



# Video Summarization with a Dual Attention Capsule Network

Hao Fu<sup>1</sup>, Hongxing Wang<sup>1</sup>, Jianyu Yang<sup>2</sup>

<sup>1</sup>Chongqing University, China

<sup>2</sup>Soochow University, China



重慶大學  
CHONGQING UNIVERSITY



蘇州大學  
SOOCHOW UNIVERSITY

# Why Need Video Summarization?

---



# Properties of a Good Summary

---

- Representative
- Diverse
- Coherent/has a sense of story



Case 2



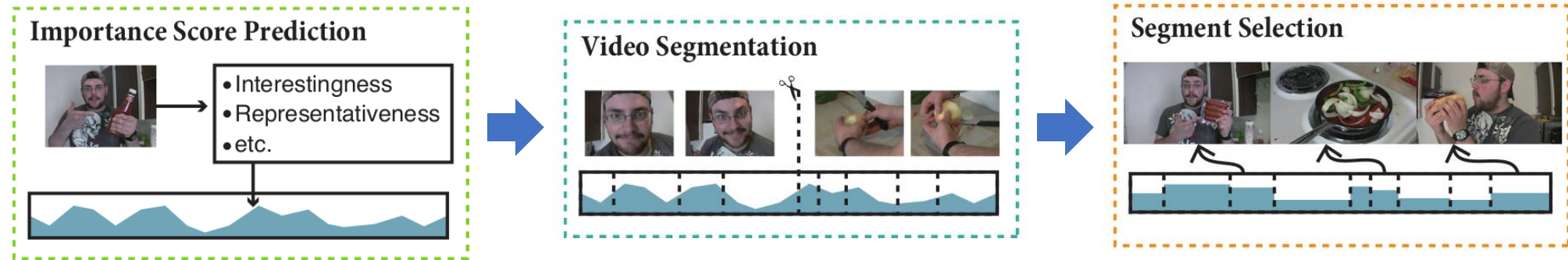
Case 1



Case 3



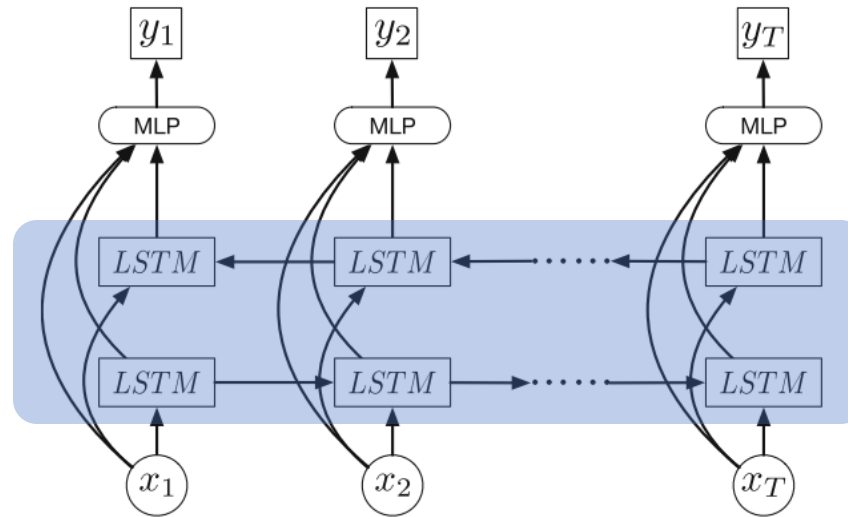
# Common Video Summarization Pipeline



Modified From: Otani et al., CVPR 2019

- Gygli et al., Creating summaries from user videos, ECCV 2014.
- Song et al., TVSum: Summarizing web videos using titles, CVPR 2015
- ...

