

Hybrid Improvements in Multimodal Analysis for Deep Video Understanding

BeibeiZhang¹, Fan Yu,^{1,2}Yaqun Fang,¹Tongwei Ren^{1,2,*}

¹State Key Laboratory for Novel Software Technology, Nanjing University

²Shenzhen Research Institute of Nanjing University





Introduction



• Deep video understanding (DVU)

- requires systems to develop a deep analysis and understanding of long video.
- use known information to reason about other, more hidden information, and to populate a knowledge graph (KG) with all acquired information.

• HLVU dataset

- 14 videos
 - 10 for development
 - 4 for test
- 1h/video in average
- shot, entity name, entity type, screenshots

Training dataset:

Honey - Romance - 86 mins.
Let's bring back Sophie - Drama - 50 mins.
Nuclear Family - Drama - 28 mins.
Shooters - Drama - 41 mins.
Spiritual Contact The Movie - Fantasy - 66 mins.
Super Hero - Fantasy - 18 mins.
The Adventures of Huckleberry Finn - Adventure - 106 mins.
The Big Something - Comedy - 101 mins.
Time Expired - Comedy / Drama - 92 mins.
Valkaama - Adventure - 93 mins.

Testing dataset:

1- Bagman - Drama / Thriller - 107 mins.

2- Manos - Horror - 73 mins.

3- Road to Bali - Comedy / Musical - 90 mins.

4- The Illusionist - Adventure / Drama - 109 mins.

Task

scene-level

- Find the unique scene.
- Fill in the graph space.
- Find next interaction in scene X between person Y and person Z.



movie-level

- Find all possible paths question.
- Fill in the part of graph question.
- Multiple choice questions.
- Find previous interaction in scene X between person Y and person Z.
- Find the 1-to-1 relationship between scenes and natural language descriptions.
- Classify scene sentiment from a given scene.





Video segmentation





- Video feature
 - C3D
- Audio feature
 - MFCC, LMFE
- Text feature
 - BERT

- Entity feature (subject, object, union)
 - CenterTrack
 - InsightFace
 - C3D





Joint learning architecture

- relationship: average of medium video feature
- interaction: medium video feature + average feature



- Low-shot, Zero-shot learning
- Joint learning

$$l = (1 - \cos(\beta, \gamma))^2 + \frac{1}{n} \sum_{i \in U} (\cos(\beta, \mu_i) + 1)^2$$

$$L = l_R + \frac{1}{n} \sum \left(l_I + l_S \right)$$



- *l* denotes loss
- denotes the feature of pair
- denotes the feature of the positive relationship
- denotes the set of negative relationships
- denotes the feature of relationship
- denotes the number of negative relationships
- denotes the total loss
- l_R denotes the loss of relationship
- *l*_{*l*} denotes the loss of interaction
- *l*_sdenotes the loss of sentiment

Query answering

movie-level

- Find all possible paths question.
- Fill in the part of graph question.
- Multiple choice questions.
 - relationship knowledge graph

scene-level

- Find the unique scene.
- Fill in the graph space.
 - interaction knowledge graph
- Find next/previous interaction in scene X between person Y and person Z.
 - split medium video into shot videos
- Find the 1-to-1 relationship between scenes and natural language descriptions
 - match with predicted interactions and sentiments.
- Classify scene sentiment from a given scene.
 - sentiment model





Improvement



Sentiment score



• Sentiment distribution

 $l_{+} = (1 - \cos(\beta, f_{+}))^{2} \cdot d_{+}$ $l_{-} = (\cos(\beta, f_{-}) + 1)^{2} \cdot d_{-}$

• Experiment result

- + denotes the positive loss
- denotes the feature of scene,
- + denotes the feature of the positive sentiment
- + denotes the distribution ratio of positive sentiment
- – denotes the negative loss
- – denotes the feature of negative sentiment
- – denotes the distribution ratio of negative sentiment

	cos_sim	emotion_score	senti_distribution		
Recall _k	27.5	29.4	29.4		



Model	C3D	I3D _{rgb}	I3D _{flow}	I3D _{avg}	I3D _{con}
Recall@50	28.8	29.1	33.3	32.8	29.6

- I3D represents rgb stream of I3D
- I3D represents optical flow stream of I3D
- I3D represents avaraging the outputs of rgb stream and optical flow stream
- I3D represents concatenating the features extracted by both of the streams

Improvement



Method	model	mat _I	mat_E	mat _{I+S}	mat_{E+O}	mat_{I+S+E}
Recall@k	4.9	4.9	12.1	5.3	12.6	6.8

- Model means description matching model
- mat means direct matching algorithm
- represents predicted interactions
- represents detected entities
- represents predicted sentiments
- represents detected objects.



THANK YOU



