

Appearance-based Gaze Estimation with Deep Learning: A Review and Benchmark

Yihua Cheng¹, Haofei Wang², Yiwei Bao¹, Feng Lu^{1,2,*}

¹State Key Laboratory of Virtual Reality Technology and Systems, SCSE, Beihang University, China.

²Peng Cheng Laboratory, Shenzhen, China.

{yihua_c, baoyiwei, lufeng}@buaa.edu.cn, wanghf@pcl.ac.cn

Abstract—Human gaze provides valuable information on human focus and intentions, making it a crucial area of research. Recently, deep learning has revolutionized appearance-based gaze estimation. However, due to the unique features of gaze estimation research, such as the unfair comparison between 2D gaze positions and 3D gaze vectors and the different pre-processing and post-processing methods, there is a lack of a definitive guideline for developing deep learning-based gaze estimation algorithms. In this paper, we present a systematic review of the appearance-based gaze estimation methods using deep learning. Firstly, we survey the existing gaze estimation algorithms along the typical gaze estimation pipeline: **deep feature extraction, deep learning model design, personal calibration and platforms**. Secondly, to fairly compare the performance of different approaches, we summarize the data pre-processing and post-processing methods, including face/eye detection, data rectification, 2D/3D gaze conversion and gaze origin conversion. Finally, we set up a comprehensive benchmark for deep learning-based gaze estimation. We characterize all the public datasets and provide the source code of typical gaze estimation algorithms. This paper serves not only as a reference to develop deep learning-based gaze estimation methods, but also a guideline for future gaze estimation research. The project web page can be found at <https://phi-ai.buaa.edu.cn/Gazehub/>.

Index Terms—gaze estimation, eye appearance, deep learning, review, benchmark.

1 INTRODUCTION

Eye gaze is an essential non-verbal communication cue that contains valuable information about human intent, enabling us to gain insights into human cognition [1, 2] and behavior [3, 4]. Eye gaze has various representations across different applications. Gaze direction serves as the universal representation in most applications. It is defined as a unit direction vector in 3D space originating from eye centers and pointing towards gaze targets. Gaze direction holds significant potential, *e.g.*, in extended reality (XR) devices [5–7], where it is employed to locate gaze targets in 3D space based on estimated gaze direction. By establishing a specific plane in 3D space, gaze direction can be converted into a point of gaze (PoG) on that plane. PoG is widely used in human-computer interaction [8–10] as it indicates the user’s attention area on a screen or display. Additionally, eye gaze can be represented as an attention map in analysis tasks [11, 12] or as target objects/people in gaze following tasks [13–15]. Accurate gaze estimation is always crucial for such applications.

Over the last decades, numerous gaze estimation methods have been proposed. These methods can be broadly categorized into three groups: 3D eye model recovery-based, 2D eye feature regression-based, and appearance-based methods. 3D eye model recovery-based methods construct a geometric 3D eye model and estimate gaze directions based on the model. Due to the diversity of human eyes, 3D eye models are usually person-specific. The 3D eye model recovery-based methods usually require personal

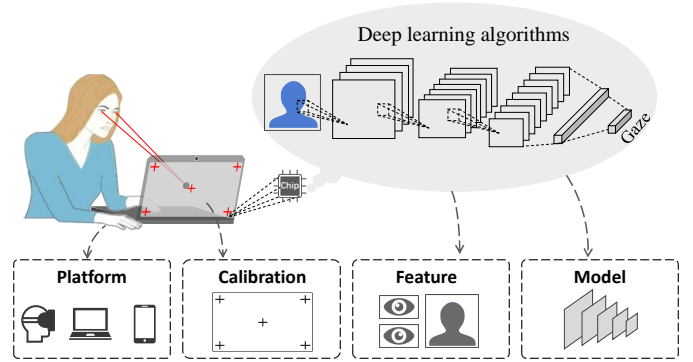


Fig. 1. Deep learning-based gaze estimation relies on simple devices but complex algorithms to estimate human gaze. It usually uses off-the-shelf cameras to capture facial appearance, and employs deep learning algorithms to regress gaze from the appearance. According to this pipeline, we survey current deep learning-based gaze estimation methods from four perspectives: deep feature extraction, deep learning model design, personal calibration, and platforms.

calibration to recover person-specific parameters such as iris radius and kappa angle. While these methods often achieve high accuracy, they require dedicated devices such as infrared cameras. 2D eye feature regression-based methods usually keep the same requirement on devices as 3D eye model recovery-based methods. They directly use detected geometric eye feature such as pupil center to regress the point of gaze (PoG). They do not require geometric calibration for converting gaze directions into PoG.

Appearance-based methods have low device requirements. They use off-the-shelf web cameras to capture human eye appear-

• Feng Lu is the Corresponding Author.

