



# MultiGap: Multi-Pooled Inception Network with Text Augmentation for Aesthetic Prediction of Photographs

Yong-Lian Hii, John See, Magzhan Kairanbay, Lai-Kuan Wong





Internet



high



Classifier

low



# Main contribution

1. Deep neural network architecture called MultiGAP that exploits features from multiple inception modules pooled by global average pooling (GAP), evaluated by prediction of (10-class) categorical distribution, before performing score binarization
2. The incorporation of textual features trained simultaneously with MultiGAP using a recurrent neural network (RNN) with gated recurrent unit (GRU) layers
3. The ability to leverage the GAP block for visualization of activation maps based on the aesthetic category and rating

# Related work

## Handcrafted low-level features

---

color, hue, saturation, light exposure, and also other heuristics driven by rule of thumbs used by professional photographers

## Generic features

---

SIFT , feature encoding method such as a Fisher Vector

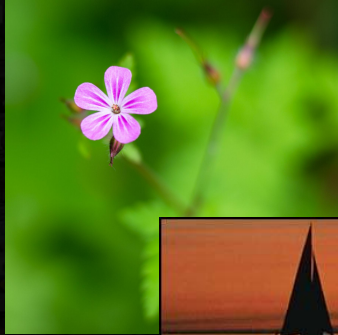
## Deep learning models

---

CNN, Double Column CNN, Multi-modal CNN



# Binary classification



High

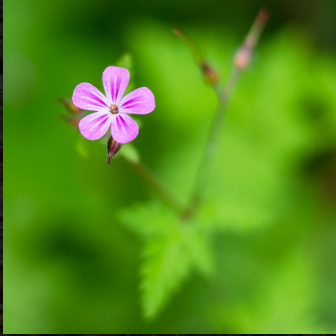
VS



Low



## Regression (Score)



6.7



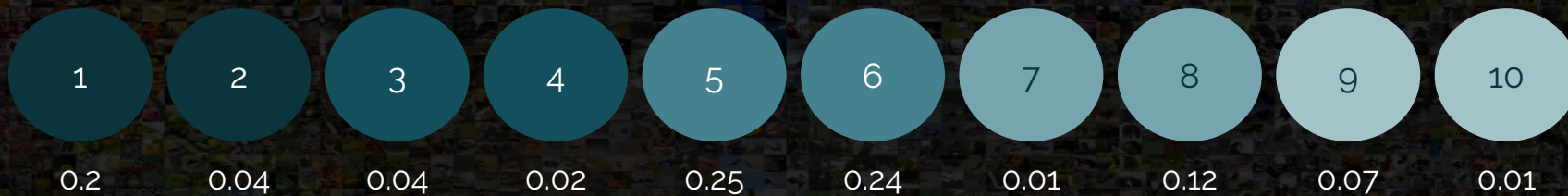
5.4



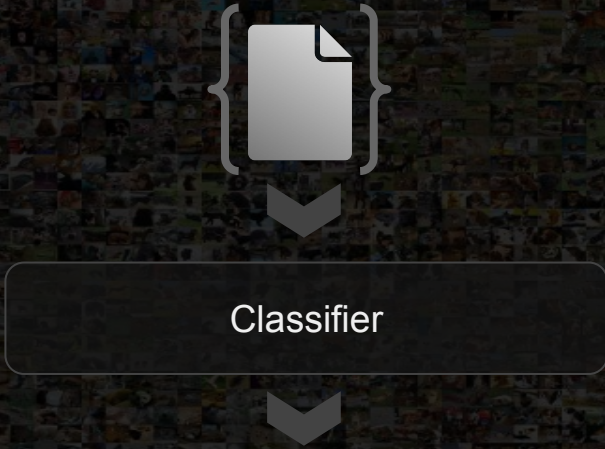
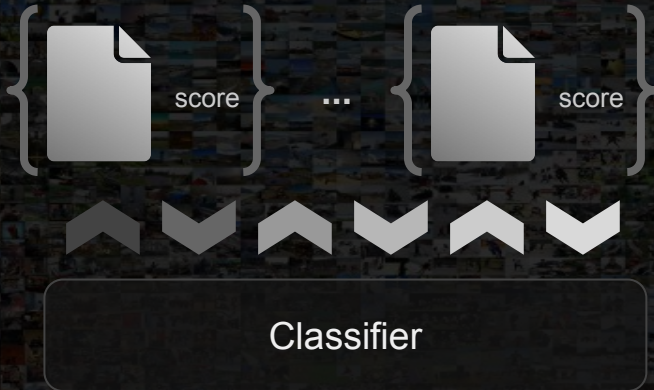
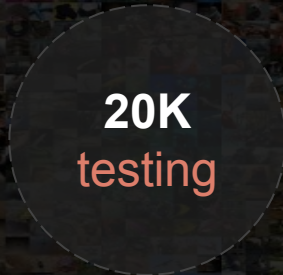
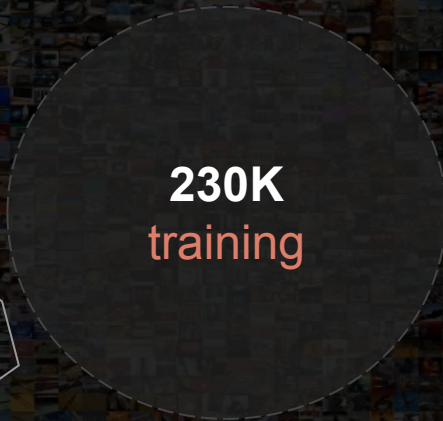
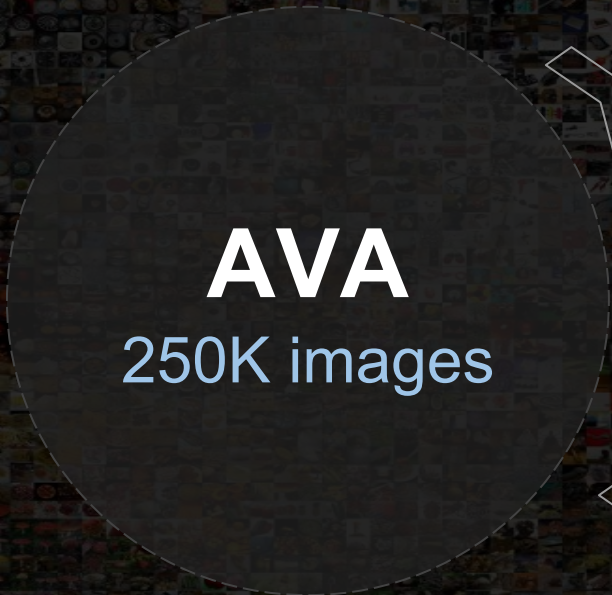
7.2



# Rating distribution







Label (aesthetically **high** or **low**)

