

Believe It or Not, We Know What You Are Looking at!

Dongze Lian^{*[0000-0002-4947-0316]}, Zehao Yu^{*[0000-0002-6559-9830]}, and
Shenghua Gao^{†[0000-0003-1626-2040]}

School of Information Science and Technology, ShanghaiTech University
{liandz, yuzh, gaoshh}@shanghaitech.edu.cn

Abstract. By borrowing the wisdom of human in gaze following, we propose a two-stage solution for gaze point prediction of the target persons in a scene. Specifically, in the first stage, both head image and its position are fed into a gaze direction pathway to predict the gaze direction, and then multi-scale gaze direction fields are generated to characterize the distribution of gaze points without considering the scene contents. In the second stage, the multi-scale gaze direction fields are concatenated with the image contents and fed into a heatmap pathway for heatmap regression. There are two merits for our two-stage solution based gaze following: i) our solution mimics the behavior of human in gaze following, therefore it is more psychological plausible; ii) besides using heatmap to supervise the output of our network, we can also leverage gaze direction to facilitate the training of gaze direction pathway, therefore our network can be more robustly trained. Considering that existing gaze following dataset is annotated by the third-view persons, we build a video gaze following dataset, where the ground truth is annotated by the observers in the videos. Therefore it is more reliable. The evaluation with such a dataset reflects the capacity of different methods in real scenarios better. Extensive experiments on both datasets show that our method significantly outperforms existing methods, which validates the effectiveness of our solution for gaze following. Our dataset and codes are released in <https://github.com/svip-lab/GazeFollowing>.

Keywords: Gaze following · Saliency · Multi-scale gaze direction fields.

1 Introduction

Gaze following is a task of following other people’s gaze in a scene and inferring where they are looking [22]. It is important for understanding the behavior of human in human-human interaction and human-object interaction. For example, we can infer the intention of persons based on their gaze points in human-human interaction. In new retailing scenario, we can infer the interest of the consumers in different products based on their eyes contact with those products (as shown in

^{*}The authors contribute equally.

[†]Corresponding author.

Figure 1 (a) (b)), and infer what kind of information (ingredients of the food, the price, expire data, *etc.*) attracts the consumers' attention most. Although gaze following is of vital importance, it is extremely challenging because of the reasons below: firstly, actually inferring the gaze point requires the depth information of the scene, head pose and eyeball movement [31,27], nevertheless it is hard to infer the depth of scene with a monocular image. Further, head pose and eyeball movements are not easy to be estimated because of occlusion (usually self-occlusion), as shown in Figure 1 (c); secondly, ambiguity exists for gaze point estimated by different third-view observers with a single view image, as shown in Figure 1 (d); thirdly, the gaze following involves the geometric relationship understanding between target person and other objects/persons in the scene as well as scene contents understanding, which is a difficult task.



Fig. 1. (a) and (b) show the application of gaze following in supermarket scenario. (c) and (d) show the challenges in gaze following (The self-occluded head and ambiguity of gaze point).

To tackle these problems, early works usually simplify the setting to avoid the handicaps for general gaze following, for example, making assumption that face is available for better head pose estimation [32], with multiple inputs for depth inference [18], with eye tracker for ground truth annotation or restricting the application scenario to people looking at each other for disambiguation. However, these simplifications restrict the applications of general gaze following. Recently, Recasens *et al.* propose to study general gaze following under most of general settings [22]. Specifically, they propose a two-pathway method (a gaze pathway and a saliency pathway) based deep neural networks for gaze following. However, in their solution, the saliency pathway and gaze pathway are independent of each other. For a third-view person, when he/she infers the gaze point of a target person, he/she infers the gaze direction first based on head pose, and then estimates the gaze point from the scene contents along the gaze direction, where gaze point denotes the position that one person is looking at in the image, and gaze direction means the direction from head position to gaze point in this paper. In other words, the saliency pathway relies on gaze direction estimation, which is neglected in [22].

In this paper, we propose a two-stage solution to mimic the behavior of a third-view person for gaze following. Specifically, in the first stage, we leverage a gaze direction pathway to predict the gaze direction based on the head image and

head position of target person. There are two motivations for such gaze direction pathway: firstly, it is more natural to infer the gaze direction rather than gaze point merely based on head image and head position; secondly, since the gaze direction can be inferred in the training phase, thus we can introduce a loss w.r.t. gaze direction to facilitate the learning of this gaze direction pathway. Next, we encode the predicted gaze direction as the multi-scale gaze direction fields. In the second stage, based on the gaze direction and the context information of the objects along the gaze direction, we can estimate a heatmap through a heatmap pathway. In this stage, we concatenate the multi-scale gaze direction fields with the original image as the input of the heatmap pathway for heatmap estimation.

