

# endAffectNet: Self-Attention based Networks for Predicting Affective Responses from Movies

Thi Phuong Thao\*, Balamurali B.T.\* Dorien Herremans\* and Gemma Roig<sup>†</sup> (thiphuongthao\_ha@mymail.sutd.edu.sg)

Singapore University of Technology and Design (SUTD), <sup>†</sup>Goethe University Frankfurt, Germany

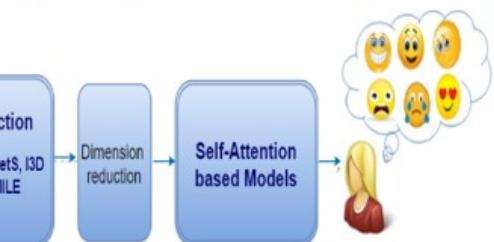
## Motivation

clips) still remains a challenge

looked in a person is hard for computers.

Affective responses of viewers from movies use low-level approaches

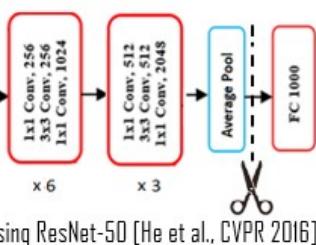
## Approach



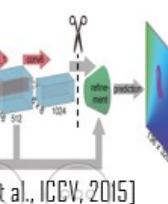
Approach for predicting emotions from movies.

Model valence and arousal separately

Attention mechanism (Vaswani et al, 2017)



Using ResNet-50 [He et al., CVPR 2016]



[a., ICCV, 2015]

## Feature Extraction

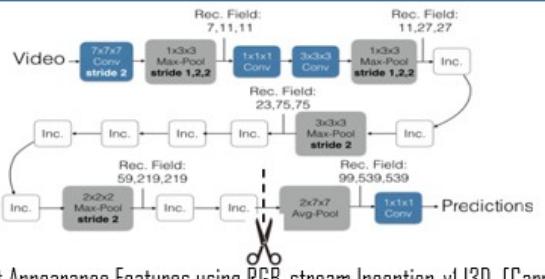


Figure 3: Extract Appearance Features using RGB-stream Inception-v1 I3D [Carreira et al., CVPR 2017]

### • Audio

- OpenSMILE: "emobase2010.conf" (INTERSPEECH 2010): window size = 320ms, hop size = 40ms
- VGGish:

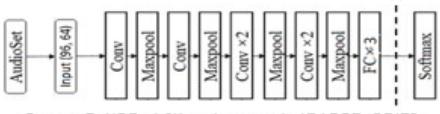
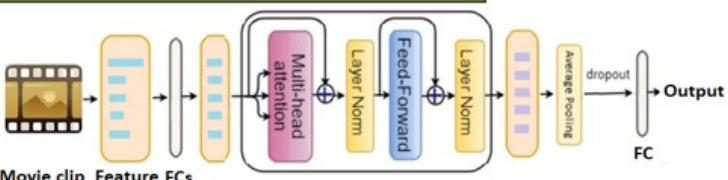


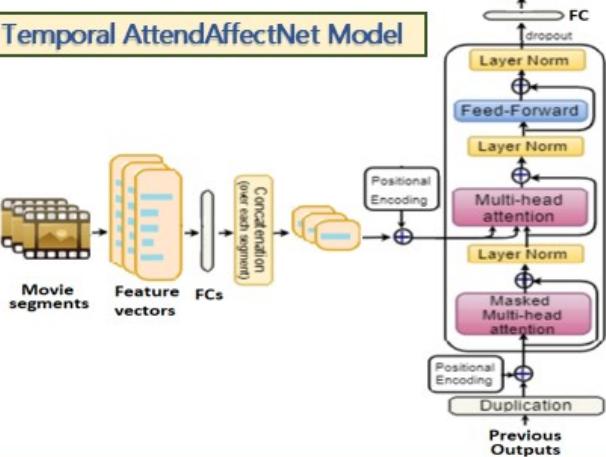
Figure 5: VGGish[Hershey et al., ICASSP, 2017]

## Proposed Models

### • Figure 6: Feature AttendAffectNet Model



### • Figure 7: Temporal AttendAffectNet Model



## Experimental Results

TABLE 1: ACCURACY OF THE PROPOSED MODELS ON THE EXTENDED COGNIMUSE DATASET.

Models	Arousal		Valence	
	MSE	PCC	MSE	PCC
Feature AAN (only video)	0.152	0.518	0.204	0.483
Feature AAN (only audio)	0.125	0.621	0.185	0.543
<b>Feature AAN (video and audio)</b>	<b>0.124</b>	<b>0.630</b>	<b>0.178</b>	<b>0.572</b>
Temporal AAN (only video)	0.178	0.457	0.267	0.232
Temporal AAN (only audio)	0.162	0.472	0.247	0.254
<b>Temporal AAN (video and audio)</b>	<b>0.153</b>	<b>0.551</b>	<b>0.238</b>	<b>0.319</b>
Sivaprasad et al. [23]	<b>0.08</b>	<b>0.84</b>	<b>0.21</b>	<b>0.50</b>

TABLE 2: ACCURACY OF THE PROPOSED MODELS ON THE STANFORD COGNIMUSE DATASET.

Models	Arousal	Valence
Feature AAN (only video)	0.152	0.518
Feature AAN (only audio)	0.125	0.621
<b>Feature AAN (video and audio)</b>	<b>0.124</b>	<b>0.630</b>
Temporal AAN (only video)	0.178	0.457
Temporal AAN (only audio)	0.162	0.472
<b>Temporal AAN (video and audio)</b>	<b>0.153</b>	<b>0.551</b>
Liu et al. [56]		
Chen et al. [55]		
Yi et al. [22]		
Yi et al. [41]		

## Visualization

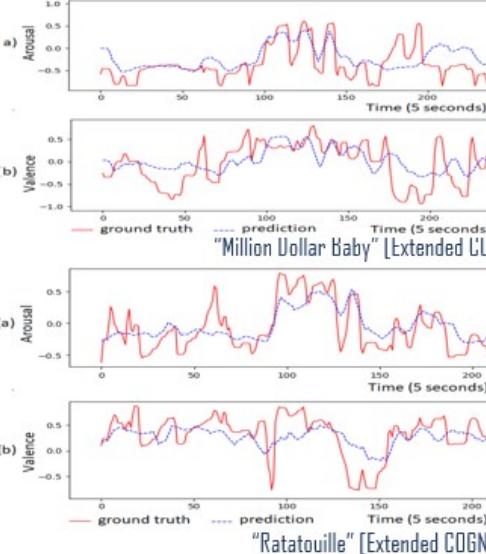


Figure 8: Visualization of the ground-truth and predicted affective responses over time for two movies.

## Conclusion

- Feature AttendAffectNet outperforms the Temporal AttendAffectNet.
- Audio-based model provides a higher performance than video-based model.
- Model combining all features (video, audio) reaches the highest accuracy.