2012, N. Murray



**AVA**

250K images

230K  
training

20K  
testing

$\delta$  used to filter out from noisy images



threshold = 5

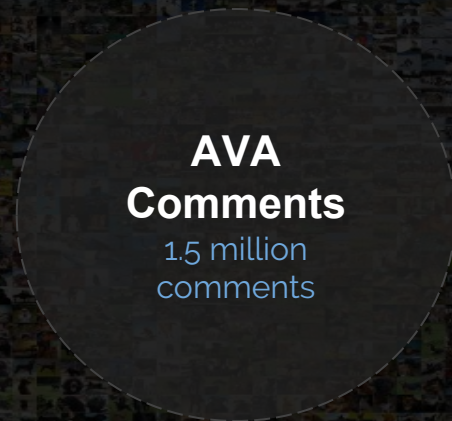
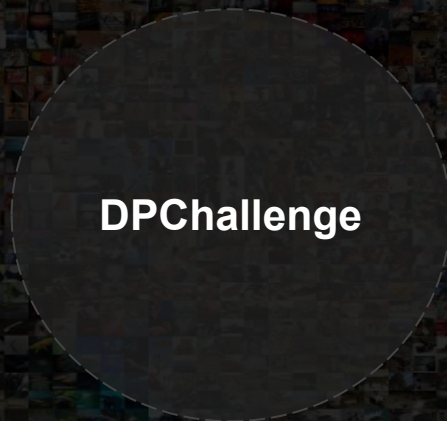
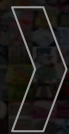
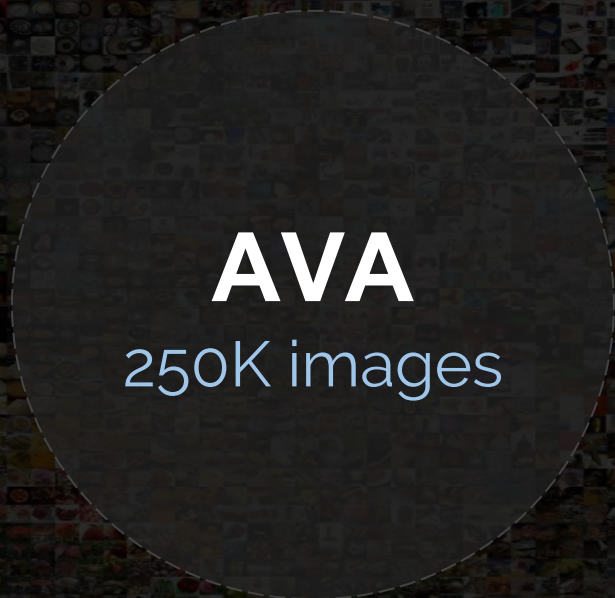
$\delta = \{2, 1.5, 1, 0.5, 0\}$



threshold = 5

$\delta = \{0\}$

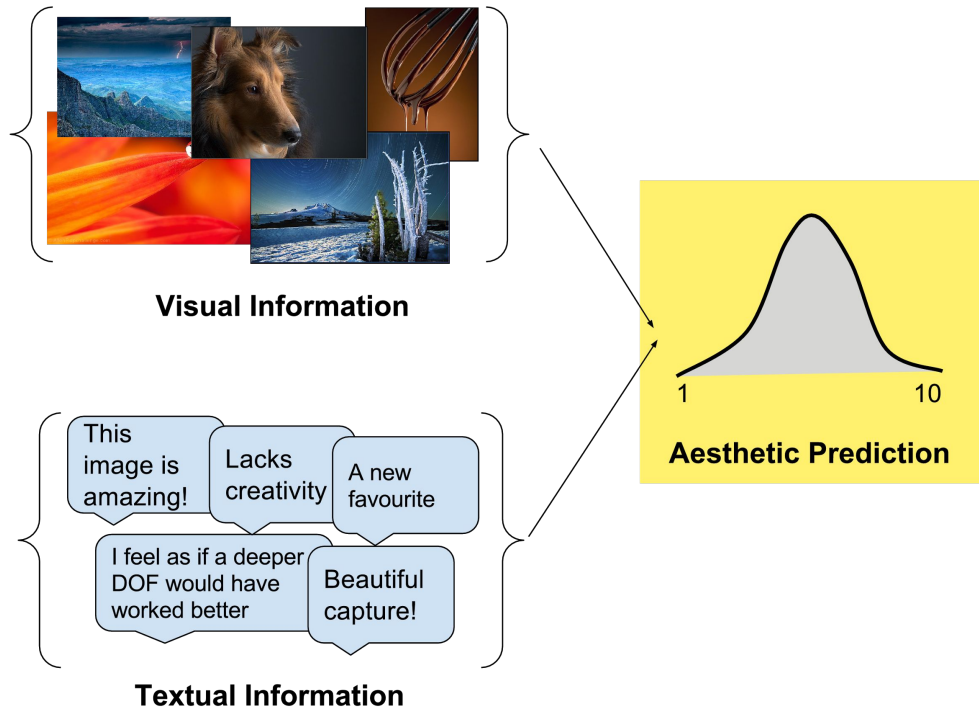




Zhou et al.