A proper dataset is important for the evaluation of gaze following. The only existing gaze following dataset (GazeFollow dataset [22]) is annotated by the third-view persons. In this paper, to evaluate the performance for real problem, we build a video-based gaze following dataset, named Daily Life Gaze dataset (DL Gaze). Particularly, we have 16 volunteers to freely move in 4 different indoor scenes, including working office, laboratory, library, corridor in the building. During the period, they can talk, read books, use their mobile phones, or freely look at other places in the scene. We record the video for them and ask the volunteer to annotate where they look later. There are 95,000 frames in total. Compared with GazeFollow, the ground truth annotated by the persons in the video is more reliable than that annotated by third-view workers. Further, it is a video-based gaze following dataset and records the gaze following for real scenes. Therefore the evaluation of gaze following on this dataset reflects the performance of different methods for real problem.

The main contributions of our paper are summarized as follows: i) we propose a two-stage solution for gaze following task. Our network architecture is inspired by the behavior of human in gaze following, therefore it is more psychological plausible; ii) we use ground truth to supervise the learning of both stages, consequently facilitates the network training. In addition, we introduce multi-scale gaze direction fields for attention prediction, which further improves the performance of gaze following; iii) we collect a video-based gaze following dataset (DL Gaze), with the ground truth annotated by the persons in the video. Therefore the evaluation on this dataset reflects the real performance for gaze following in real problem; iv) Extensive experiments on both datasets validate the effectiveness of our solution for gaze following.

2 Related Work

Gaze following. Previous work about gaze following paid attention to restricted scenes, which added some priors for specific applications. In [32], a face detector was employed to extract face, which was limited for the people looking away from the camera. [17] detected whether people were looking at each other in a movie, which was helpful for interaction. Eye tracker was utilized to predict the next object in order to improve action recognition in [3]. [20] only estimated the gaze direction from head position, but not the specific gaze point. These methods were

applied to a particular scene. Recent works [22,18,23] focused on general gaze following, which had wider applications. Given a single picture containing one or more people, the gaze points of some people in the image were estimated, without any restrictions in [22]. Some extensive works [18,23] focused on multi-modality image or predicted gaze point in videos. The RGB-D image was introduced to predict gaze in images and videos [18] because the multi-modality data provided 3D head pose information in order to find more accurate gaze point. In [23], the cross-frame gaze point in videos could be predicted for the people in a frame.

Eye tracking. Eye tracking is strongly related to gaze following. Different from gaze following, eye tracking technology inferred which direction or which point on the screen one person was looking at [29]. Previous work [33,6] built the geometry model to infer the gaze point on the screen target. Recently, many appearance-based methods [11,31] solved the problem by learning a complex function from the eye images to gaze point, which needed large-scale dataset. These methods took the eye images and face image as inputs because gaze direction could be determined according to the eye movement and head pose [31]. However, the eye images could not be utilized to predict gaze point because they were occluded or very noisy in gaze following. Thus, gaze following direction is almost obtained from the head image.

Saliency. Saliency detection and gaze following are two different tasks [22,23] even though they were closely related. Saliency detection predicts fixation map from observers out of the original images [7,9,14]. Gaze following in image predicts the position that people in a scene were looking at. Previous works about saliency prediction considered the low-level features and saliency maps at different scales [8]. Subsequently, the features from different levels were combined to model a bottom-up, top-down architecture [9]. Recently, deep neural networks have been applied to saliency prediction and achieve great success [13,19]. However, the object in the gaze point region may be not salient, which reveals that it is hard to find the gaze point through a saliency algorithm directly.

3 Approach

Inspired by the behavior of human in gaze following, we propose a two-stage solution. Specifically, when a third-view person estimates the gaze of the target person, he/she first estimates the gaze direction of the target based on the head image, then the gaze point is predicted based on the scene content along the gaze direction. Similarly, we feed the head image and its position in the image into a gaze direction pathway for gaze direction prediction in the first stage, and then the multi-scale gaze direction fields are encoded. In the second stage, the gaze direction fields are concatenated with the original image as the input of heatmap pathway for heatmap regression. It is worth noting that all components in our network are differentiable and the whole network can be trained with an end-to-end learning. The network architecture is shown in Figure 2.

3.1 Gaze direction pathway

Gaze direction pathway takes head image and head position as inputs for gaze direction prediction. We feed the head image into a ResNet-50 for feature extraction, and then concatenate head features with head position features encoded by a network with three fully connected layers for gaze direction prediction. Different from work of Recasens *et al.* [22], which takes head image and head position for gaze mask prediction, our network only estimates the gaze direction. There are two reasons accounting for our solution. Firstly, it is easier to infer the gaze direction than gaze mask merely based on head image and its position. Secondly, we can use gaze direction to supervise the learning of gaze direction pathway to make it more robustly trained. It is also worth noting that the predicted gaze direction would be used to generate gaze direction fields, which is further used for heatmap regression in the heatmap pathway, and the optimization of heatmap would also update the parameters in the gaze direction pathway.