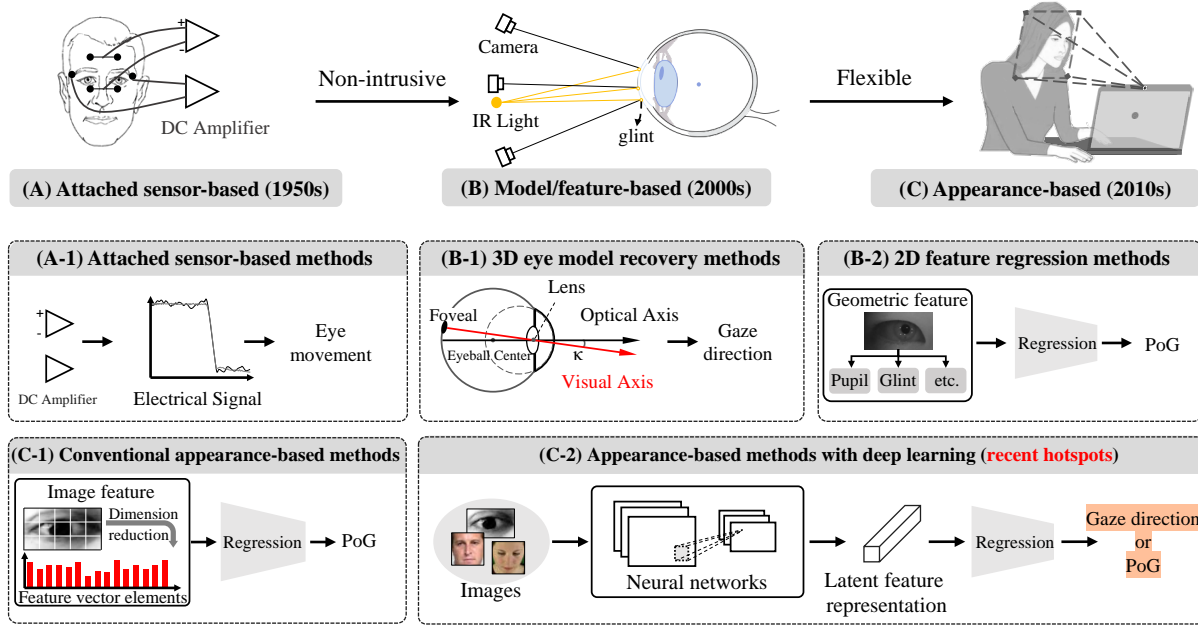


Fig. 2. From intrusive skin electrodes [16] to off-shelf web cameras [17], gaze estimation is more flexible. Gaze estimation methods are also updated with the change of devices. We illustrate five kinds of gaze estimation methods. (1). Attached sensor-based methods. The method samples the electrical signal of skin electrodes. The signal indicates the eye movement of subjects [18]. (2) 3D eye model recovery methods. The method usually builds a geometric eye model to calculate the visual axis, *i.e.*, gaze directions. The eye model is fitted based on the light reflection. (3) 2D eye feature regression methods. The method relies on IR cameras to detect geometric eye features such as pupil center, glints, and directly regress the PoG from these features. (4) Conventional appearance-based methods. The method use entire images as feature and directly regress human gaze from features. Some feature reduction methods are also used for extracting low-dimensional feature. For example, Lu *et al.* divide eye images into 15 subregion and sum the pixel intensities in each subregion as feature [19]. (5) Appearance-based gaze estimation with deep learning, which is the recent hotspots. Face or eye images are directly inputted into a designed neural network to learn latent feature representation, and human gaze is regressed from the feature representation.

ance and regress gaze from the appearance. Although the setup is simple, they have strict requirements on the gaze estimation algorithm. They usually require 1) An effective feature extractor to extract gaze features from high-dimensional raw images. Some feature extractors such as histograms of oriented gradients are used in the conventional method [20]. 2) A robust regression function to learn the mappings from appearance feature to human gaze. It is non-trivial to map the high-dimensional eye appearance to the low-dimensional gaze. Many regression functions have been used to regress gaze from appearance, *e.g.*, local linear interpolation [21] and adaptive linear regression [19]. 3) A large number of training samples to learn the regression function. They usually collect personal samples with a time-consuming personal calibration, and learn a person-specific gaze estimation model. Some studies seek to reduce the number of training samples [19].

Recently, deep learning-based methods have gained popularity as they offer several advantages over conventional appearance-based methods. These methods use convolution layers or transformers [22] to automatically extract high-level gaze features from images. Deep learning models are also highly non-linear and can fit the mapping function from eye appearance to gaze direction even with large head motion. These advantages make deep learning-based methods more accurate and robust than conventional methods. Deep learning-based methods also improve cross-subject gaze estimation performance significantly, reducing the need for time-consuming person calibration. These improvements expand the application range of appearance-based gaze estimation.

In this paper, we provide a systematic review of appearance-based gaze estimation methods using deep learning algorithms. As shown in Fig. 1, we discuss these methods from four perspectives: 1) deep feature extraction, 2) deep neural network architecture

design, 3) personal calibration, and 4) device and platform. From the deep feature extraction perspective, we describe the strategies for extracting features from eye images, face images and videos. Under the deep neural network architecture design perspective, we first review methods based on the supervised strategy, containing the supervised, self-supervised, semi-supervised and unsupervised methods. Then, We describe different deep neural networks in gaze estimation including multi-task CNNs, recurrent CNNs. Furthermore, we introduce methods that integrate CNN models and prior knowledge of gaze. From the personal calibration perspective, we describe how to use calibration samples to further improve the performance of CNNs. We also introduce the method integrating user-unaware calibration sample collection mechanism. Finally, from the device and platforms perspective, we consider different cameras, *i.e.*, RGB cameras, IR cameras and depth cameras, as well as different platforms, *i.e.*, computers, mobile devices and head-mount devices. We review the advanced methods using these cameras and proposed for these platforms.

Besides deep learning-based gaze estimation methods, we also summarize the practices of gaze estimation. We first review the data pre-processing methods of gaze estimation including face and eye detection methods and data rectification methods. Then, considering various forms of human gaze, *e.g.*, gaze direction and PoG, we further provide data post-processing methods. These methods describe the geometric conversion between various representations of human gaze. We also build gaze estimation benchmarks. We collect and implement the codes of typical gaze estimation methods, and evaluate them on various datasets. For the different kinds of gaze estimation methods, we convert their result for fair comparisons with data post-processing methods. Our benchmarks provide comprehensive comparison between state-of-

the-art gaze estimation methods.