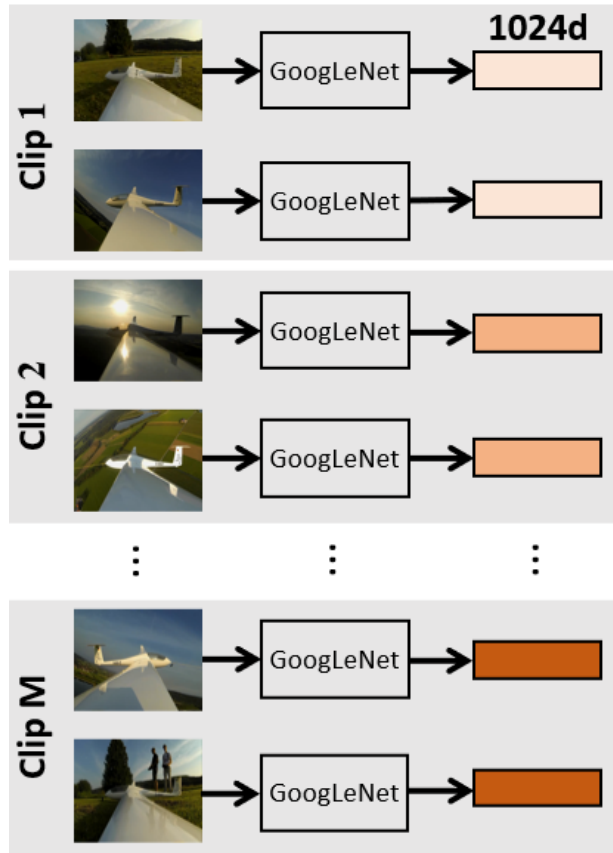
# Popular Solution and Existing Challenges



Modified From: Zhang et al., ECCV 2016

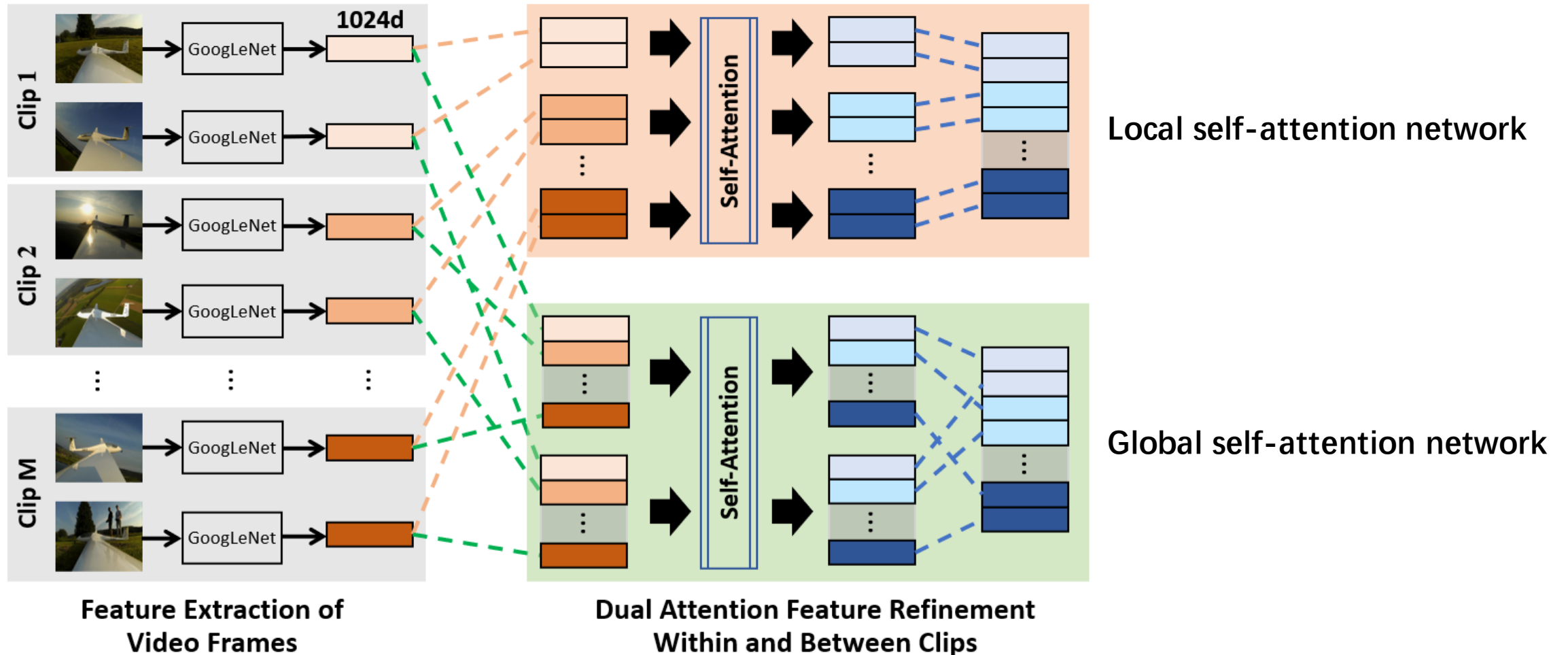
- Unfriendliness for parallel processing
- Considerable computational burdens
- Not capable enough to capture long-term dependencies