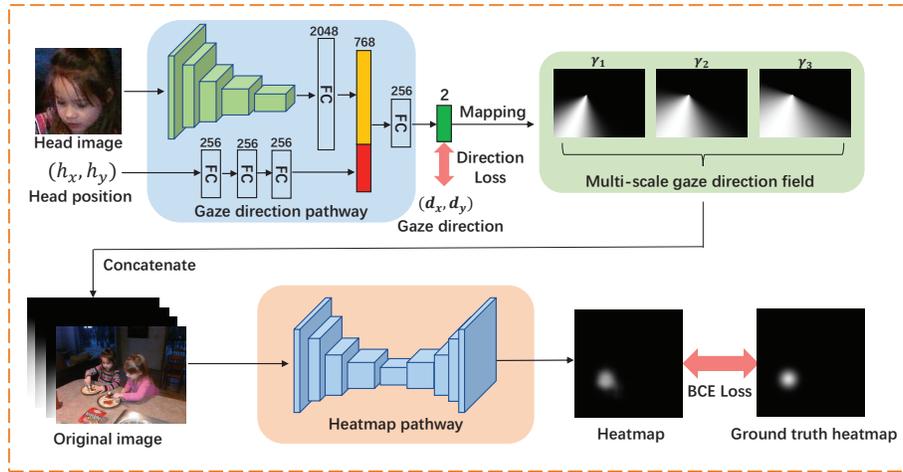


Fig. 2. The network architecture for gaze following. There are two modules in this network: gaze direction pathway and heatmap pathway. In the first stage, a coarse gaze direction is predicted through gaze direction pathway, and then it is encoded as multi-scale gaze direction fields. We concatenate the multi-scale fields and the original image to regress heatmap of final gaze point through heatmap pathway.

3.2 Gaze direction field

Once the gaze direction is estimated, the gaze point is likely to be along the gaze direction. Usually, the field of view (FOV) of the target person is simplified as a cone with the head position as the apex of the cone. So given a point

$P = (p_x, p_y)$, if we do not consider the scene contents, then the probability of the point P being the gaze point should be proportional to the angle θ between line L_{HP} and predicted gaze direction, here $H = (h_x, h_y)$ is the head position, as shown in Figure 3 (a). If θ is small, then the probability of the point being gaze point is high, otherwise, the probability is low. We utilize the cosine function to describe the mapping from the angle to the probability value. We denote the probability distribution of the points being gaze point without considering the scene contents as the **gaze direction field**. Thus, gaze direction field is a probability map, where intensity value of each point shows the probability that this point is the gaze point. Its size is the same as the scene image.

Particularly, the line direction of L_{HP} can be calculated as follows:

$$G = (p_x - h_x, p_y - h_y) \quad (1)$$

Given an image with size $W \times H$ (here W and H are the width and height of image, respectively), we denote the predicted gaze direction as $\hat{d} = (\hat{d}_x, \hat{d}_y)$, then the probability of the point P being the gaze point can be calculated as follows:

$$Sim(P) = \max\left(\frac{\langle G, \hat{d} \rangle}{|G||\hat{d}|}, 0\right) \quad (2)$$

Here we let the probability of P being gaze point to be 0 when the angle between gaze direction and line L_{HP} is larger than 90° , which means the real gaze direction should not contradict with the predicted gaze direction. We depict the calculation of gaze direction field.

If the predicted gaze direction is accurate, it is desirable that the probability distribution is sharp along the gaze direction, otherwise, it is desirable that the probability changes smoothly. In practice, we leverage multi-scale gaze direction fields with different sharpness for heatmap prediction. Specifically, we use the following way to control the sharpness of the gaze direction field:

$$Sim(P, \gamma) = [Sim(P)]^\gamma \quad (3)$$

Here γ controls aperture of the FOV cone. Larger γ corresponds to a FOV cone with smaller aperture, as shown in Figure 2. In our implementation, considering the change rate of $Sim(P, \gamma)$, we empirically set $\gamma_1 = 5, \gamma_2 = 2, \gamma_3 = 1$. More details about γ can be found in the supplementary material.

It also worth noting that the gaze direction fields are differentiable w.r.t. the network parameters of gaze direction pathway, so that the whole architecture can be trained with an end-to-end learning strategy.

3.3 Heatmap pathway

Gaze direction fields encode the distribution of gaze points inferred from gaze direction, together with scene contents, we can infer the gaze point. Specifically, we concatenate the original image and the multi-scale gaze direction fields, and feed them into a heatmap pathway for heatmap regression. The point corresponding

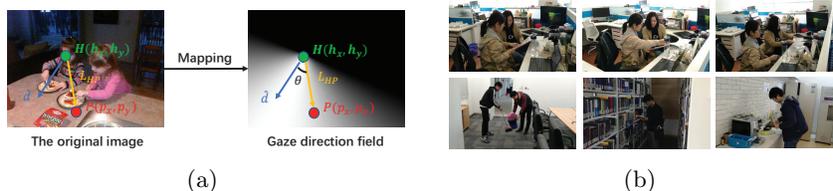


Fig. 3. (a) The original image: the blue line shows gaze direction of the left girl inside the image, and the green dot shows the head position. Gaze direction field, which measures the probability of each point being gaze point with cosine function between the line direction of L_{HP} and predicted gaze direction \hat{d} . (b) Our DL Gaze dataset.

to the maximum value of the heatmap is considered as the final gaze point. In practice, we leverage a feature pyramid network (FPN) [15] for the heatmap pathway in light of its success in object detection. The last layer of heatmap pathway is followed with a Sigmoid activation function, which guarantees the probability of each pixel falls into $[0, 1]$.

