

Where Are They Really Looking?

Ambika Verma Brady Zhou

University of Texas at Austin

{ambika, brady.zhou}@utexas.edu

December 2, 2016

Abstract

In this paper, we present a framework to predict where actors in an image are looking by utilizing joint attention and/or interaction characteristics. This approach is specifically useful in case of more than one person/actor being present in a scene. We use a combination of CNN [1] and MRF [2] models to accomplish this task. The approach is motivated from recent state of the art work in Gaze Following (CNN) and Social Attention (MRF) domains.

1 Introduction

In our day to day life we tend to interact with people and/or objects in a fairly predictable manner specifically in terms of social attention and/or interaction such as we tend to notice things that a group of people might be attending to. Additionally our gaze (direction we look in) is a strong indication of what we might do next. Thus, it is imperative for a computer or robot trying to predict the type of actions or interactions occurring in an image or video, to follow gaze of actors in a scene effectively. As a result, gaze following can be used for different high level tasks such as predicting visual saliency, activity recognition, active perception and behavioural analysis.

As noted earlier we solely concentrate on gaze following task for multiple actors in an image or video, with the idea in mind that joint gaze prediction can benefit the task of gaze following since humans interact in a predictive manner. This leads to some of the assumptions made in this paper (which are derived from social attention do-

main) [2], namely head orientation is a robust indicator of gaze direction and if another actor in the scene is looking at a particular location it is more likely for another person (with head orientation in the same direction) to be looking at it as well.

1.1 Related Work

As noted by [1] there are only a few works which build upon the task of gaze following. There are three different scenarios where this task has been studied - Free-viewing Saliency, Social Attention/Interaction and true Gaze Following, we review each of tasks briefly in this section.

Identifying types of social attention and interactions occurring in a scene inherently limits itself to scenes with multiple people only [2] [3] [4], [5], though, it boasts of intuitive methods to perform joint gaze prediction. [2] develops MRF and CRF methods for their task of predicting type of interaction taking place in a scene, which can be applied to long first-person videos to filter out useful clips or subsets (similar to video summarization). This work utilizes an MRF model to predict where people are looking in a scene to later classify the type of interaction taking place using CRF. Thus, gaze following is utilized in such techniques in an implicit manner. The model developed utilizes joint predictions effectively, though is heavily biased to look at other faces in the scene. This bias is acceptable for identifying the type of interaction amongst people but does not generalize well for the gaze following task, wherein actors can be looking at another actor or an object.



Figure 1: Test set example images with annotations.

On the other hand works in free-viewing saliency and gaze following tasks do include scenes with both multiple people and only one person. [2] is a CNN based model which is specifically tasked for gaze following utilizing head position and orientation along with saliency to accomplish the task.

This method has several advantages, such as it is robust in case of single actor being present in a scene, does not suffer from any evident bias of looking at people only (such as in [2]) and the network architecture does not require extensive hand-crafted functions (such as those required in probabilistic models such as MRF or CRF). Even though, these merits are far reaching the model uses a naive approach in cases where multiple actors are present in the scene, it resorts to providing predictions for each actor individually.

Another relevant work which utilizes gaze following is [6] in which an actors gaze is used as an additional feature to generate more accurate free-viewing saliency maps. It has been noted in numerous works [6], [3], [7], [8], [9] that the gaze of an actor in a scene influences where passive observers look during free-viewing. In this work the authors use head pose [9] along with low level saliency [10] to predict more accurately the semantically salient regions in a scene. This work does utilize joint predic-

tions in case of multiple actors in the scene and it is worthy to note that considerable performance improvement is achieved over state of the art saliency methods [10]. Even though this method can be extended to gaze following task, it has not been evaluated on the same.

Works such as [2] and [6] utilize gaze following and joint predictions to achieve a secondary task. Thus, our approach is to combine the works of [2] and [1] in a constructive manner to add joint prediction capability to the model developed in [2] with the goal of achieving higher performance for cases with multiple actors in the scene.

2 Technical Approach

As in [1] we also assume that the head positions are provided and as demonstrated in [6] a head detector with pose estimation can be effectively used with a few modifications.

