

Adaptive Feature Fusion Network for Gaze Tracking in Mobile Tablets

Yiwei Bao¹, Yihua Cheng¹, Yunfei Liu¹, Feng Lu^{1, 2}

¹State Key Laboratory of Virtual Reality Technology and Systems, School of CSE, Beihang University, Beijing, China. ²Peng Cheng Laboratory, Shenzhen, China



Introduction Gaze Estimation



Facial image of the subject

2D gaze location on screen

Related Works



Limitation

- 1. Neglect the special similarity of two eye structures
- 2. Treat face and eye images separately.

Motivation

Two Eye Similarity





Right (flipped)

Structures like canthus, eyelid and outline of two eyes are almost identical.

Relationship between eye appearance and face





Eye appearance varies under different circumstances. Face appearance gives some clues (head pose, glasses, et al.) about how eyes would look like.

Four stream inputs: eye images, face image, face and eye bounding boxes locations (Rects)



Feature Fusion Scheme: eye feature fusion via stacking, attention weights and conv layers.



Adaptive Group Normalization: recalibrate eye features according to facial features





Experiments GazeCapture Dataset

Method	Phone error (cm)	Tablet error (cm)	3.5 3 iTracker [20] 3 SAGE [38] TAT [39] AFF-Net (ours) 2.81 2.72 2.66 2.30
iTracker [20]	1.86	2.81	2 1.86 1.78 1.77 1 62
SAGE [40]	1.78	2.72	1.5 I I I I I I I I I I I I I I I I I I I
TAT [41]	1.77	2.66	l l l l l l l l l l l l l l l l l l l
AFF-Net	1.62	2.30	·
			Phone Tablet

Output performs current SOTA methods on mobiles and tablets.

Experiments MPIIGaze Dataset



Also achieves SOTA performance on MPIIGaze dataset.

Experiments Ablation study—GazeCapture dataset

Method	GazeCapture		
Wiediod	Phone (cm)	Tablet (cm)	
AFF-Net	1.62	2.30	
without ST	1.67	2.39	
without SE	1.68	2.31	
without AdaGN	1.69	2.33	

Ablation study proves effectiveness of each module.

Experiments Ablation study—GazeCapture dataset



Further analysis shows that proposed modules reduce estimation errors. Especially on those **remote locations.**

Experiments Ablation study—GazeCapture dataset



Proposed network gives more accurate estimation when eyes and face are small in the image.



Thanks

Iufeng@buaa.edu.cnhttp://phi-ai.org