There are two reasons to predict probability heatmap instead of a direct gaze point coordinate:

- As pointed in [25], mapping from image to the coordinates of gaze point directly is a highly non-linear function. Compared with gaze point estimation, heatmap prediction is more robust, which means even some entries of heatmap are not accurately predicted, the gaze point prediction based on heatmap can still be correct. Thus heatmap regression is more commonly used in many applications, including pose estimation [21] and face alignment. The experimental results in section 4 also validate the advantage of heatmap regression over gaze point regression.
- Gaze following in an image is sometimes ambiguous [22] due to the lack of the ground truth. Different workers may vote for different gaze points, especially when the eye images are invisible, the head image is low-resolution or occluded. Thus, the gaze point is usually multimodal, and the output of network is expected to support the multimodal prediction. Heatmap regression satisfies such requirement.

Following [21], the heatmap of ground truth gaze point is generated by centering a Gaussian kernel at the position of gaze point as follows:

$$H(i, j) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(i-g_x)^2 + (j-g_y)^2}{2\sigma^2}} \quad (4)$$

where $g = (g_x, g_y)$ and $H(i, j)$ are the ground truth gaze point and its heatmap, respectively. σ is the variance of Gaussian kernel. We empirically set $\sigma = 3$ in our implementation.

3.4 Network training

The inputs of our network consist of three parts: head image, head position and the original image. The head and original image are resized to 224×224 , and the head position is the coordinate when the original image size is normalized to 1×1 . The outputs of network consist of two parts: gaze direction and visual attention. The gaze direction is the normalized vector from the head position to gaze point and the visual attention is a heatmap with size 56×56 , whose values indicate the probability that the gaze point falls here.

Specifically, the gaze direction loss is:

$$\ell_d = 1 - \frac{\langle d, \hat{d} \rangle}{|d| |\hat{d}|} \quad (5)$$

where d and \hat{d} are the ground truth and predicted gaze direction, respectively.

We employ the binary cross entropy loss (BCE Loss) for heatmap regression, which is written as follows:

$$\ell_h = -\frac{1}{N} \sum_{i=1}^N H_i \log(\hat{H}_i) + (1 - H_i) \log(1 - \hat{H}_i) \quad (6)$$

where H_i and \hat{H}_i are the i -th entry of ground truth heatmap and predicted visual heatmap, respectively. N is with the size 56×56 .

The whole loss function consists of gaze direction loss and heatmap loss:

$$\ell = \ell_d + \lambda \ell_h \quad (7)$$

where λ is the weight to balance ℓ_d and ℓ_h . We set $\lambda = 0.5$ in our experiments.

4 Experiments

4.1 Dataset and evaluation metric

Dataset. The GazeFollow dataset [22] is employed to evaluate our proposed method. The images of this dataset are from different source datasets, including SUN [26], MS COCO [16], Actions 40 [28], PASCAL [2], ImageNet [24] and Places [30], which is challenging due to variety of scenarios and amounts of people. The whole dataset contains 130,339 people and 122,143 images. The gaze points of people are inside the image. There are 4,782 people of dataset used for testing and the rest for training. To keep the evaluation consistency with existing work, we follow the standard training/testing split in [22].

To validate the performance of different gaze following algorithms for real scenarios, we also build a video-based Daily Life Gaze following dataset (DL Gaze). Specifically, DL Gaze contains the activities of 16 volunteers in 4 scenes, and these scenes include working office, laboratory, library and corridor in the building. They can freely talk, read books, use their mobile phones, and look

at other places in the scene, as shown in Figure 3 (b). We record the video for these volunteers with an iPhone 6s. Then we ask each volunteer to annotate what he/she looks at. Two frames are annotated per second. So the ground truth annotation is more reliable. It worth noting the occlusion also exists and there is severe change of illumination. Therefore our dataset is very challenging. The performance of gaze following in our dataset reflects the capability of gaze following in real scenarios. There are 86 videos, 95,000 frames (30 fps) in total. We test the model trained on GazeFollow with our dataset directly.

Evaluation metric. Following [22], we employ these metrics (**AUC**, **Dist**, **MDist**, **Ang**) to evaluate the difference between the predicted gaze points and their corresponding ground truth. Details can be found in the supplementary material. In addition, we also introduce **Minimum angular error (MAng)**, which measures the minimum angle between the predicted gaze direction and all ground truth annotations:

4.2 Implementation details

We implement the proposed method based on the PyTorch framework. In the training stage, We employ a ResNet-50 to extract head image feature and encode the original image feature. The network is initialized with the model pretrained with ImageNet [1]. When the first stage of training converges, we train the heatmap pathway and finally we finetune the whole network with an end-to-end learning strategy. The hyper-parameters of our network are listed as follows: batch size (128), learning rate ($1e^{-4}$), weight decay (0.0005). Adaptive moment estimation (Adam) algorithm [10] is employed to train the whole network.