We plan to use the CNN network from [1] with post-processing added through an MRF model which is designed on the lines of model used in [2] to account for joint gaze predictions. The overall architecture of our approach is illustrated in Fig 1.

A given image along with head position is provided as input to the CNN. We only use images where there

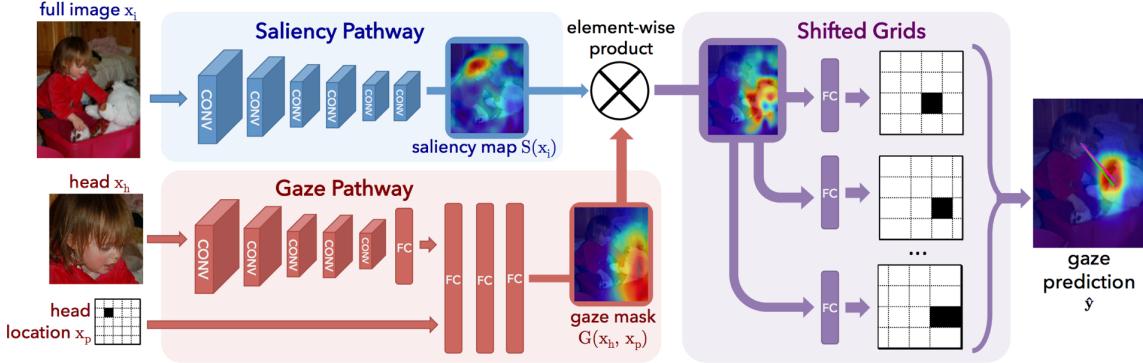


Figure 2: The network architecture from [1].

are multiple actors in the scene. The CNN then iterates through the image for each actor and gives as an output a two-dimensional position (g_x, g_y) normalized by the image dimensions to lie between 0 and 1. This prediction is then used as an initial seed for MRF model. In addition the MRF model is also supplied with head location and orientation vectors, where orientations are obtained using the head position along with predicted gaze location and is jittered to replicate errors one would observe using a head detector. The MRF model then optimizes the unary and pairwise potentials to jointly refine the CNN predictions based on the predictions obtained for other actors in the scene.

Following sections discuss the CNN and MRF models in detail.

2.1 CNN Model

The CNN model consists of two pathways, namely Gaze pathway and Saliency pathway. As noted in [1], the gaze pathway tends to learn predicting the direction of gaze based on the head orientation (this is the reason why we use the CNN prediction and head position to generate orientation vectors for the MRF model) and the saliency pathway tends to learn to find salient regions in the scene. The model is implemented in Caffe [11] and the two pathways are equivalent to the first five layers of AlexNet [12]. The saliency pathway is initialized with weights from Places-CNN [13] and gaze pathway is initialized with ImageNet-CNN [14]. The saliency and gaze masks obtained are then combined through an elementwise layer

and final prediction is made through shifted grids. More details can be found in [1].

2.2 MRF Model

Markov random fields are a form of probabilistic graph models that have been heavily used throughout the computer vision field to solve classification problems [15], [16]. Markov random fields have proved to be effective in certain applications like image segmentation and parsing text (NLP).

The goal of using an MRF is find a labeling $\{p_1 = l_1, p_2 = l_2, \dots, p_n = l_n\}$ such that the joint probability $P(p_1 = l_1, p_2 = l_2, \dots, p_n = l_n)$ is maximized.

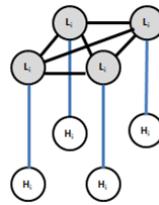


Figure 3: Markov random field depiction.

We use the graphical model from [2] to refine the gaze predictions from [1]. The MRF model is used to utilize the relationship between the multiple actors in a scene. The motivation for this approach stems from example from scenes where multiple actors are looking at the same location. In case of individual predictions there is no way for

benefitting from the information we have about the gaze of other actors in the scene, whereas joint prediction allows us to give higher weights to locations which lie in the direction of gaze of an actor and are also looked at by other actors in the scene.