# Proposed Method

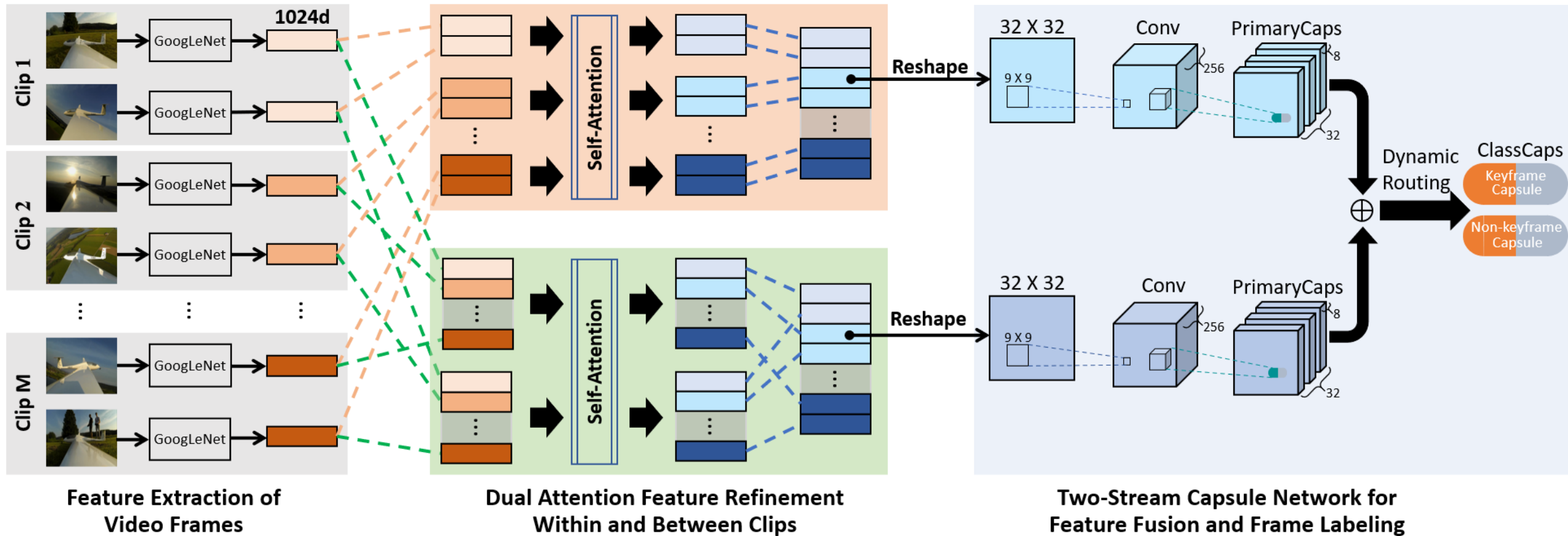


**Feature Extraction of  
Video Frames**

# Proposed Method



# Proposed Method





# Experiments

---

## Datasets

- SumMe: 25 videos with various events such as travel, cooking, and sports
- TVSum: 50 videos with various genres such as news, how-to, and egocentric

## Evaluation metrics

- F-score
- Kendall's  $\tau$
- Spearman's  $\rho$

# Comparison with SOTA (F-score)

---

Method	SumMe			TVSum		
	C	A	T	C	A	T
-						
dppLSTM [5]	38.6	41.6	40.7	54.2	57.9	56.9
SUM-GAN <sub>sup</sub> [6]	41.7	43.6	-	56.3	<b>61.2</b>	-
DR-DSN <sub>sup</sub> [7]	42.1	43.9	42.6	58.1	59.8	58.9
SASUP [35]	45.3	-	-	58.2	-	-
CSNet <sub>sup</sub> [21]	<b>48.6</b>	48.7	44.1	58.5	57.1	57.4
Ours	47.5	<b>49.3</b>	<b>45.2</b>	<b>59.4</b>	59.8	<b>59.2</b>