The paper is organized as follows. Section 2 introduces the background of gaze estimation. We introduce the development and category of gaze estimation methods. Section 3 reviews the state-of-the-art deep learning-based method. In Section 4, we introduce the public datasets as well as data pre-processing and post-processing methods. We also build the benchmark in this section. In Section 5, we conclude the development of current deep learning-based methods and recommend future research directions. This paper can not only serve as a reference to develop deep learning-based gaze estimation methods, but also a guideline for future gaze estimation research.

2 GAZE ESTIMATION BACKGROUND

2.1 Categorization

Figure 2 illustrates the development of gaze estimation methods. Early gaze estimation methods detect eye movement patterns such as fixation, saccade and smooth pursuit [16]. They attach sensors around eyes and measure eye movement using potential differences [23, 24]. With the development of computer vision technology, modern eye-tracking devices have emerged. They usually estimate gaze using eye/face images captured by cameras. In general, there are two types of such devices, remote eye tracker and head-mounted eye tracker. The remote eye tracker usually keeps a certain distance from the user, *e.g.*, ~ 60 cm. The head-mounted eye tracker usually mounts the cameras on a frame of glasses. Compared to the intrusive eye tracking devices, the modern eye tracker greatly enlarges the range of application.

Computer vision-based methods can be further divided into three types: 2D eye feature regression methods, 3D eye model recovery methods and appearance-based methods. The first two types of methods estimate gaze based on geometric features such as contours, reflection and eye corners. The geometric features can be accurately extracted with the assistance of dedicated devices, *e.g.*, infrared cameras. More concretely, the 2D eye feature regression method learns a mapping function from geometric feature to point of gaze, *e.g.*, the polynomials [25, 26] and the neural networks [27]. The 3D eye model recovery method builds subject-specific geometric eye models to estimate human gaze directions. The eye model is fitted with geometric features, such as the infrared corneal reflections [28, 29], pupil center [30] and iris contours [31]. However, they usually require a personal calibration process for each subject, since the eye model contains subject-specific parameters such as cornea radius, kappa angles.

Appearance-based methods directly learn a mapping function from images to human gaze. Different from previous methods, appearance-based methods do not require dedicated devices for detecting geometric features. They use image features such as image pixel [19] or deep features [17] to regress gaze. Various regression models have been used, *e.g.*, neural networks [32], gaussian process regression [33], adaptive linear regression [19], convolutional neural networks [17] and transformers [34]. However, it is still a challenging task due to complex facial appearance.

2.2 Appearance-based Gaze Estimation

Appearance-based methods directly learn mapping function from eye appearance to human gaze. As early as 1994, Baluja *et al.* propose a neural network and collect 2,000 samples for training [32]. Tan *et al.* use a linear function to interpolate

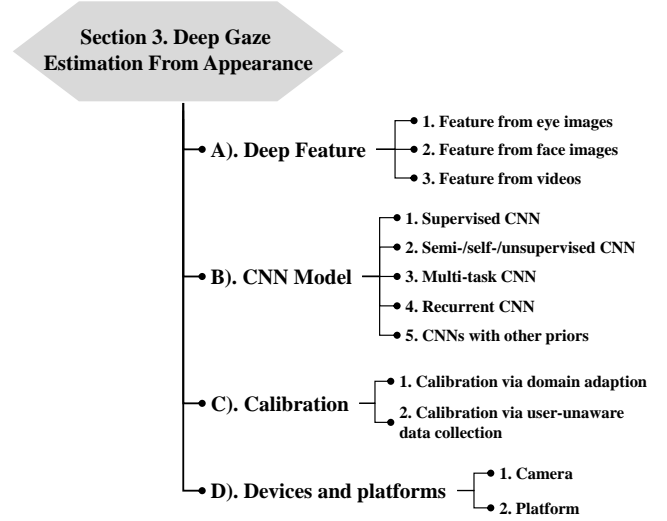


Fig. 3. The architecture of section 3. We introduce gaze estimation with deep learning from four perspectives.

unknown gaze position using 252 training samples [21]. These methods usually learn a subject-specific mapping function. They require a time-consuming data collection for the specific subject. To reduce the number of training samples, Williams *et al.* introduce semi-supervised gaussian process regression methods [33]. Sugano *et al.* propose a method that combines gaze estimation with saliency [35]. Lu *et al.* propose an adaptive linear regression method to select an optimal set of sparsest training sample for interpolation [19]. However, these methods only show reasonable performance in a constrained environment, *i.e.*, fixed head pose and the specific subject. Their performance significantly degrades when tested on an unconstrained environment. This problem is always challenging in appearance-based gaze estimation.

To address the performance degradation across subjects, Funes *et al.* present a cross-subject training method [36]. However, the reported mean error is larger than 10 degrees. Sugano *et al.* introduce a learning-by-synthesis method [37]. They use a large number of synthetic cross-subject data to train their model. Lu *et al.* employ a sparse auto-encoder to learn a set of bases from eye image patches and reconstruct the eye image using these bases [38]. On the other hand, to tackle the head motion problem, Sugano *et al.* cluster the training samples with similar head poses and interpolate the gaze in local manifold [39]. Lu *et al.* initiate the estimation with the original training images and compensating for the bias via regression [40]. They further propose a novel gaze estimation method that handles the free head motion via eye image synthesis using a single camera [41].