The MRF model implementation uses the following unary and pairwise potentials. The unary potentials allows the model to predict the direction of the gaze based on head position and orientation vectors and on the other hand the pairwise potentials capture the interaction between actors in the scene.

$$\phi_u = \phi_1 \cdot \phi_2 \cdot \phi_3 \quad (1)$$

The unary potential is a product over three different potential functions. Each of these potentials seeks to score the probability of a gaze (higher is better).

$$\phi_1(l_i = l, d_i, f_i) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left\{ -\frac{\|d_i - (l - f_i)\|^2}{2\sigma^2} \right\} \quad (2)$$

The first unary term is a guassian distribution of how well the predicted gaze aligns with the head orientation. This scores the labeling of a person i 's gaze to a location l . l is a label, the cell of the predicted gaze. d_i is a unit vector, the head orientation. f_i is the position of the subject's face. σ is a constant.

$$\phi_2(l_i = l, f_i) = \frac{1}{1 + \exp \{-c_t \cdot \|l - f_i\|\}} \quad (3)$$

The second unary term is a thresholding function that assigns the score to 0 if the predicted gaze is too close to a person's face. This is used to prevent extremely short gazes (labeling a person's gaze to their own face). c_t is a constant.

$$\phi_3(l_i = l) = \begin{cases} c_p, & \text{if } l \in \{f_1, f_2, \dots, f_n\}. \\ 1, & \text{otherwise.} \end{cases} \quad (4)$$

The third unary term is to bias subjects' gazes towards other faces in the scene. c_p is a constant.

Next, MRFs, have a pairwise potential, a measure of how the assigned labels are related. In our model, the pairwise potentials model how one person's gaze affects another.

$$\phi_P(l_i = l_1, l_j = l_2) = \begin{cases} c_e & \text{if } l_1 = l_2. \\ 1 - c_e & \text{otherwise.} \end{cases} \quad (5)$$

The pairwise potential provides a bias towards looking at the same position in a scene. The constant $\frac{1}{2} \leq c_e \leq 1$ biases two gazes to the same cell.

The parameters σ , c_t , c_p , c_e are learned through training images with more than one actor in the scene through cross validation. Additionally, the image space which is the set of all possible gaze points is discretized into cells, we chose 50×50 empirically for performance reasons. With more cell sizes, we can have more fine grained predictions, but at $n = 50$ and greater, we see no significant difference in performance, and we can run the full pipeline in less than a second.

To perform energy minimization for the MRF, direct inference is intractable. The number of possible labelings for n people using a grid size of 25×25 is $(25 \times 25)^n$. The complexity grows exponentially with the number of people with a high constant and a direct solver is not possible to use. We use a graph cut algorithm [17] for approximate energy minimization that requires an initial labeling for each person. The method, alpha expansion, then iterate upon this initial labeling until convergence.

3 Experiments

3.1 Dataset

GazeFollow dataset [1] is used as part of this work. The dataset contains 122,143 training images and 4,782 test images drawn from varied other datasets such as SUN [18], MS COCO {lin2014microsoft, Actions 40 [19], PASCAL [20], ImageNet challenge [14], and Places dataset [13].

Each image in training set is annotated with ground truth head position and ground truth gaze point.

Each image in test set is annotated with ground truth head position and 10 ground truth gaze points to account

for human consistency. We use average of these 10 gaze points in all further experiments.

To assess performance on joint gaze prediction, we test our model on a subset of the GazeFollow dataset that contains multiple subjects.

The CNN model [1] is trained on the GazeFollow dataset. 3500 images with multiple actors in the scene are used to train the MRF model and 1000 images are used for testing purpose.

3.2 Evaluation Metrics

We use two error metrics to quantify the performance of our method and to conduct comparative studies.

L_2 distance - The images are normalized to 1x1 and the distance between ground truth and predicted gaze point is taken. For a single image with multiple people, we report the average of the errors, that is -

$$\frac{1}{n} \sum_{i=1}^n \sqrt{(p_x^{(i)} - g_x^{(i)})^2 + (p_y^{(i)} - g_y^{(i)})^2} \quad (6)$$

Angular error - The predicted gaze point is connected with the eye position to create a unit vector v , and the ground truth gaze point is connected with the eye position to create a unit vector u . The angular error is defined as $u^T v$, and converted to degrees.