# Proposed method



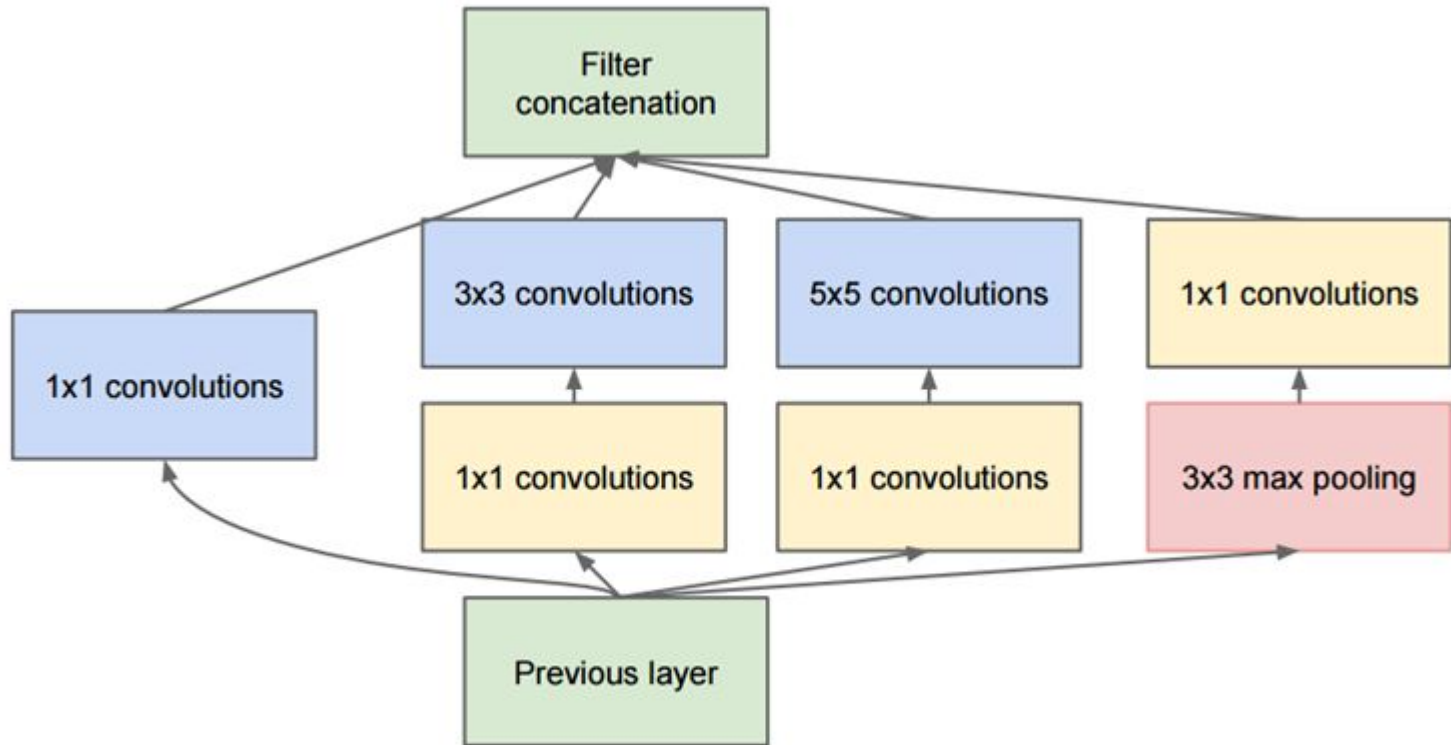






# Visual features

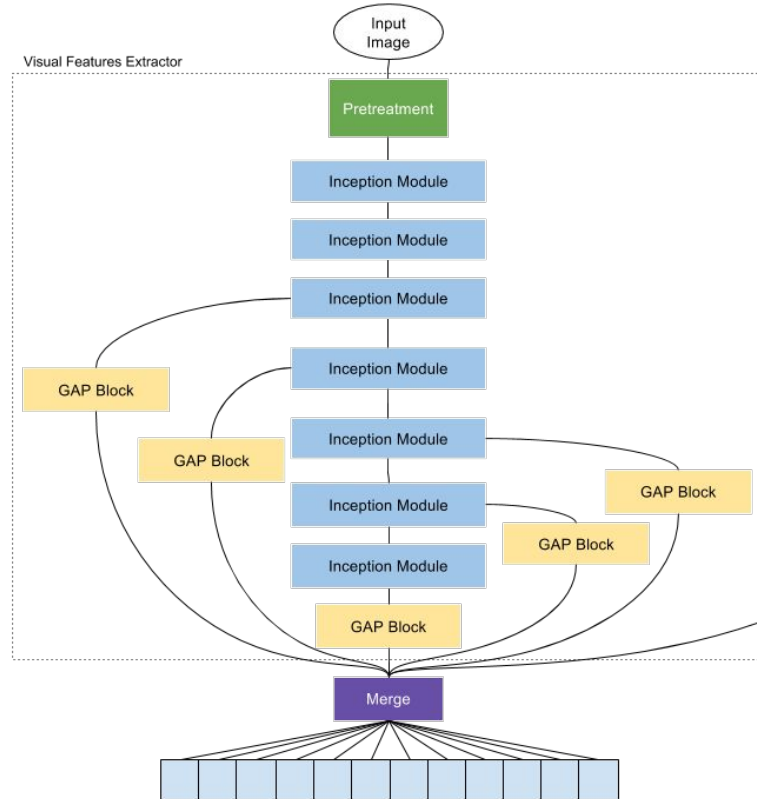
## Inception module



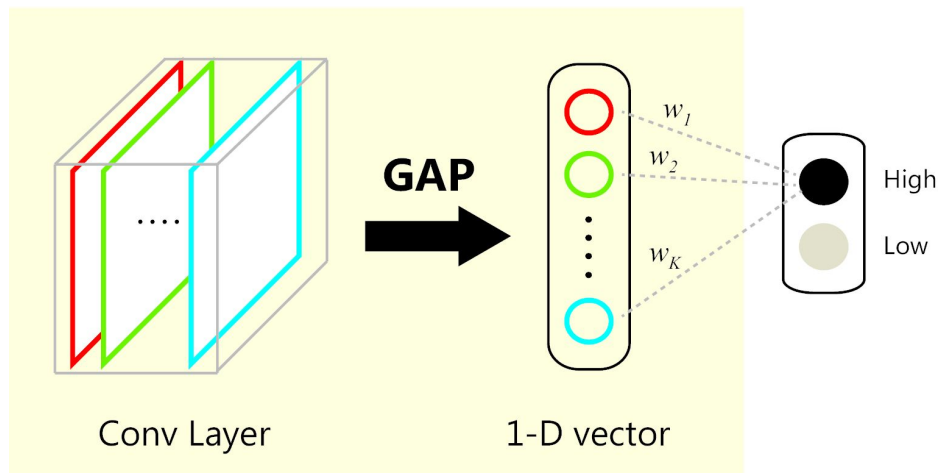


# Visual features

## Proposed method

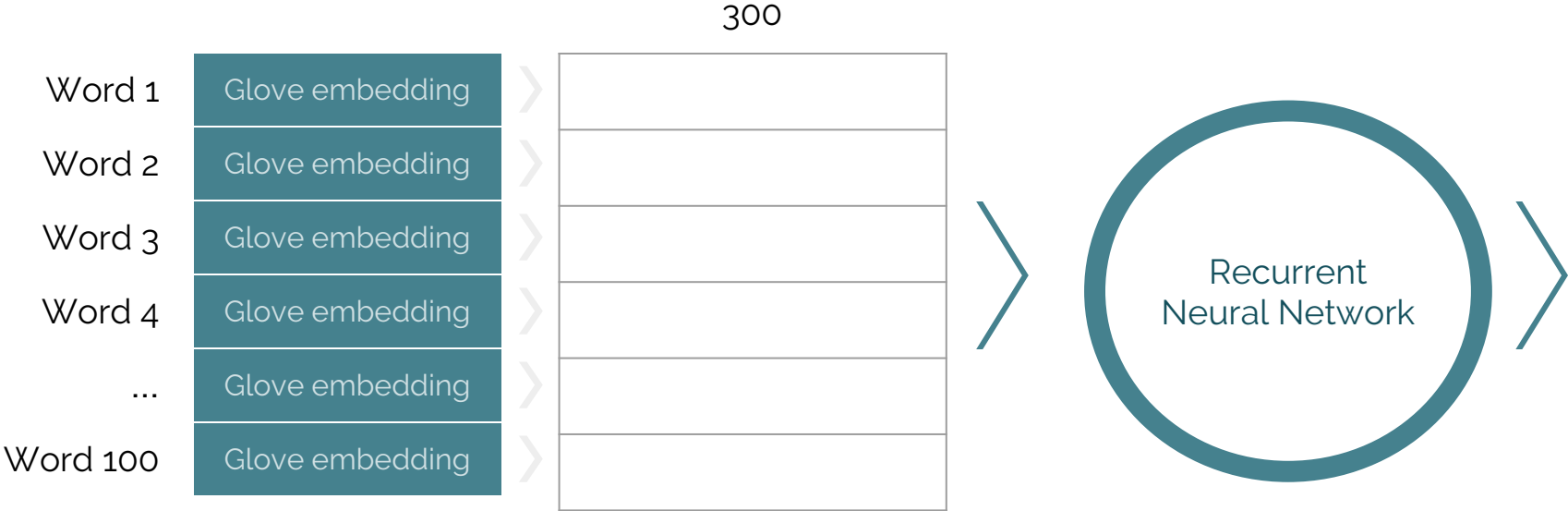


