

# SRTube: Video-Language Pre-Training with Action-Centric Video Tube Features and Semantic Role Labeling

## Supplementary Material

### 001 1. Implementation Details

#### 002 1.1. Model Architecture.

003 **Video and Text Encoders.** In our SRTube, we adopt the  
004 BEiT [2] architecture as our video encoder which comprises  
005 12 layers with 768 hidden units. Following VindLU [5], two  
006 temporal attention layers [3] are added before self-attention  
007 layers in a video encoder. Initially, the video encoder was  
008 trained for action recognition and trajectory prediction using  
009 the AVA [8] dataset. In this stage, visual encoder is  
010 initialized using the BEiT [2] model, which is pre-trained  
011 on ImageNet[15]. The weights obtained from this stage are  
012 then used to further initialize the visual encoder in the SR-  
013 Tube pre-training stage to enhance its capabilities for advanced  
014 action-centric video analysis. The Tube builder consists of a 6-layer  
015 transformer that follows the DETR [4] decoder architecture and is  
016 initialized with DETR. For text encoding, we use the BERT [6]  
017 architecture and the cross fusion encoder consists of the last three  
018 layers of the BERT model, as suggested in [5].  
019

#### 020 1.2. Downstream Implementation Details

021 **Text to video retrieval.** For retrieval tasks, we jointly optimize  
022 the VTC loss for video-text alignment during fine-tuning. During  
023 inference, we select top-k candidates by computing the dot-product  
024 similarity between the output embedding of VTA and TSA. Then TSA  
025 scores are added to VTA scores and then rerank the selected  
026 candidates based on their total scores.  
027

028 **Video QA.** For video question and answering tasks, we jointly  
029 optimize the VTC, VTM, MAM, and ANM loss for fine-tuning. The  
030 text decoder receives input from the output embeddings of both  
031 VTA and TSA to generate answers. During inference, a high-  
032 confidence answer is selected from the output of the text decoder.  
033

034 **Video Captioning.** We concatenate the video context feature  
035 with the video tube feature and directly input this combination  
036 into a text decoder for caption generation. The caption is generated  
037 from a [MASK] embedding and continues until an [END] token is  
038 reached. To optimize the model, we utilize language modeling loss.  
039 The hyperparameters used for fine-tuning are detailed in the  
040 Table 3.

#### 041 1.3. Datasets descriptions.

042 In this section, we describe all of the video-text pairs used  
043 for evaluation. The details of the datasets are described in  
044 table 4.

### 045 2. Additional Experimental Results.

046 In this section, we present additional experimental results  
047 for text-to-video retrieval, video question and answering for  
048 additional ablation studies.

049 **Tube feature and SRL phrase.** We present the effectiveness  
050 of tube features and SRL features in a fine-tuned model for  
051 ablation studies in Table 1. We present experimental results  
052 on text-to-video retrieval on MSR-VTT and LSMDC datasets,  
053 achieving R@1 metrics. These experiments are conducted using  
054 the same settings as in the main paper.

Method	MSR-VTT	LSMDC
	R@1	R@1
Baseline(V,T)	40.6	21.3
Baseline + Tube (V,U,T)	43.1	24.5
Baseline + SRL (V,T,S)	42.7	24.3
Baseline + Tube + SRL (V, U, T, S)	43.9	25.7

Table 1. Ablation tests of the tube and semantic phrase features.

055 **Video encoder.** We compare the performance of models  
056 with different video encoders, including both traditional (C3D  
057 [16]) and conventional method (VideoSwin [14], BEiT [2]) in  
058 Table 2. We compare results using the TVQA dataset for visual  
059 question answering. The results demonstrate that our model,  
060 initialized with weights derived from the tube builder training  
061 stage, consistently outperforms alternative approaches.  
062

Method	TVQA
C3D [16]	75.1
VideoSwin [14]	76.7
BEiT [2]	78.2
Ours	78.6

Table 2. Ablation tests of video encoder initialization.