3.3 Baselines

We evaluate the performance of the following models on the GazeFollow [1] dataset. The goal of these experiments is to infer the learning carried out by the MRF model. This is achieved by varying the input applied to the MRF.

1. **Random** Random labels are assigned for each subject gaze point.
2. **MRF-Chance** The input to the MRF model is set to random instead of using the CNN predictions. This allows us to gauge how useful the CNN predictions are.
3. **MRF-Unary** The initial labels are obtained from optimizing the unary potentials. This is our implementation of [2].
4. **CNN** Equivalent to results obtained in [1].
5. **Ours** We use the predictions from [1] and use these in the MRF model from [2].

6. **Ours-Oracle** Our full model utilizing ground truth head position and orientations to quantify the models capability when provided with pristine inputs.

3.4 Quantitative Results

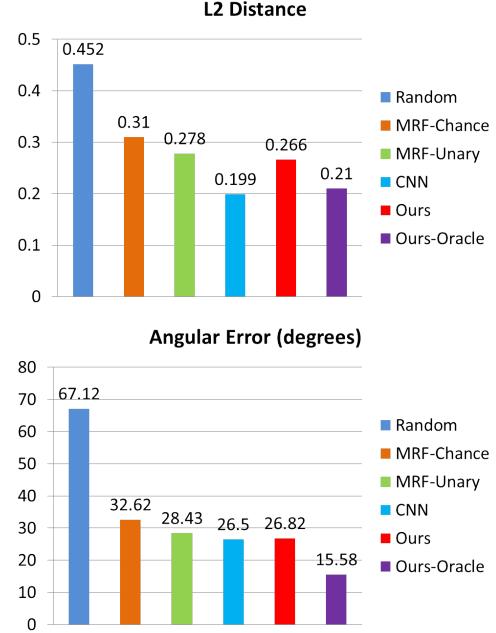


Figure 4: Quantitative comparisons between the gaze prediction models.

As expected random chance of assigning any point in the image space as the gaze point leads to high L2 and angular errors of 0.452 and 67.12 degrees respectively. The CNN model achieves a performance of 0.1991 L2 distance and 26.5 degrees angular error. This is computed by iterating through each actor in a scene for the test images.

To test the MRF model we provide it with different inputs in terms of the initial labels as well as test its sensitivity to the head position and orientation information. Comparing the results of MRF Chance (which has random initializations) and MRF Unary (initial labels computed using unary potentials) with our model (CNN predictions as initial labels) strongly indicates the performance improvement by including the CNN predictions as initial labels

and is illustrated in Fig. This also points in the direction that the MRF alone is not capable of capturing enough information from the head position and orientation only. The features learnt by CNN in terms of saliency and gaze maps capture orthogonal scene information which results in more robust predictions. Therefore, an MRF proves to be an ideal post-processing solution to account for joint gaze following predictions.

Furthermore, we compare the effect of head position and orientation accuracy on the overall performance of our full model. In one case we provide ground truth head position and orientation (computed using ground truth head position and gaze point) to our model and in another we add some gaussian random noise to the head position and orientation (computed using jittered head position and CNN prediction). For the oracle system we observe 0.21 L2 error and 15 degrees angular error compared to 0.27 and 26.8 degrees with noisy inputs. The accuracy of predictions is clearly sensitive to the accuracy of the head position and orientation values. Additionally the angular error is affected to a higher extent than L2 error, illustrating the close relationship between head orientation and gaze direction. This result is in accordance to the work in [1] wherein the authors compare the predictions when no head information is used versus no image information.

