

Video Summarization with a **Dual Attention Capsule Network**

Hao Fu¹, Hongxing Wang¹, Jianyu Yang²



OOCHOW UNIVERSITY

¹Chongqing University, China ²Soochow University, China

Contributions at a Glance

Generating a compact and non-redundant summary for a given video without missing significant information.



- We propose a novel dual attention capsule network model, which can effectively incorporate the short- and long-term temporal dependencies among video frames for summarization.
- Our proposed video summarization is parallelizable, which can easily handle longer-term dependencies among video frames than the RNN/LSTM-based approaches.
- Experimental results show that our proposed method owns stronger learning ability, and is competitive with existing state-of-the-art methods.

Method

Problem

Feature Extraction

- Extractor: GoogLeNet
- Divided frames into M clips
- Each clip is viewed as a local block

Dual Attention Feature Refinement

Self-attention mechanism

- Learn the short- and long-term dependencies within and between clips
- Two-Stream Capsule Network Learning
- Later fusion mechanism
- Scalar \rightarrow Vector
- Margin loss



Experimental Results

Ablation study of our method

τ and ρ results compared with SOTA

	-								
Method	SumMe	TVSum	Dataset	SumMe		SumMe		TVSum	
Ours-local	45.2	58.3	Metric	Kendall's $ au$	Spearman's $ ho$	Kendall's $ au$	Spearman's $ ho$		
Ours-global	45.4	58.5	Random	0.000	0.000	0.000	0.000		
Ours-fc	46.6	58.7	DR-DSN _{sup} [7]	0.034	0.041	0.025	0.039		
Ours-Ic	40.0	50.7	dppLSTM [5]	0.040	0.049	0.042	0.055		
Ours	47.5	59.4	Ours	0.063	0.059	0.058	0.065		
			dppLSTM [5] 0.040 0.049 0.042 0.055 Ours 0.063 0.059 0.058 0.065						

An example summary generated by our approach on the "Bearpark climbing" video



F-score(%) results compared with SOTA

Method	SumMe			TVSum		
-	С	А	Т	С	А	Т
dppLSTM [5]	38.6	41.6	40.7	54.2	57.9	56.9
SUM-GAN _{sup} [6]	41.7	43.6	-	56.3	61.2	-
DR-DSN _{sup} [7]	42.1	43.9	42.6	58.1	59.8	58.9
SASUP [35]	45.3	-	-	58.2	-	-
CSNet _{sup} [21]	48.6	48.7	44.1	58.5	57.1	57.4
Ours	47.5	49.3	45.2	59.4	59.8	59.2

References

- K. Zhang, W. Chao, F. Sha, and K. Grauman, Video
- summarization with long short-term memory, ECCV'16. B. Mahasseni, M. Lam, and S. Todorovic, Unsupervised video summarization with adversarial LSTM networks, CVPR'17.
- K. Zhou, Y. Qiao, and T. Xiang, Deep reinforcement learning for unsupervised video summarization with diversityrepresentativeness reward, AAAI'18.
 - C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, Going deeper with convolutions, CVPR'15
- Y. Jung, D. Cho, D. Kim, S. Woo, and I. Kweon, Discriminative feature learning for unsupervised video summarization, AAAI'19.
- H. Wei, B. Ni, Y. Yan, H. Yu, X. Yang, and C. Yao, Video summarization via semantic attended networks, AAAI'18.