2.3 Deep Learning for Gaze Estimation

Appearance-based gaze estimation suffers from many challenges, including head motion and subject differences, particularly in the unconstrained environment. Traditional appearance-based methods often struggle to effectively address these challenges due to their limited fitting ability. Deep learning have been used in many computer vision tasks and demonstrated outstanding performance. Zhang *et al.* propose the first CNN-based gaze estimation method to regress gaze directions from eye images [17]. They use a simple CNN and the performance surpasses most of the conventional appearance-based approaches. Following this study, an increasing

number of improvements and extensions on CNN-based gaze estimation methods emerged. Face images [42] and videos [43] have also been used for gaze estimation. These inputs provide more valuable information than using eye images alone. Some methods are proposed for handling the challenges in an unconstrained environment. For example, Cheng *et al.* use asymmetric regression to handle the extreme head pose and illumination condition [44]. Park *et al.* learn a pictorial eye representation to alleviate the personal appearance difference [45]. The calibration-based methods learn a subject-specific CNN model [46, 47]. Xu *et al.* investigated the vulnerability of appearance-based gaze estimation [48].

3 DEEP GAZE ESTIMATION FROM APPEARANCE

We survey deep learning-based gaze estimation methods in this section. We introduce these methods from four perspectives, deep feature extraction, deep neural network architecture design, personal calibration as well as device and platform. Figure 3 gives an overview of this section.

3.1 Deep Feature from Appearance

Feature extraction plays a crucial role in most of the learning-based tasks. It is challenging to effectively extract features from complex eye appearance due to identity, illumination and *etc.* The quality of the extracted features determines the gaze estimation accuracy. In this section, we summarize feature extraction mechanisms according to the types of input into the deep neural network, including eye images, face images and videos.

3.1.1 Feature from Eye Images

Human gaze has a strong correlation with eye appearance. Even a minor perturbation in gaze direction can result in noticeable changes in eye appearance. For instance, when the eyeball rotates, the position of the iris and the shape of the eyelid undergo alterations, leading to corresponding changes in gaze direction. This relationship between gaze and eye appearance enables the gaze estimation based on the visual feature of eyes. Conventional methods typically estimate gaze using high-dimensional raw image features [21, 51]. These features are obtained by raster scanning all the pixels in eye images, resulting in a representation that contains a significant amount of redundancy. Moreover, these features are highly sensitive to environmental changes, which can pose challenges in achieving accurate gaze estimation.

Deep learning-based methods automatically extract deep features from eye images. Zhang *et al.* propose the first deep learning-based gaze estimation method [17]. They employ a CNN to extract features from grey-scale eye images and concatenate the features with head pose. As with most deep learning tasks, the deeper network structure and larger receptive field, the more informative features can be extracted. Zhang *et al.* [49] further extend their previous work [17] and present a GazeNet which is inherited from a 16-layer VGG network [52]. Chen *et al.* [53] use dilated convolutions to extract high-level eye features, which efficiently increases the receptive field size of the convolutional filters without reducing spatial resolution.

Recent studies found that concatenating the features of two eyes helps to improve the gaze estimation accuracy [54, 55]. Fischer *et al.* [54] employ two VGG-16 networks to extract individual features from two eye images, and concatenate two eye features for regression. Cheng *et al.* [55] build a four-stream CNN network for extracting features from two eye images.

Two streams of CNN are used for extracting individual features from left/right eye images, the other two streams are used for extracting joint features of two eye images. They claim that the two eyes are asymmetric, and propose an asymmetric regression and evaluation network to extract different features from two eyes. More recent studies propose to use attention mechanism to fuse two eye features. Cheng *et al.* [56] argue that the weights of two eye features are determined by face images due to the specific task in [56], so they assign weights for two eye features with the guidance of facial features. Bao *et al.* [57] propose a self-attention mechanism to fuse two eye features. They concatenate the feature maps of two eyes and use a convolution layer to generate the weights of the feature map. Murthy *et al.* [58] simultaneously estimate feature vectors and weights for each eye image and concatenate the left and right eye feature. They also propose a network which obtains the difference between left and right eye feature to circumvent person-dependent features.

Above methods extract the general features from eye images while other works explored extracting specific features to handle the head motion and subject difference. Several studies have attempted to extract subject-invariant features from eye images [45, 47, 59]. Park *et al.* [45] convert the original eye images into a unified gaze representation, which is a pictorial representation of eyeball, iris and pupils. They regress gaze directions from the pictorial representation. Wang *et al.* propose an adversarial learning approach to extract the domain/person-invariant feature [59]. They feed the features into an additional classifier and design an adversarial loss function to handle the appearance variations. Park *et al.* use an autoencoder to learn the compact latent representation of gaze, head pose and appearance [47]. They introduce a geometric constraint on gaze representations, *i.e.*, the rotation matrix between the two given images transforms the gaze representation of one image to another. In addition, some methods use generative adversarial networks (GAN) to pre-process eye images to handle specific environment factors. Kim *et al.* [60] convert low-light eye images into bright eye images. Rangesh *et al.* [61] use a GAN to remove eyeglasses.

Besides the supervised approaches for extracting gaze features, unannotated eye images have also been used for learning gaze representations. Yu *et al.* input the difference of gaze representations from two eyes into pre-trained network for gaze redirection [62]. They learn 2-D representations from unannotated eye images which can be seemed as unaligned gaze. Sun *et al.* propose a cross encoder to disentangle gaze feature and appearance feature. They improve the few-shot performance using learned gaze feature.

