}} // \text{ the one with highest score is considered as label.}$  $\mathcal{L}_{a} \leftarrow \text{CROSSENTROPY}(confidence, label)$  $\theta_{1} \leftarrow \theta_{1} - \eta_{a1} \frac{\partial \mathcal{L}_{a}}{\partial \theta_{1}} // \text{ cross entropy loss}$ 29: 30: 31: end for 32: 33: end if 34:  $step \leftarrow step + 1$ 35: end while 36: return  $\theta_1, \theta_2$ 

Algorithm 2 Testing pipeline for the WS-DEC model Input: V // a video for testing(never seen in the training set) **Input:**  $l_{\theta_1}, g_{\theta_2}$ // the learned sentence localizer, caption generator **Input:**  $\{S_i^{(r)}\}_{i=0}^{N_r}$ // initial random bunch temporal segments(fixed for all testing example) **Output:**  $\{C_i, S_i\}_{i=0}^{N_e^{(v)}}$  // the predicted events within the video 1:  $step \leftarrow 0$  // iteration round 2:  $N_e^{(v)} \leftarrow N_r$  // initial number of event equals the predefined one 2:  $N_e \leftarrow N_r$  // initial induct of event equals the predefined temporal segments 3:  $\{S_i\}_{i=0}^{N_e^{(v)}} \leftarrow \{S_i^{(r)}\}_{i=0}^{N_r}$  // initial temporal segment equals the predefined temporal segments 4:  $\{C_i\}_{i=0}^{N_e^{(v)}} \leftarrow \{g_{\theta_2}(V, S_i)\}_{i=0}^{N_e^{(v)}}$  // generate initial description 5: while  $step < iteration\_step$  do // iterate until max iteration step tile step < iteration\_step do // iterate until max iteration step  $\{S'_i\}_{i=0}^{N_e^{(v)}} \leftarrow \{l_{\theta_1}(V, Ci)\}_{i=0}^{N_e^{(v)}} // \text{ re localize the generated sentence}$   $\{S''_i\}_{i=0}^{N_e^{(v)'}} \leftarrow \text{SEGMENTMERGE}(\{S_i\}_{i=0}^{N_e^{(v)}}, \{S'_i\}_{i=0}^{N_e^{(v)}}) // \text{ merge segments based on§ 4.3.3}$   $N_e^{(v)} \leftarrow N_e^{(v)'} // \text{ update the predicted number of events}$   $\{S_i\}_{i=0}^{N_e^{(v)}} \leftarrow \{S''_i\}_{i=0}^{N_e^{(v)}} // \text{ update events temporal segments}$   $\{C_i\}_{i=0}^{N_e^{(v)}} \leftarrow \{g_{\theta_2}(V, S_i)\}_{i=0}^{N_e^{(v)}} // \text{ update description sentences}$   $step \leftarrow step + 1$  d while6: 7: 8: 9: 10: 11: 12: end while 13: return  $\{Ci, S_i\}_{i=0}^{N_e^{(v)}}$ 

# References

- [1] Shuiwang Ji, Wei Xu, Ming Yang, and Kai Yu. 3d convolutional neural networks for human action recognition. *IEEE transactions on pattern analysis and machine intelligence*, 35(1):221–231, 2013.
- [3] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 706–715, 2017.