# GAP (Global Average Pooling) Layer

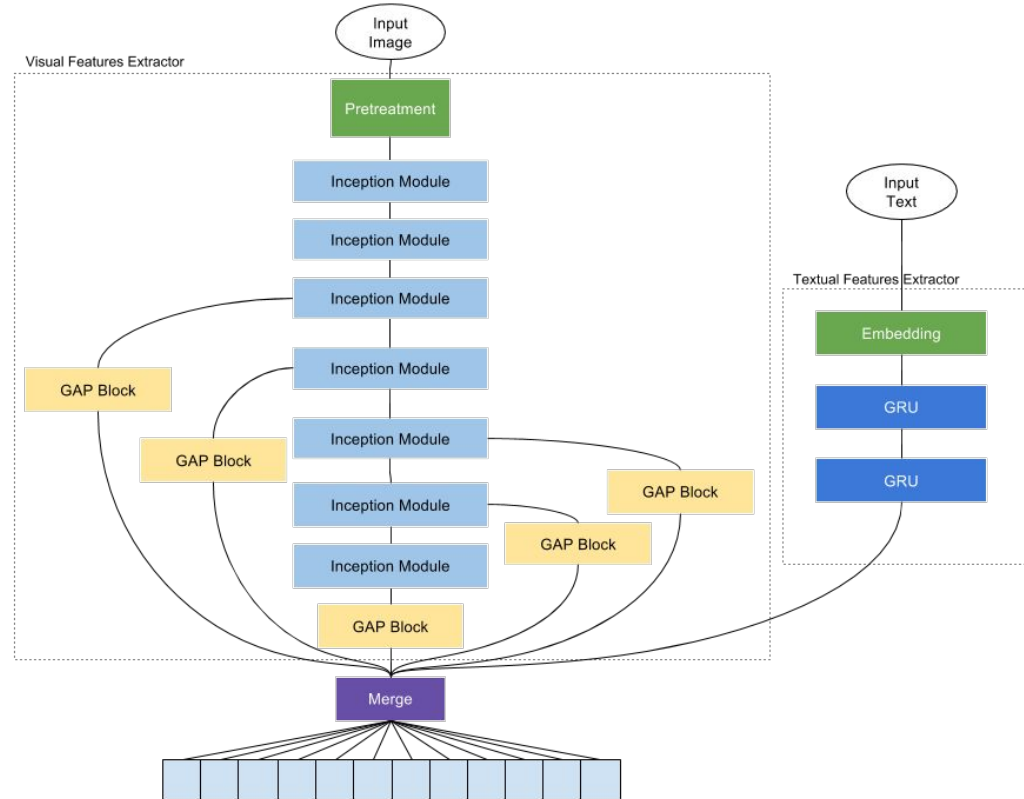




# Textual Features



# Proposed method

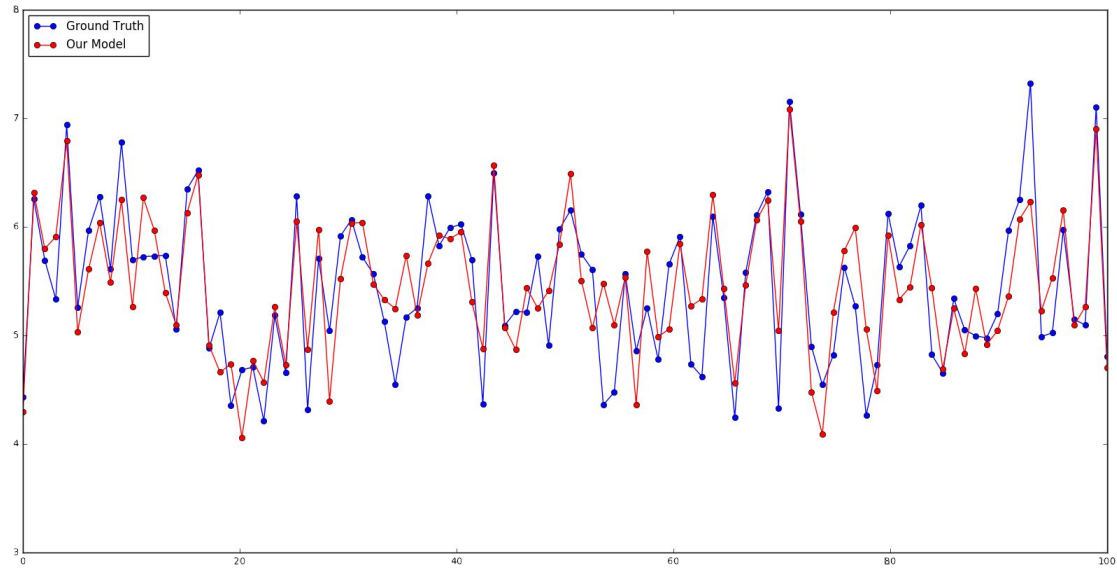




# Results

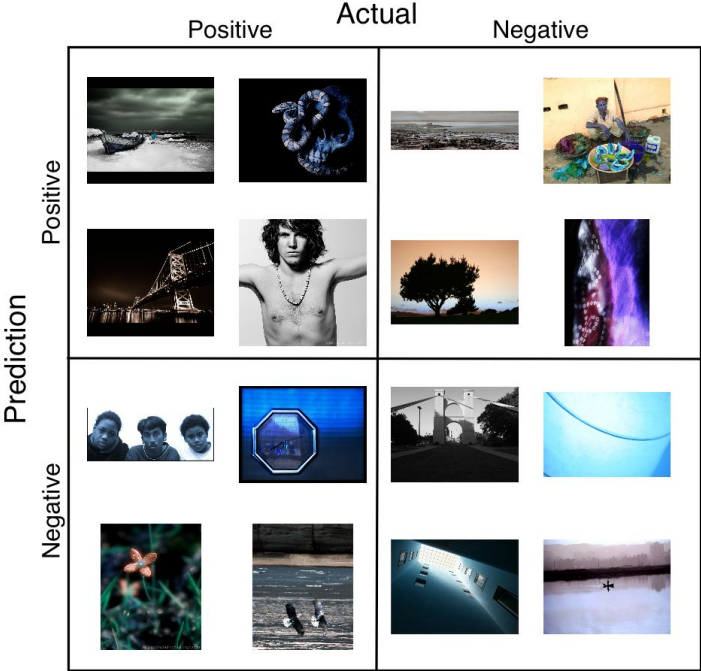
	Model	Accuracy
Image	DCNN [2]	73.25