3.1.2 Feature from Face Images

Face images contain head pose information that also contributes to gaze estimation. Conventional methods extract features such as head pose [41] and facial landmarks [63–65] from face images. Eye image-based methods typically use head pose vectors as an additional input [17, 55]. Nevertheless, the impact of head pose appears to be marginal [49], particularly when the basic network has already achieved high accuracy. One possible rationale for this observation lies in the fact that a single head pose often corresponds to a broad gaze range [66, 67], thereby providing only a coarse indication of gaze direction rather than precise gaze information. Deep facial feature performs better than the head pose. Recent studies directly use face images as input and employ a CNN to extract deep facial features [42, 50, 68, 69] as shown in Fig. 4 (b). It demonstrates an improved performance than the

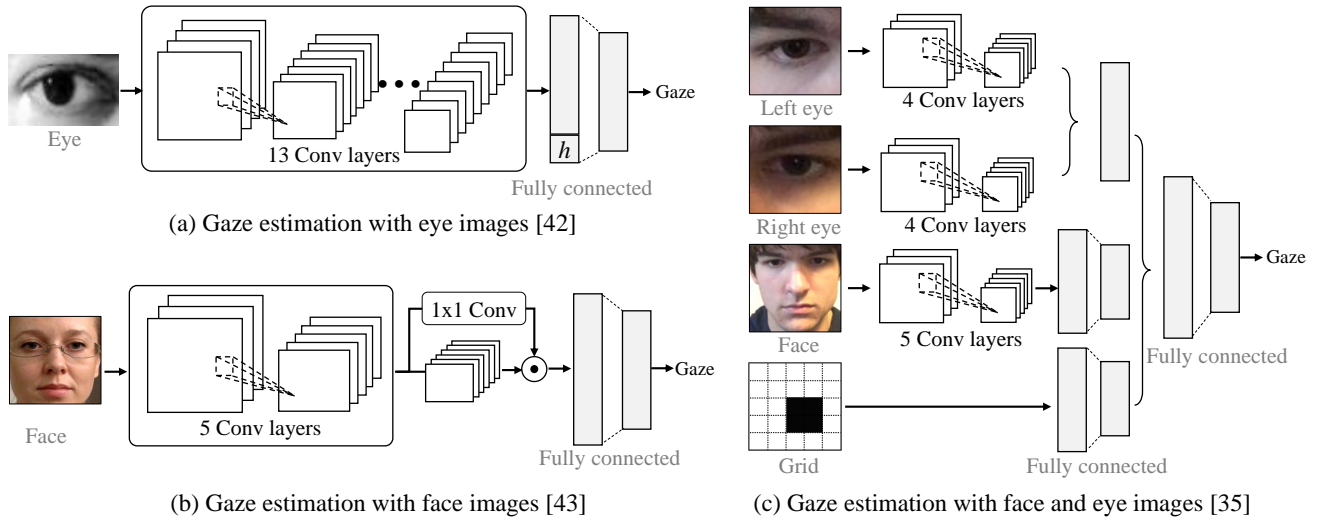


Fig. 4. Some typical CNN-based gaze estimation networks. (a). Gaze estimation with eye images [49]. (b) Gaze estimation with face images [50]. (c). Gaze estimation with face and eye images [42].

approaches that only use eye images. Cheng *et al.* [34] explore the transformer for gaze estimation. They use CNN to extract feature maps from face images and input the feature map into transformer encoder for gaze estimation.

Face images contain redundant information. Researchers have attempted to filter out the useless features in face image [50, 70]. Zhang *et al.* [50] propose a spatial weighting mechanism to efficiently encode the location of the face into a standard CNN architecture. The system learns spatial weights based on the activation maps of the convolutional layers. This helps to suppress the noise and enhance the contribution of the highly activated regions. Zhang *et al.* [71] propose a learning-based region selection method by dynamically selecting suitable sub-regions from facial region. Cheng *et al.* [72] propose a plug-and-play self-adversarial network to purify facial features. They remove gaze-irrelevant image features while preserve gaze-relevant features, so the robustness of gaze estimation network has been improved.

Some studies crop the eye image out of the face images and directly feed it into the network. These works usually use a three-stream network to extract features from face images, left and right eye images, respectively as shown in Fig. 4 (c) [42, 53, 73–75]. Besides, Deng *et al.* [76] decompose gaze directions into the head rotation and eyeball rotation. They use face images to estimate the head rotation and eye images to estimate the eyeball rotation. These two rotations are aggregated into a gaze vector through a gaze transformation layer. Cheng *et al.* [56] propose a coarse-to-fine gaze estimation method. They use face feature to estimate basic gaze directions, then refine the basic gaze direction with eye features. They use GRU [77] to build the network. Cai *et al.* [78] use a transformer encoder [22] to aggregate face and eye features. They feed face and two eye features into the transformer encoder and concatenate the outputs of the encoder for gaze estimation.

Facial landmarks have also been used as additional features to model the head pose and eye position. Palmero *et al.* combine individual streams (face, eyes region and face landmarks) in a CNN [79]. Dias *et al.* extract the facial landmarks and directly regress gaze from the landmarks [80, 81]. The network outputs the gaze direction as well as an estimation of its own prediction uncertainty. Jyoti *et al.* further extract geometric features from the facial landmark locations [82]. The geometric feature includes

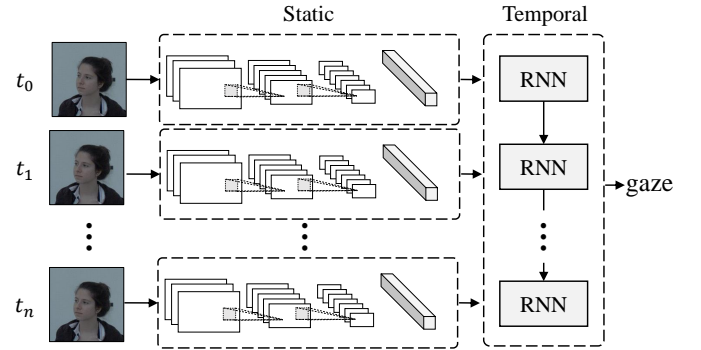


Fig. 5. Gaze estimation with videos. It first extracts static features from each frame using a typical CNN, and feeds these static features into RNN for extracting temporal information.

the angles between the pupil center as the reference point and the facial landmarks of the eyes and the tip of the nose. The detected facial landmarks can also be used for unsupervised gaze representation learning. Dubey *et al.* [83] collect the face images from the web and annotate their gaze zone based on the detected landmarks. They perform gaze zone classification tasks on the dataset for unsupervised gaze representation learning. In addition, since the cropped face image does not contain face position information, Krafka *et al.* [42] propose the iTracker, combining the information from left/right eye images, face images as well as face grid information. The face grid indicates the position of the face region in images and it is usually used in PoG estimation.