### 063 3. Semantic phrase

064 In Fig. 1, We show the results of applying semantic role  
065 labeling to the original video descriptions to get semantic  
066 phrases. Each phrase corresponds to a specific semantic label  
067 highlighted in different colors. The final extracted semantic  
068 phrases are connected using “and” to generate our proposed  
069 semantic phrase which is shown as “filtered description” in  
070 Fig.1.

Datasets	Optimize & Fine-tuning Confgs	# Query	Epoch	BS $\times$ GPUs
MSR-VTT-Ret		8	20	24 X 2
MSR-VTT-QA		8	20	24 X 2
MSR-VTT-Cap	optimizer : AdamW [10]	8	20	24 X 2
LSMDC-Ret	weight decay : 0.02	8	25	24 X 2
DiDeMo-Ret	learning rate scheduler : Cosine Decay	8	20	24 X 2
MSVD-QA	learning rate : $10^{-5}$	8	15	24 X 2
MSVD-Cap	frame resolution : 224 X 224	8	15	24 X 2
TVQA-QA	number of frame : 16	8	20	24 X 2
Activity net-Ret		8	20	48 X 1
SSv2-Template-Ret		4	20	48 X 1
SSv2-Lable-Ret		4	20	48 X 1

Table 3. Fine-tuning configurations for video-language downstream tasks. **BS**: batch size, **Ret** : retrieval task, **QA** : Question and answering task, **Cap** : captioning task. # denotes the number of.

Dataset	Source	# Clip	# Sentence
<b>MSR-VTT</b> [17]	YouTube	10K	200K
<b>MSVD</b> [17]	YouTube	2K	10K
<b>LSMDC</b> [1]	Movie	118K	118K
<b>DiDeMo</b> [9]	Flickr	10K	40K
<b>TVQA</b> [12]	TV show	22K	153K
<b>SSv2-Template</b> [13]	SSv2 [7]	220K	174
<b>SSv2-Label</b> [13]	SSv2	220K	174
<b>ActivityNet-Captions</b> [11]	YouTube	20K	100K

Table 4. Comparison of video-language benchmarks. We describe each benchmark in terms of the source of the original videos, as well as the number of clips and text. # denotes the number of.

## 071 4. Visualization.

072 To illustrate the functionality of our proposed model, we  
073 display examples of predicted object trajectories generated  
074 by the tube builder, alongside results from various down-  
075 stream tasks. We indicate the predicted object trajectory  
076 with red bounding boxes in the image. For example, in  
077 Fig. 2, we show input video and predicted trajectories  
078 of an object and video captioning results compared with  
079 ground truth. In cases where videos lack target objects, our  
080 model adeptly captures the background region as a video  
081 feature, as also depicted in Fig. 2. Furthermore, we present  
082 inference results for Visual Question Answering (VQA)  
083 in Figs. 3 and 4. In addition to successful outcomes, we  
084 analyze error cases in trajectory prediction. We note that  
085 the model tends to struggle with complex video inputs, such  
086 as those involving scene changes or frames of poor quality,  
087 often leading to incorrect trajectory predictions. Lastly, we  
088 demonstrate zero-shot retrieval results for MSR-VTT video  
089 data through a demo video, with a sample showcased in  
090 Fig. 5. The demo video linked is below:

091  
092 <https://youtu.be/g0RkfwlMhcI>



Video description	Man putting cardboard into recycling bin on suburban street	
	Semantic Role Labeling	Semantic Groups
	[ARG0: Man] [V: putting] [ARG1: cardboard] [ARG2: into recycling bin on suburban street]	Man putting cardboard into recycling bin on suburban street recycling bin
Filtered description	Man putting cardboard into recycling bin on suburban street recycling bin	