4.3 Performance evaluation

We compare our proposed method with the following state-of-the-art gaze following methods:

- Judd *et al.* [9]: Such a method uses a saliency model as a predictor of gaze and the position with maximum saliency value is used as predicted gaze point inside the image.
- SalGAN [19]: SalGAN [19] is the latest saliency method, and it takes the original image as input to generate visual heatmap. The position of maximum in visual attention is regarded as gaze point.
- SalGAN for heatmap: We replace the FPN with SalGAN in heatmap pathway, and all the rest components are the same with our method.
- Recasens *et al.* [22]: The gaze pathway and saliency pathway are introduced to extract the image and head feature, and both features are fused to get the final gaze point. The supervision is introduced in the last layer.
- Recasens *et al.**: For a fair comparison, we modify the backbone of Recasens *et al.* [22] from the AlexNet [12] to ResNet-50 [5] to extract head feature and image feature. All other parts remain the same as Recasens *et al.* [22].

- One human [22]: A third-view observer is employed to predict gaze points on the testing set in [22]. It is desirable that machine can achieve the human level performance.

Table 1. Performance comparison with existing methods on the GazeFollow dataset. One-scale and multi-scale correspond to the number of gaze direction fields in our model. For one-scale model, $\gamma = 1$.

Methods	AUC	Dist	MDist	Ang	MAng
Center [22]	0.633	0.313	0.230	49.0°	-
Random [22]	0.504	0.484	0.391	69.0°	-
Fixed bias [22]	0.674	0.306	0.219	48.0°	-
SVM + one grid [22]	0.758	0.276	0.193	43.0°	-
SVM + shift grid [22]	0.788	0.268	0.186	40.0°	-
Judd <i>et al.</i> [9]	0.711	0.337	0.250	54.0°	-
SalGAN [19]	0.848	0.238	0.192	36.7°	22.4°
SalGAN for heatmap	0.890	0.181	0.107	19.6°	9.9°
Recasens <i>et al.</i> [22]	0.878	0.190	0.113	24.0°	-
Recasens <i>et al.*</i> [22]	0.881	0.175	0.101	22.5°	11.6°
One human [22]	0.924	0.096	0.040	11.0°	-
Ours (one-scale)	0.903	0.156	0.088	18.2°	9.2°
Ours (multi-scale)	0.906	0.145	0.081	17.6°	8.8°

The experiment results in Table 1 and Table 2 show that our model outperforms all baselines in terms of all evaluation metrics. We also have the following findings: (1) Recasens *et al.** outperforms Recasens *et al.* shows the importance of the basic network. (2) SalGAN has the better performance than Judd *et al.*, which shows that better saliency detection method agrees with visual attention better. (3) Although employing the same basic network (ResNet-50), our method (one-scale) still achieve the better performance than Recasens *et al.**, which proves the soundness of our human behavior inspired a two-stage solution for gaze following. (4) The multi-scale model achieves better performance than that of one-scale, which validates the importance of multi-scale fields fusion. (5) Performance on our dataset is worse than that on GazeFollow, which shows the challenge of gaze following in real applications. (6) The improvement of our method over SalGAN for heatmap validates the effectiveness of FPN for heatmap regression.

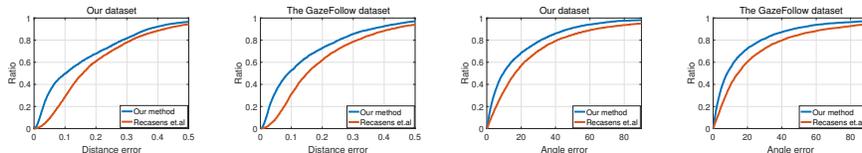
We further compare our method with Recasens *et al.* using accumulative error curve and the results are shown in Figure 4. We can see that our method usually achieves better prediction than the work of Recasens *et al.* [22].

Ablation study. In order to evaluate the effectiveness of every component of different inputs and network. We design the following baselines:

- original image: We directly feed the original image into heatmap pathway for heatmap regression.

Table 2. Performance comparison with existing methods on our dataset. Each frame only contains one gaze point, so only Dist and Ang are used for performance evaluation.

Methods	Dist	Ang
Recasens <i>et al.</i> [22]	0.203	26.9°
Recasens <i>et al.</i> * [22]	0.169	21.4°
Ours (multi-scale)	0.157	18.7°

**Fig. 4.** Accumulative error curves of different methods on both datasets.

- original image + ROI head: We directly feed the original image into heatmap pathway for heatmap regression. Further, we directly extract the features corresponding to Region of Interest (ROI, the region of head) from the heatmap pathway and use it for gaze direction regression. Then we train the whole network with multi-task learning.
- w/o mid-layer supervision: The gaze direction supervision is removed, and both pathways are trained with an end-to-end learning strategy. Only one-scale gaze direction field is concatenated to the original image.

Table 3. The results of ablation study.