3.1.3 Feature from Videos

Temporal information from videos also contributes to better gaze estimates. Recurrent Neural Network (RNN) has been widely used in video processing, *e.g.*, long short-term memory (LSTM) [43, 84]. As shown in Fig. 5, they usually use a CNN to extract features from face images at each frame, and then input these features into a RNN. The temporal relations between each frame are automatically captured by the RNN for gaze estimation.

Temporal features such as the optical flow and eye movement dynamics have been used to improve gaze estimation accuracy. The optical flow provides the motion information between the frames. Wang *et al.* [85] use the optical flow constraints with

2D facial features to reconstruct the 3D face structure based on the input video frames. Eye movement dynamics, such as fixation, saccade and smooth pursuits, have also been used to improve gaze estimation accuracy. Wang *et al.* [86] propose to leverage eye movement to generalize eye tracking algorithm to new subjects. They use a dynamic gaze transition network to capture underlying eye movement dynamics and serve as prior knowledge. They also propose another static gaze estimation network, which estimates gaze based on the static frame. They finally combine the two networks for better gaze estimation accuracy. The combination method of the two networks is solved as a standard inference problem of linear dynamic system or Kalman filter [87].

3.2 CNN Models

Convolutional neural networks have been widely used in many computer vision tasks [88]. They also demonstrate superior performance in the field of gaze estimation. In this section, we first review the existing gaze estimation methods from the learning strategy perspective, *i.e.*, the supervised CNNs and the semi-/self-/un-supervised CNNs. Then we introduce the different network architectures, *i.e.*, multi-task CNNs and the recurrent CNNs for gaze estimation. In the last part of this section, we discuss the CNNs that integrate prior knowledge to improve performance.

3.2.1 Supervised CNNs

Supervised CNNs are the most commonly used networks in appearance-based gaze estimation [17, 89–91]. Fig. 4 shows the typical architecture of supervised gaze estimation CNN. The network is trained using image samples with ground truth gaze directions. The gaze estimation problem is essentially learning a mapping function from raw images to human gaze. Therefore, similar to other computer vision tasks [92], the deeper CNN architecture usually achieves better performance. A number of CNN architectures that have been proposed for typical computer vision tasks also show great success in gaze estimation task, *e.g.*, LeNet [17], AlexNet [50], VGG [49], ResNet18 [43] and ResNet50 [66]. Besides, some well-designed modules also help to improve the estimation accuracy [53, 56, 93, 94]. Chen *et al.* use a dilated convolution to extract features from eye images [53]. Cheng *et al.* propose an attention module for fusing two eye features [56]. Cheng *et al.* integrate the CNN and the transformer encoder [22] to improve the estimation performance [34].

To supervise the CNN during training, the system requires the large-scale labeled dataset. Several large-scale datasets have been proposed [17, 42]. However, it is difficult and time-consuming to collect enough gaze data in practical applications. Inspired by the physiological eye model [95], some researchers propose to synthesize labeled photo-realistic image [37, 96, 97]. These methods usually build eye-region models and render new images from these models. One of such methods is proposed by Sugano *et al.* [37]. They synthesize dense multi-view eye images by recovering the 3D shape of eye regions, where they use a patch-based multi-view stereo algorithm [98] to reconstruct the 3D shape from eight multi-view images. Wood *et al.* propose to synthesize the close-up eye images for a wide range of head poses, gaze directions and illuminations to develop a robust gaze estimation algorithm [99]. Following this work, Wood *et al.* further propose another system named UnityEye to rapidly synthesize large amounts of eye images of various eye regions [100]. To make the synthesized images more realistic, Shrivastava *et al.* propose

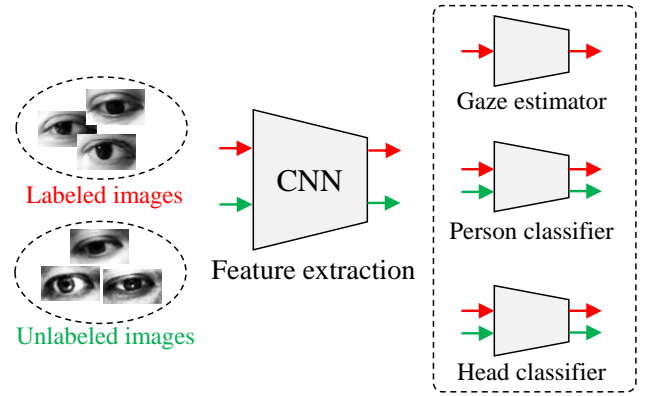


Fig. 6. A semi-supervised CNN [59]. It uses both labeled images and unlabeled images for training. It designs an extra appearance classifier and a head pose classifier. The two classifiers align the feature of labeled images and unlabeled images.