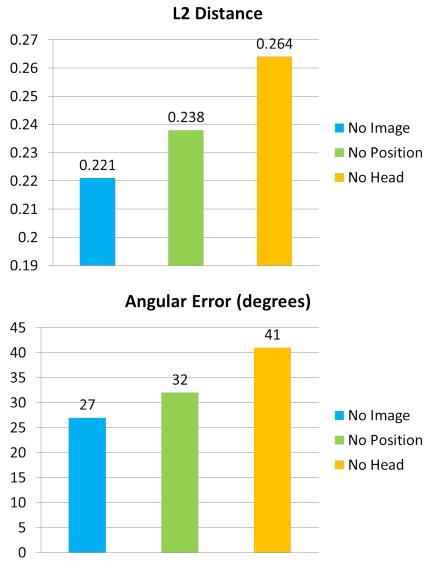


Figure 5: CNN performance given different head information [1].

Looking at the above mentioned performance results, specifically CNN versus our full model, the overall performance is similar for both cases. Thus, to investigate this further and to ascertain if the MRF model has learnt some intelligent features we divide the set of test images into two subsets, namely images with good CNN predictions and those with bad CNN predictions. A threshold on L2 and angular error is applied to populate these two sets. The corresponding performance results for these subsets are shown in Figure 5.

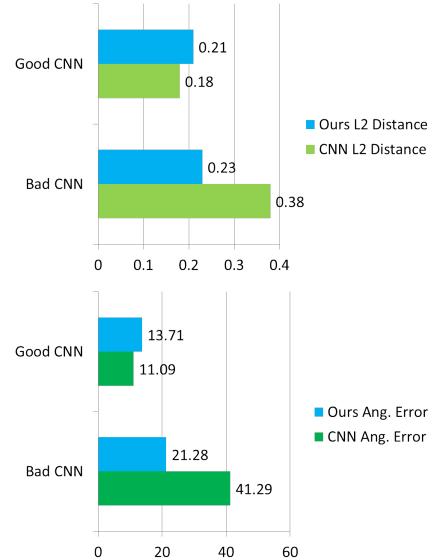


Figure 6: Results on various thresholded subsets of the data. A corresponds to the CNN results and B corresponds to our method.

The performance results on these subsets lead to the following observations. For the set with bad CNN predictions the MRF is clearly able to refine the gaze point predictions based on joint distributions. In the case of good CNN predictions the MRF model still tries to correct these even though the results from CNN model are fairly accurate with respect to the ground truth.

From the above observations we can conclude that the full model suffers when the CNN predictions are fairly accurate and therefore does not know when to stop correcting the inputs provided. This can also possibly point towards an overfitting scenario for the MRF model wherein

it always implicitly assumes the input gaze point labels to be inaccurate.

3.5 Qualitative Results

Qualitative results are presented on the last page. We can see that the success cases consist of images where people are looking at other people, which validates that the unary potentials work well. The failure cases consist of images where there is a face closely aligned to the head orientation vector. In these cases, the third unary term that biases faces would pull the predicted gaze closer to the face.

4 Conclusion

We explored an approach to combine and utilize two state of the art methods, GazeFollow model [1] and Social Attention [2] for gaze following task specifically in the case of multiple actors in a scene. In the MRF model unary and pairwise potentials are formulated to account for head pose of an actor and relationship between actors respectively in a scene.

From the experiments conducted it is clear that the MRF model directly benefits from initial labels sourced from CNN predictions. Additionally, our model is more sensitive to the head position and orientation accuracy, similar to results obtained in [1] (shown in Fig 4.)

4.1 Further Extensions

In performance analysis we noticed that most gazes are overshooting the ground truth gaze point, and we have pinpointed the cause to one of the unary potential terms. The potential function ϕ_1 seeks to align the gaze prediction point with the head orientation will always choose to pick a further point along a given gaze vector due to the discretization of the grid. Further points will be more aligned with the head orientation and therefore have a higher score. Future work would include exploring robust potential functions such as Gaussian distribution of gaze vector length.

An important direction to go in for further enhancing the method is to allow the model to dynamically decide between further refining the input CNN predictions and

going back to the given input if further refinement is deleterious. Additionally, a holistic approach would be to integrate joint prediction mechanism with the CNN model from [1], which would result in a more robust and accurate network than combining two different techniques.