Methods	AUC	Dist	MDist	Ang	MAng
Original image	0.839	0.212	0.146	32.6°	21.6°
Original image + ROI head	0.887	0.182	0.118	22.9°	10.7°
W/O mid-layer supervision	0.875	0.178	0.101	24.4°	12.5°
Ours (one-scale)	0.903	0.156	0.088	18.2°	9.2°

The experiment results are listed in Table 3. We can see that predicting heatmap merely based on the scene image is not easy, even the head and its position are already included in the image. With gaze direction as supervision (original image+ ROI head) to aid the heatmap pathway learning, the performance can be boosted. With a gaze direction pathway to predict gaze direction, our method greatly outperforms original image-based solution for gaze following, which further validates the importance of two-stage solution. Further the improvement of our method (one scale) over w/o mid-layer supervision validates the importance of gaze direction prediction, which is an advantage of our solution, *i.e.*, our two-stage method benefits from gaze direction prediction.

The information fusion. In the second stage, we combine the gaze direction field and image content information. However, how to choose the position (early, middle, late fusion) and way (multiplication or concatenation) of fusing? Here we compare our method with the following information fusion strategies:

- Middle fusion (mul): Fuse the gaze direction field and the image content feature map (7×7) with multiplication in the middle layer.
- Middle fusion (concat): Fuse the gaze direction field and the image content feature map (7×7) with concatenation in the middle layer.
- Early fusion (mul): Fuse the gaze direction field and the image content feature map (28×28) with multiplication in the early layer of encoder in heatmap pathway.
- Late fusion (mul): Fuse the gaze direction field and the image content feature map (28×28) with multiplication in the last layer of decoder in heatmap pathway.
- Image fusion (mul): Directly multiply the original image with gaze direction field.

Table 4. Different information fusion strategies.

Methods	AUC	Dist	MDist	Ang	MAng
Middle fusion (mul)	0.882	0.183	0.118	21.7°	10.7°
Middle fusion (concat)	0.884	0.177	0.105	21.0°	10.5°
Early fusion (mul)	0.898	0.160	0.098	18.7°	9.6°
Late fusion (mul)	0.888	0.176	0.102	20.1°	10.1°
Image fusion (mul)	0.895	0.163	0.096	19.3°	9.7°
Ours (concat)	0.903	0.156	0.088	18.2°	9.2°

Table 4 shows the results of different information fusion strategies. We can find that early fusion usually obtains higher performance than middle and late fusion, which implies early suppression of useless scene contents is important for gaze following. Furthermore, we find that usually concatenating the gaze direction field with image or feature achieves slightly better results than the multiplication. The possible reason is that the predicted gaze direction may not be very accurate, and the multiplication between image and the gaze direction field would lead to the change of intensities of pixels and cause information loss. While for concatenation, the information is still there and the heatmap pathway can tackle the heatmap prediction, even the gaze direction fields are not accurate.

Objective. Since the predicted gaze point may be multimodal, we introduce heatmap as the ground truth. Here, we also compare our method with networks based on other types of outputs, including:

- Point: We employ two ResNet-50 to extract features for both original image and head image. Such a comparison is fair because the encoder part of FPN is also ResNet-50. In this baseline, we only predict gaze point.

Table 5. The evaluation of different objectives.

Methods	AUC	Dist	MDist	Ang	MAng
Point	0.892	0.173	0.103	21.9°	10.5°
Multi-task point	0.900	0.165	0.097	20.4°	10.1°
Shifted grid [22]	0.899	0.171	0.096	21.4°	10.3°
Heatmap (our)	0.903	0.156	0.088	18.2°	9.2°

- Multi-task regression: The network architecture is the same as point regression, but it predicts both gaze direction and gaze point simultaneously.
- Shifted grid: Based on our network architecture, shifted grid (10×10) [22] is utilized to classify the gaze point into different grids.

**Fig. 5.** Some prediction results on the testing set, the red lines indicate the ground truth gaze and the yellow ones are the predicted gaze.

The comparison results of different objectives are listed in Table 5. In our network architecture, heatmap regression achieves the best results than both point and shifted grid based objectives. As aforementioned, heatmap regression is more robust than directly point prediction because even a portion of heatmap values is incorrect, it is still possible to correctly predict the gaze point. Thus such a heatmap regression strategy is commonly used for human pose estimation [4]. Our experiments also validate its effectiveness for gaze following.