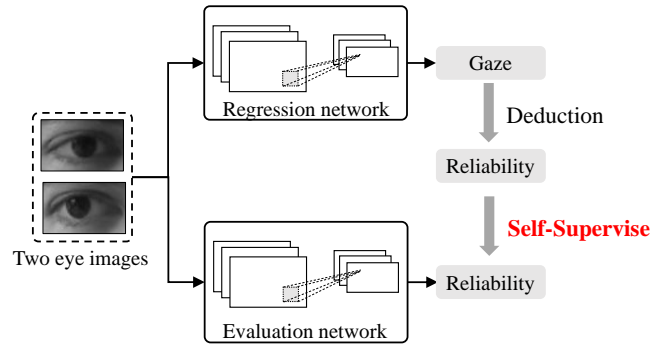


Fig. 7. A self-supervised CNN [55]. The network is consisted of two sub-networks. The regression network estimates gaze from two eye images and generates the ground truth of the other network for self-supervision.

an unsupervised learning paradigm using generative adversarial networks to improve the realism of the synthetic images [101]. Wang *et al.* plot eye shapes based on geometric models and use GAN to render eye images [102]. These methods serve as data augmentation tools to improve the performance of gaze estimation.

Gaze redirection has also been used as a data augmentation tool. It generates face images with target gaze based on given face images. Recently, many gaze redirection methods have been proposed [103–106] and bring significant performance improvement. NeRF [107] shows great multi-view consistency and is used to learn implicit face model from multi-view images. It can also renders face images under novel gaze for gaze redirection [108, 109].

3.2.2 Semi-/Self-/Un-supervised CNNs

Semi-/self-/un-supervised CNNs attract much attention recently and also show large potential in gaze estimation. There are typically two main topics in recent research. 1) Gaze data collection is time-consuming and expensive. To reduce the requirement on annotated images, some methods leverage unannotated images to learn robust feature representation [62, 110]. 2) Gaze estimation methods show performance drop in new environments/domains. Researchers use annotated images in source domains and unannotated images in target domains to improve the performance in target domains [111, 112]. The second topic is more systematic than the first topic with recent development. It is defined as unsupervised domain adaption, where the “unsupervised” aspect refers to the lack of labelled data in the target domain.

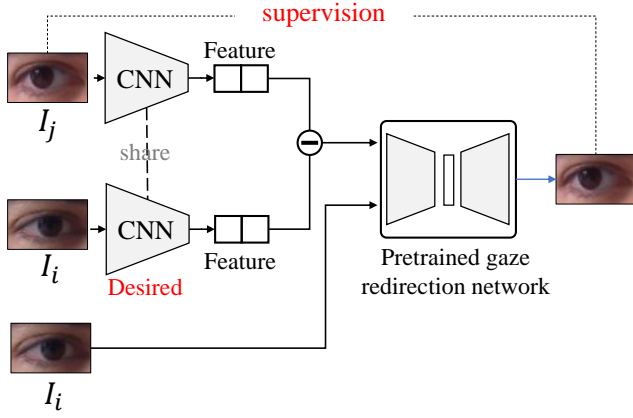


Fig. 8. An unsupervised CNN [62]. It extracts 2D feature from eye images. The feature difference and one eye image are fed into a pretrained gaze redirection network to generate the other eye image.

Semi-supervised CNNs require both labeled and unlabeled images for optimizing networks. Wang *et al.* propose an adversarial learning approach to improve the model performance on the target subject/dataset [59]. As shown in Fig. 6, it requires labeled images in the training set as well as unlabeled images of the target subject/dataset. They use the labeled data to supervise the gaze estimation network and design an adversarial module for semi-supervised learning. Given these features used for gaze estimation, the adversarial module tries to distinguish their source and the gaze estimation network aims to extract subject/dataset-invariant features to cheat the module. Kothari *et al.* [110] found the strong gaze-related geometric constraints when people “look at each other” (LAEO). They estimate 3D and 2D landmarks in the images of LAEO dataset [113], and generate pseudo gaze annotation for gaze estimation. While it cannot bring competitive performance, therefore, they further integrate labeled images and LAEO datasets for semi-supervised gaze estimation.

Self-supervised CNNs aim to formulate a pretext auxiliary learning task to improve the estimation performance. Cheng *et al.* propose a self-supervised asymmetry regression network for gaze estimation [55]. As shown in Fig. 7, the network consists of a regression network to estimate the two eyes’ gaze directions and an evaluation network to assess the reliability of two eyes. During training, the result of the regression network is used to supervise the evaluation network, the accuracy of the evaluation network determines the learning rate in the regression network. They simultaneously train the two networks and improve the regression performance without additional inference parameters. Xiong *et al.* introduce a random effect parameter to learn the person-specific information in gaze estimation [114]. They utilize the variational expectation-maximization algorithm [115] and stochastic gradient descent [116] to estimate the parameters of the random effect network during training. They use another network to predict the random effect based on the feature representation of eye images. The self-supervised strategy predicts the random effects to enhance the accuracy for unseen subjects. He *et al.* introduce a person-specific user embedding mechanism [117]. They concatenate the user embedding with appearance features to estimate gaze. They also build a teacher-student network, where the teacher network optimizes the user embedding during training and the student network learns the user embedding from the teacher network.