Video description	Tired bearded man in a casual look sits on a couch, uses remote control, switches channels, nothing interests him, stands up and goes away. cozy living-room stuff on the background. having break	
	Semantic Role Labeling	Semantic Groups
	[ARG0: Tired bearded man in a casual look] [V: sits] [ARG1: on a couch] ,uses remote control , switches channels , nothing interests him , stands up and goes away . cozy living - room stuff on the background . [ARG0: Tired bearded man in a casual look] sits on a couch , [V: uses] [ARG1: remote control] , switches channels , nothing interests him , stands up and goes away . cozy living-room stuff on the background. having break	Tired bearded man in a casual look sits on a couch  Tired bearded man in a casual look uses remote control  ⋮
Filtered description	Tired bearded man in a casual look sits on a couch and Tired bearded man in a casual look uses remote control ⋮ and Tired bearded man in a casual look having break	

Figure 1. Example of SRL pre-processing on WebVid-2M datasets.



GT : A man kicking away an animal  
 SRTube : A man is kicking at a goat



GT : A man in a green shirt stands in a parking lot  
 SRTube : A guy is talking in a parking lot



GT : A sea creature is emerging from the water  
 SRTube : A beautiful ocean



GT : A man is driving a car for review  
 SRTube : A red car is driving

Figure 2. Visualization of the MSR-VTT video captioning results. We also show tube trajectory with red bounding box for  $i_t h$  tube query.



Question : What is a man doing?

GT : fall Baseline : **jump** Ours : **fall**



Question : What is the little girl doing in front of a woman?

GT : swing Baseline : **stand** Ours : **swing**

Figure 3. Visualization of the VQA results on MSVD. We show tube trajectory with red bounding boxes for tube query. The baseline is the same setting as the ablation study on features.



Question : What are people doing?

GT : throw    Baseline : **exercise**    Ours : **stand**



Question : What are two people doing beside a tree ?

GT : train    Baseline : **run**    Ours : **kick**

Figure 4. Error cases of VQA on MSVD dataset. We show tube trajectory with red bounding boxes for tube query. The baseline is the same setting as the ablation study on features.



Figure 5. Zero-shot text-to-video retrieval results on MSR-VTT. We type text and find most related video. We present demo video and this is sample of demo.

093

**References**

094

095

096

097

098

099

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

- [1] Rohrbach Anna, Rohrbach Marcus, Tandon Niket, and Schiele Bernt. A dataset for movie description. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3202–3212, 2015. 2
- [2] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021. 1
- [3] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *ICML*, page 4, 2021. 1
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Eur. Conf. Comput. Vis.*, pages 213–229. Springer, 2020. 1
- [5] Feng Cheng, Xizi Wang, Jie Lei, David Crandall, Mohit Bansal, and Gedas Bertasius. Vindlu: A recipe for effective video-and-language pretraining. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 10739–10750, 2023. 1
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 1
- [7] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska and Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The” something something” video database for learning and evaluating visual common sense. In *Proceedings of the IEEE international conference on computer vision*, pages 5842–5850, 2017. 2
- [8] Chunhui Gu, Chen Sun, David A Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sudheendra Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, et al. Ava: A video dataset of spatio-temporally localized atomic visual actions. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 6047–6056, 2018. 1
- [9] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. Localizing moments in video with natural language. In *Int. Conf. Comput. Vis.*, pages 5803–5812, 2017. 2
- [10] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 2
- [11] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Niebles Carlos Juan. Dense-captioning events in videos. In *Int. Conf. Comput. Vis.*, pages 706–715, 2017. 2
- [12] Jie Lei, Licheng Yu, Mohit Bansal, and Tamara L Berg. Tvqa: Localized, compositional video question answering. *arXiv preprint arXiv:1809.01696*, 2018. 2
- [13] Jie Lei, Tamara L Berg, and Mohit Bansal. Revealing single frame bias for video-and-language learning. *arXiv preprint arXiv:2206.03428*, 2022. 2
- [14] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3202–3211, 2022. 1

- [15] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252, 2015. 1 149 150 151 152 153
- [16] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *Int. Conf. Comput. Vis.*, pages 4489–4497, 2015. 1 154 155 156 157
- [17] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5288–5296, 2016. 2 158 159 160 161