4.4 Visualization of predicted results

We show the predicted gaze points and their ground truth in Figure 5, and predicted heatmaps in Figure 6 (a). We can see that for most of points, our

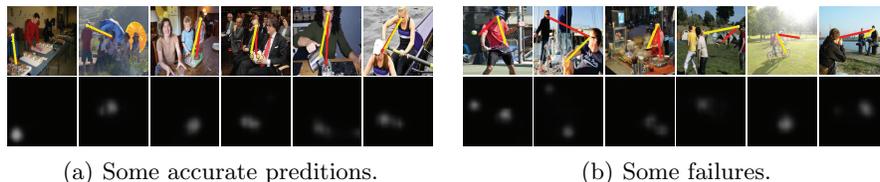


Fig. 6. The first row: ground truth (red lines) and predicted gaze (yellow lines). The second row: predicted heatmaps. (Please zoom in for details.)

method can predict gaze points accurately (As shown in Figure 4, when the distance error is 0.1, the portions of correctly predicted points is 50% and 45% on the GazeFollow dataset and DL Gaze, respectively.). There are two reasons contribute the good performance of our method: i) our two-stage solution agrees with the behavior of human, and gaze direction field would help suppress the regions falling out of gaze direction, which consequently improves the heatmap regression; ii) the supervision on gaze direction helps train a more robust network for gaze following. We also show some failures in Figure 6 (b). The first three columns of examples show that our predictions can be multimodal. Although the position of heatmap maximum is not right, some others peaks can also predict the gaze point. Regarding the last three columns of failures, we can see that the predicted heatmap is inaccurate. This probably caused by the small head or head occlusion, which makes gaze direction and gaze point prediction extremely difficult, even for us human.

5 Conclusion

In this paper, we proposed a two-stage solution for gaze tracking. In stage I, we feed the head image and its position for gaze direction prediction. Then we use gaze direction field to characterize the distribution of gaze points without considering the scene contents. In stage II, the gaze direction fields are concatenated with original image, and fed into a heatmap pathway for heatmap regression. The advantages of our solution are two-fold: i) our solution mimics the behavior of human in gaze following, therefore it is more psychological plausible; ii) besides leverage heatmap to supervise the training of our network, we can also leverage gaze direction to facilitate the training of gaze direction pathway, therefore our network can be more robustly trained. We further build a new DL Gaze dataset to validate the performance of different gaze following methods in real scenarios. Comprehensive experiments show that our method significantly outperforms existing methods, which validates the effectiveness of our solution.

Acknowledgement. This project is supported by NSFC (No. 61502304).

References

1. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. pp. 248–255. IEEE (2009)
2. Everingham, M., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. *International journal of computer vision* **88**(2), 303–338 (2010)
3. Fathi, A., Li, Y., Rehg, J.M.: Learning to recognize daily actions using gaze. In: *European Conference on Computer Vision*. pp. 314–327. Springer (2012)
4. Fragkiadaki, K., Levine, S., Felsen, P., Malik, J.: Recurrent network models for human dynamics. In: *Computer Vision (ICCV), 2015 IEEE International Conference on*. pp. 4346–4354. IEEE (2015)
5. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 770–778 (2016)
6. Hennessey, C., Nouredin, B., Lawrence, P.: A single camera eye-gaze tracking system with free head motion. In: *Proceedings of the 2006 symposium on Eye tracking research & applications*. pp. 87–94. ACM (2006)
7. Itti, L., Koch, C.: Computational modelling of visual attention. *Nature reviews neuroscience* **2**(3), 194 (2001)
8. Itti, L., Koch, C., Niebur, E.: A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on pattern analysis and machine intelligence* **20**(11), 1254–1259 (1998)
9. Judd, T., Ehinger, K., Durand, F., Torralba, A.: Learning to predict where humans look. In: *Computer Vision, 2009 IEEE 12th international conference on*. pp. 2106–2113. IEEE (2009)
10. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014)
11. Krafska, K., Khosla, A., Kellnhofer, P., Kannan, H., Bhandarkar, S., Matusik, W., Torralba, A.: Eye tracking for everyone. *arXiv preprint arXiv:1606.05814* (2016)
12. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: *Advances in neural information processing systems*. pp. 1097–1105 (2012)
13. Kümmerer, M., Theis, L., Bethge, M.: Deep gaze i: Boosting saliency prediction with feature maps trained on imagenet. *arXiv preprint arXiv:1411.1045* (2014)
14. Leifman, G., Rudoy, D., Swedish, T., Bayro-Corrochano, E., Raskar, R.: Learning gaze transitions from depth to improve video saliency estimation. In: *Proc. IEEE Int. Conf. on Computer Vision*. vol. 3 (2017)
15. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: *CVPR*. vol. 1, p. 4 (2017)
16. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: *European conference on computer vision*. pp. 740–755. Springer (2014)
17. Marín-Jiménez, M.J., Zisserman, A., Eichner, M., Ferrari, V.: Detecting people looking at each other in videos. *International Journal of Computer Vision* **106**(3), 282–296 (2014)
18. Mukherjee, S.S., Robertson, N.M.: Deep head pose: Gaze-direction estimation in multimodal video. *IEEE Transactions on Multimedia* **17**(11), 2094–2107 (2015)