# Comparison with SOTA ( $\tau$ and $\rho$ )

---

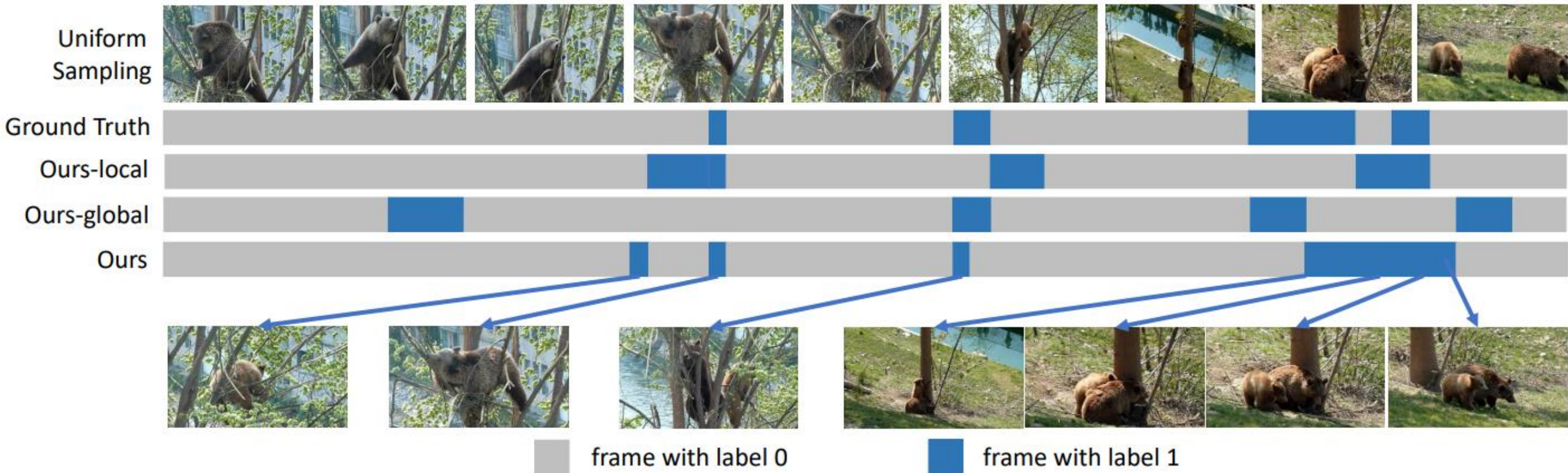
Dataset	SumMe		TVSum	
Metric	Kendall's $\tau$	Spearman's $\rho$	Kendall's $\tau$	Spearman's $\rho$
Random	0.000	0.000	0.000	0.000
DR-DSN <sub>sup</sub> [7]	0.034	0.041	0.025	0.039
dppLSTM [5]	0.040	0.049	0.042	0.055
Ours	<b>0.063</b>	<b>0.059</b>	<b>0.058</b>	<b>0.065</b>

# Results under Unbalanced Train-Test Duration

---

Dataset	SumMe		TVSum	
Metric	F-score	Kendall's $\tau$	F-score	Kendall's $\tau$
dppLSTM [5]	36.6(-2.0)	0.036(-0.004)	51.5(-3.2)	0.032(-0.010)
DR-DSN <sub>sup</sub> [7]	39.7(-2.4)	0.029(-0.005)	56.2(-1.9)	0.019(-0.006)
Ours	<b>46.1(-1.4)</b>	<b>0.076(+0.013)</b>	<b>59.0(-0.4)</b>	<b>0.056(-0.002)</b>

# Qualitative Results





# Ablation Study

---

Method	SumMe	TVSum
Ours-local	45.2	58.3
Ours-global	45.4	58.5
Ours-fc	46.6	58.7
Ours	<b>47.5</b>	<b>59.4</b>

# Conclusions

---

- 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.

# Thanks

**20151757@cqu.edu.cn**