	RDCNN [2]	74.46
	Kao et al. [12]	74.51
	AlexNet [10] – finetuned	75.11
	DMA [5]	75.41
	GoogLeNet [16] – finetuned	75.60
	MultiGAP	<b>75.76</b>
	SingleGAP	<b>76.31</b>
	BDN [13]	76.80
Text	word2vec [19]	78.40
	1D-CNN [20]	79.48
	Naive Bayes SVM [14]	80.90
	RNN (1-layer GRU)	<b>81.09</b>
	RNN (2-layer GRU)	<b>81.79</b>
Joint	Multimodal DBM [14]	78.88
	SingleGAP + RNN (2-layer GRU)	<b>80.54</b>
	MultiGAP + RNN (2-layer GRU)	<b>82.27</b>

# Results

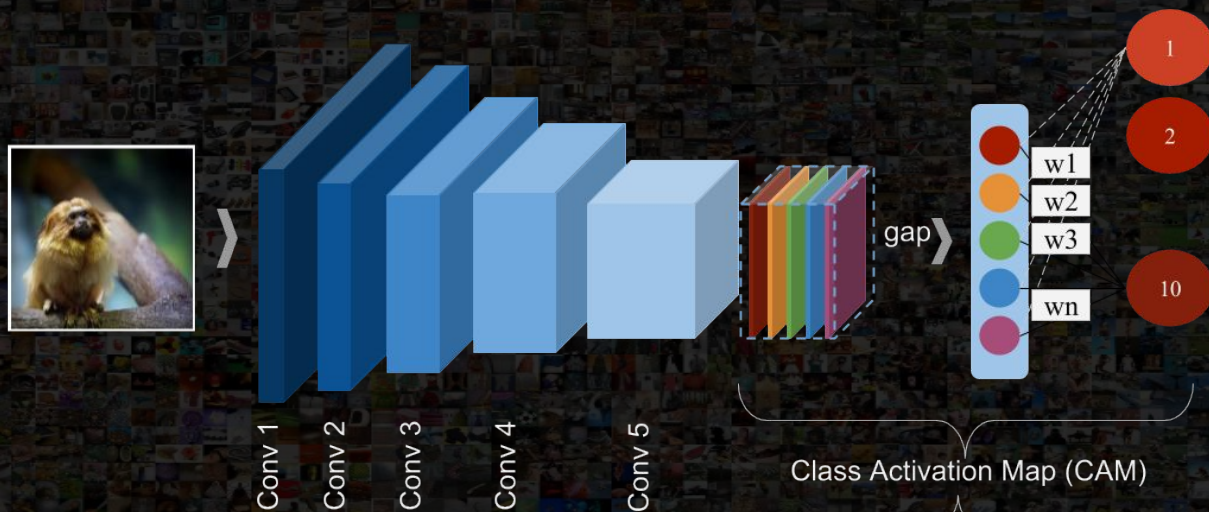




# Confusion matrix



# Class Activation Map (CAM)



$$w_1 * \text{[red box]} + w_2 * \text{[orange box]} + w_3 * \text{[green box]} + w_4 * \text{[blue box]} + \dots + w_n * \text{[pink box]} = \text{[CAM heatmap]}$$

# Class Activation Maps (CAM)



(a)

(b)



Thank you for  
your attention

Q&A

