

- nna

Supplementary Material **Global-and-Local Relative Position Embedding** for Unsupervised Video Summarization

Anonymous ECCV submission

Paper ID 4920

Baseline Architecture

Firstly, we elaborate on the baseline architectures such as VAE-GAN and CSNet. The base networks generally sample input videos by 2 fps, and it becomes a sequence of T images with a size of $H \times W$. Note that the T is determined by the length of the video ranging from 200 frames to 1000 frames when using 2 fps. We extract the 1024 sized feature of every single image by GoogLeNet pre-trained on ImageNet. The features are stacked along the temporal dimension and become $T \times 1024$. Next, the feature is embedded to $T \times 256$ using a fully connected (FC) layer followed by the bidirectional LSTM with a hidden size of 256. The output features from both directions are concatenated to $T \times 512$ and reduced to $T \times 1$ using FC and sigmoid function to predict a score per image.

The VAE-GAN framework enables the network to learn a video summariza-tion without ground truth labels. The main idea of unsupervised video sum-marization is generating original features of a video by giving summarized fea-tures through the reconstruction network. The discriminator then distinguishes whether it is from the original feature or the generated feature. If the generator performs enough well to deceive the discriminator, the selected features might have represented the video content successfully, and these frames are considered as kevframes.

The CSNet tackles with mode collapse of the unsupervised framework (i.e., VAE-GAN) and enables more discriminative feature learning. In addition to the feature reconstructing framework of VAE-GAN, they give a constraint that penalizes the low discrepancy of output scores by maximizing a variance of the scores. To ease the temporal dependency problem, they split the features into several temporal segments. Also, a temporal difference in the feature space is utilized to highlight fast scene changes.

In-depth Analysis of Network Design

In this section, we present the experimental results which are detailed but less essential to write on the main paper. It can be regarded as a design choice for our global-and-local relative position embedding (GL-RPE).

Table 1: The experiment by varying the feature dimension of RPE. The TVSum dataset is used in this table.

2.1 Relational matrix in self-attention embedding

We discover the optimal value of the feature dimension in the relative position embedding in this experiment. As shown in Table 1, the metrics using rank order statistics are measured by varying the feature dimension. The size of 64 shows the best result. Note that $RPE + GL_8$ is used in this experiment.



Fig. 1: (a) The output of LSTM is passed through the GL-RPE in a serial manner.(b) The output of LSTM is added to the output of GL-RPE in a parallel manner.

2.2 Parallel relative position embedding

In the main paper, the feature obtained from long short-term memory (LSTM) is passed to our GL-RPE in a serial manner as shown in Fig. 1-(a). We also tested the parallel case of adding each output of LSTM and RPE as illustrated in Fig. 1-(b), and the corresponding results are shown in Table 2. The serial case shows a better result and is adopted to our main submission.

3 Visualization for Keyframes

In our main submission, only a few selected frames are visualized. Here, we
present a dense visualization of the predicted summary. As shown in Fig. 2, the
representative frames are well selected.

Method	Kendall's τ	Spearman's ρ	Method	Kendall's τ	Spearman's ρ
$RPE_s + GR$	0.039	0.051	$RPE_p + GL_2$	0.039	0.052
$RPE_s + GI$	4 0.058	0.076	$RPE_p + GL_4$	0.042	0.056
$RPE_s + GI$	6 0.062	0.081	$\operatorname{RPE}_p + GL_6$	0.043	0.056
$RPE_s + GI$	0.070	0.091	$\operatorname{RPE}_p + GL_8$	0.047	0.062
	(a) Serial			(b) Parallel	
able 2: Res	ults for the p	arallel RPE. Th	ne TVSum data	aset is used	l in this tab
	VO TWO	TWO TWO			
AVD AN	LORIDHOVA TWOMHEEL DRIDDOVA	The distinguistics			
			A COMPANY A		
· 🍂 - 🍂					
	AL DEAL	The State of			
			Contraction of the other distance		
					IPARTA CALL AND A THE A
		- 63° 5 - 620°	NA ZA		
	STORE OF STREET				
		THE REPORT OF THE PARTY OF THE		1100 Participanti (1100	
			EXAN ESAN		
		Contraction of the			
ARA	× 1 1				
i 🏊 i 🛃 i 🔊					
NE ON DE	CN FON	TRAN AND A		-Carlo	to the
ALL IN					
Phone of			0 0	9-0-11	
	ON LOVON				
** ***				\sim · · ·	
The second se			00 000		0-0-0
	AL AL				
" In the second second	Remark and Income of the			-4	
				There is the second second	0 - OP
ZOK					
		A CONTRACT			
		Andrew Contraction of Contraction			
			1		
	doo in domast	r compled and	Iroutromod aro	highlighto	d with blue

Fig. 2: A video is densely sampled, and keyframes are highlighted with blue.