References

- [1] A. Recasens*, A. Khosla*, C. Vondrick, and A. Torralba, “Where are they looking?” in *Advances in Neural Information Processing Systems (NIPS)*, 2015, * indicates equal contribution. 1, 2, 3, 4, 5, 6, 7
- [2] A. Fathi, J. K. Hodgins, and J. M. Rehg, “Social interactions: A first-person perspective,” in *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE, 2012, pp. 1226–1233. 1, 2, 3, 5, 7
- [3] M. J. Marín-Jiménez, A. Zisserman, M. Eichner, and V. Ferrari, “Detecting people looking at each other in videos,” *International Journal of Computer Vision*, vol. 106, no. 3, pp. 282–296, 2014. 1, 2
- [4] H. Soo Park and J. Shi, “Social saliency prediction,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 4777–4785. 1
- [5] H. Soo Park, E. Jain, and Y. Sheikh, “Predicting primary gaze behavior using social saliency fields,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2013, pp. 3503–3510. 1
- [6] D. Parks, A. Borji, and L. Itti, “Augmented saliency model using automatic 3d head pose detection and learned gaze following in natural scenes,” *Vision research*, vol. 116, pp. 113–126, 2015. 2
- [7] R. E. Flom, K. E. Lee, and D. E. Muir, *Gaze-following: Its development and significance*. Lawrence Erlbaum Associates Publishers, 2007. 2
- [8] N. J. Emery, “The eyes have it: the neuroethology, function and evolution of social gaze,” *Neuroscience & Biobehavioral Reviews*, vol. 24, no. 6, pp. 581–604, 2000. 2
- [9] R. Newman, Y. Matsumoto, S. Rougeaux, and A. Zelinsky, “Real-time stereo tracking for head pose and gaze estimation,” in *Automatic Face and Gesture Recognition, 2000. Proceedings. Fourth IEEE International Conference on*. IEEE, 2000, pp. 122–128. 2
- [10] A. Garcia-Diaz, X. R. Fdez-Vidal, X. M. Pardo, and R. Dosil, “Saliency from hierarchical adaptation through decorrelation and variance normalization,” *Image and Vision Computing*, vol. 30, no. 1, pp. 51–64, 2012. 2

- [11] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, “Caffe: Convolutional architecture for fast feature embedding,” in *Proceedings of the 22nd ACM international conference on Multimedia*. ACM, 2014, pp. 675–678. [3](#)
- [12] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105. [3](#)
- [13] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva, “Learning deep features for scene recognition using places database,” in *Advances in neural information processing systems*, 2014, pp. 487–495. [3, 4](#)
- [14] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, “Imagenet large scale visual recognition challenge,” *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015. [3, 4](#)
- [15] C. Wang, N. Komodakis, and N. Paragios, “Markov random field modeling, inference & learning in computer vision & image understanding: A survey,” *Computer Vision and Image Understanding*, vol. 117, no. 11, pp. 1610–1627, 2013. [3](#)
- [16] A. Quattoni, S. Wang, L.-P. Morency, M. Collins, and T. Darrell, “Hidden conditional random fields,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 29, no. 10, pp. 1848–1852, 2007. [3](#)
- [17] Y. Boykov, O. Veksler, and R. Zabih, “Fast approximate energy minimization via graph cuts,” *IEEE Transactions on pattern analysis and machine intelligence*, vol. 23, no. 11, pp. 1222–1239, 2001. [4](#)
- [18] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba, “Sun database: Large-scale scene recognition from abbey to zoo,” in *Computer vision and pattern recognition (CVPR), 2010 IEEE conference on*. IEEE, 2010, pp. 3485–3492. [4](#)
- [19] B. Yao, X. Jiang, A. Khosla, A. L. Lin, L. Guibas, and L. Fei-Fei, “Human action recognition by learning bases of action attributes and parts,” in *2011 International Conference on Computer Vision*. IEEE, 2011, pp. 1331–1338. [4](#)
- [20] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, “The pascal visual object classes (voc) challenge,” *International journal of computer vision*, vol. 88, no. 2, pp. 303–338, 2010. [4](#)



Table 1: Ground truth plotted in red. CNN predictions plotted in blue. Ours plotted in green. The top two rows show success cases of capturing mutual gazes. The bottom two rows show failure cases.