19. Pan, J., Canton, C., McGuinness, K., O'Connor, N.E., Torres, J., Sayrol, E., Giro-i Nieto, X.: Salgan: Visual saliency prediction with generative adversarial networks. arXiv preprint arXiv:1701.01081 (2017)
20. Parks, D., Borji, A., Itti, L.: Augmented saliency model using automatic 3d head pose detection and learned gaze following in natural scenes. *Vision research* **116**, 113–126 (2015)
21. Pfister, T., Charles, J., Zisserman, A.: Flowing convnets for human pose estimation in videos. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 1913–1921 (2015)
22. Recasens*, A., Khosla*, A., Vondrick, C., Torralba, A.: Where are they looking? In: *Advances in Neural Information Processing Systems (NIPS)* (2015), * indicates equal contribution
23. Recasens, A., Vondrick, C., Khosla, A., Torralba, A.: Following gaze in video. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 1435–1443 (2017)
24. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. *International Journal of Computer Vision* **115**(3), 211–252 (2015)
25. Tompson, J.J., Jain, A., LeCun, Y., Bregler, C.: Joint training of a convolutional network and a graphical model for human pose estimation. In: *Advances in neural information processing systems*. pp. 1799–1807 (2014)
26. Xiao, J., Hays, J., Ehinger, K.A., Oliva, A., Torralba, A.: Sun database: Large-scale scene recognition from abbey to zoo. In: *Computer vision and pattern recognition (CVPR), 2010 IEEE conference on*. pp. 3485–3492. IEEE (2010)
27. Xiong, X., Liu, Z., Cai, Q., Zhang, Z.: Eye gaze tracking using an rgbd camera: a comparison with a rgb solution. In: *Proceedings of the 2014 ACM International Joint Conference on Pervasive and Ubiquitous Computing: Adjunct Publication*. pp. 1113–1121. ACM (2014)
28. Yao, B., Jiang, X., Khosla, A., Lin, A.L., Guibas, L., Fei-Fei, L.: Human action recognition by learning bases of action attributes and parts. In: *Computer Vision (ICCV), 2011 IEEE International Conference on*. pp. 1331–1338. IEEE (2011)
29. Zhang, X., Sugano, Y., Fritz, M., Bulling, A.: Appearance-based gaze estimation in the wild. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 4511–4520 (2015)
30. Zhou, B., Lapedriza, A., Xiao, J., Torralba, A., Oliva, A.: Learning deep features for scene recognition using places database. In: *Advances in neural information processing systems*. pp. 487–495 (2014)
31. Zhu, W., Deng, H.: Monocular free-head 3d gaze tracking with deep learning and geometry constraints. In: *The IEEE International Conference on Computer Vision (ICCV)* (Oct 2017)
32. Zhu, X., Ramanan, D.: Face detection, pose estimation, and landmark localization in the wild. In: *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. pp. 2879–2886. IEEE (2012)
33. Zhu, Z., Ji, Q.: Eye gaze tracking under natural head movements. In: *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*. vol. 1, pp. 918–923. IEEE (2005)

Supplementary Material

Dongze Lian^{*[0000-0002-4947-0316]}, Zehao Yu^{*[0000-0002-6559-9830]}, and
Shenghua Gao^{†[0000-0003-1626-2040]}

School of Information Science and Technology, ShanghaiTech University
{liandz, yuzh, gaoshh}@shanghaitech.edu.cn

This supplementary material provides some supplementary notes and results. First, we give the detailed information of evaluation metric in Section 4.1 (**AUC**, **Dist**, **MDist**, **Ang**). Secondly, we analyze the running time of these methods. Finally, we explain the choice of parameter γ in the Eq. 3.

1 Evaluation metric

Area Under Curve (AUC): The area under ROC curve, which is generated according to [9].

L_2 distance (Dist): The Euclidean distance between predicted gaze point and the average of ground truth annotations. The original image size is normalized to 1×1 .

Minimum L_2 distance (MDist): The minimum Euclidean distance between predicted gaze point and all ground truth annotations. The original image size is normalized to 1×1 .

Angular error (Ang): The angular error between predicted gaze direction and ground truth direction corresponding to average gaze point.

2 Running time

We compare the performance and the running time of different baselines in the testing phase. All the methods are tested with an NVIDIA Titan X GPU. The metric used here is the L_2 distance between predicted gaze point and the average of ground truth annotation. The running time refers to the running time for one testing sample, which is an average running time of 1000 test samples (The unit is millisecond). We repeat the experiments ten times.

3 The choice of γ

If the predicted gaze direction is accurate, it is desirable that the probability distribution is sharp along the θ , otherwise, it is desirable that the probability changes smoothly. Thus, we choose three different γ to represent the decay rate in the Eq. 3 and our experiments verify the effectiveness of multi-scale gaze

^{*}The authors contribute equally.

[†]Corresponding author.

direction field. About the choice of γ , considering the decay rate, we empirically set γ to 1, 2, 5, respectively. That is because when the value of the cosine function is 0.5, the angle is about 60° , 45° and 30° . The gaze direction fields of different γ are shown in Fig. 2 in the paper.

Table 1. The performance and running time of different methods on GazeFollow.

Methods	Dist	time (ms)
SalGAN [19]	0.238	8.8
Recasens <i>et al.</i> [22]	0.190	10.4
Recasens* <i>et al.</i> [22]	0.175	15.8
Our method (one-scale)	0.156	16.1
Our method (multi-scale)	0.145	16.3