Unsupervised CNNs only require unlabeled data for training. Nevertheless, it is hard to optimize CNNs without the ground

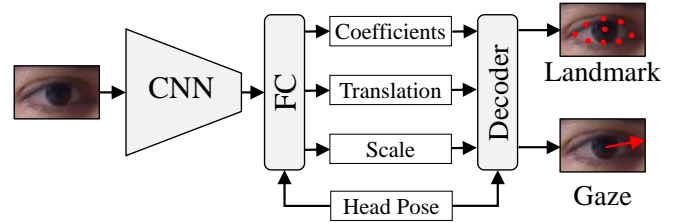


Fig. 9. A multitask CNN [119]. It estimates the coefficients of a landmark-gaze model as well as the scale and translation parameters. The three results are used to calculate eye landmarks and estimated gaze.

truth. Many specific tasks are designed for unsupervised CNNs. Dubey *et al.* [83] collect unlabeled facial images from webpages. They roughly annotate the gaze region based on the detected landmarks. Therefore, they can perform the classical supervised task for gaze representation learning. Yu *et al.* utilize a pre-trained gaze redirection network to perform unsupervised gaze representation learning [62]. As shown in Fig. 8, they use the gaze representation difference of the input and target images as the redirection variables. Given the input image and the gaze representation difference, the gaze network aims to reconstruct the target image. Therefore, the reconstruction task supervises the optimization of the gaze representation network. Sun *et al.* propose a cross-encoder for unsupervised learning [118]. They acquire paired eye images for training where the paired images have the same gaze or appearance. They use an encoder to extract appearance and gaze features from eye images. They exchange the two features of selected paired images and aim to reconstruct the original image based on the exchanged feature. Note that, these approaches learn the gaze representation, but they also require a few labeled samples to fine-tune the final gaze estimator.

3.2.3 Multi-task CNNs

Multi-task learning usually contains multiple tasks that provide related domain information as inductive bias to improve model generalization [120]. Some auxiliary tasks are proposed for improving model generalization in gaze estimation.

Lian *et al.* propose a multi-task multi-view network for gaze estimation [121]. They estimate gaze directions based on single-view eye images and PoG from multi-view eye images. They also propose another multi-task CNN to estimate PoG using depth images [122]. They design an additional task to leverage facial features to refine depth images. The network produces four features for gaze estimation, which are extracted from the facial images, the left/right eye images and the depth images.

Some works seek to decompose the gaze into multiple related features and construct multi-task CNNs to estimate these features. Yu *et al.* introduce a constrained landmark-gaze model for modeling the joint variation of eye landmark locations and gaze directions [119]. As shown in Fig. 9, they build a multi-task CNN to estimate the coefficients of the landmark-gaze model as well as the scale and translation information to align eye landmarks. Finally, the landmark-gaze model serves as a decoder to calculate gaze from estimated parameters. Deng *et al.* decompose the gaze direction into eyeball movement and head pose [76]. They design a multi-tasks CNN to estimate the eyeball movement from eye images and the head pose from facial images. The gaze direction is computed from eyeball movement and head pose using geometric transformation. Wu *et al.* propose a multi-task CNN that simultaneously segments the eye part, detects the IR LED glints,

and estimates the pupil and cornea center [123]. The gaze direction is covered from the reconstructed eye model.

Other works perform multiple gaze-related tasks simultaneously. Recasens *et al.* present an approach for following gaze in video by predicting where a person (in the video) is looking, even when the object is in a different frame [124]. They build a CNN to predict the gaze location in each frame and the probability containing the gazed object of each frame. Also, visual saliency shows strong correlation with human gaze in scene images [125, 126]. In [127], they estimate the general visual attention and human's gaze directions in images at the same time. Kellnhofer *et al.* propose a temporal 3D gaze network [43]. They use bi-LSTM [128] to process a sequence of 7 frames to estimate not only gaze directionS but also gaze uncertainty.

3.2.4 Recurrent CNNs

Human eye gaze is continuous. This inspires researchers to improve gaze estimation performance by using temporal information. Recently, recurrent neural networks have shown great capability in handling sequential data. Some researchers employ recurrent CNNs to estimate the gaze in videos [43, 79, 84].

We first give a typical example of the data processing workflow. Given a sequence of frames $\{X_1, X_2, \dots, X_N\}$, a united CNN f_U is used to extract feature vectors from each frame, *i.e.*, $x_t = f_U(X_t)$. These feature vectors are fed into a recurrent neural network f_R and the network outputs the gaze vector, *i.e.*, $g_i = f_R(x_1, x_2, \dots, x_N)$. Palmero *et al.* set $N = 4$ and $i = 4$ in their method. They input four frames to estimate the gaze of the last frame [79]. Kellnhofer *et al.* set $N = 7$ and $i = 4$ [43]. They consider extra three frames after the target frame compared with Palmero. These methods both select the nearest three frames (including the previous and the next three frames) for additional vision feature. Besides, some methods utilize the past gaze trajectory for gaze prediction [86, 129]. They select a larger time range, *e.g.*, $1 \sim 2s$ (30 \sim 90 frames), in the gaze prediction task. We visualize a expample network architecture in Fig. 5.

Different types of input have been explored to extract features. Kellnhofer *et al.* directly extract features from facial images [43]. Zhou *et al.* combine the feature extracted from facial and eye images [84]. Palmero *et al.* use facial images, binocular images and facial landmarks to generate the feature vectors [79]. Different RNN structures have also been explored, such as GRU [77] in [79], LSTM [130] in [84] and bi-LSTM [128] in [43]. Cheng *et al.* leverage the recurrent CNN to improve gaze estimation performance from static images rather than videos [56]. They generalize the gaze estimation as a sequential coarse-to-fine process and use GRU to relate the basic gaze direction estimated from facial images and the gaze residual estimated from eye images.

3.2.5 CNNs with Other Priors

Prior information also helps to improve gaze estimation accuracy, such as decomposition of gaze direction, anatomical eye models and eye movement patterns [45, 55, 76, 86, 114, 131].

Decomposition of Gaze Direction. Human gaze can be decomposed into the head pose and the eyeball pose. Deng *et al.* use two CNNs to estimate head pose from facial images and eyeball pose from eye images. They integrate these two results into final gaze directions using geometric transformation [76].

Anatomical Eye Model. The human eye is composed of eyeball, iris, pupil center and etc. Park *et al.* propose a pictorial gaze representation based on the eye model [45]. They render the eye

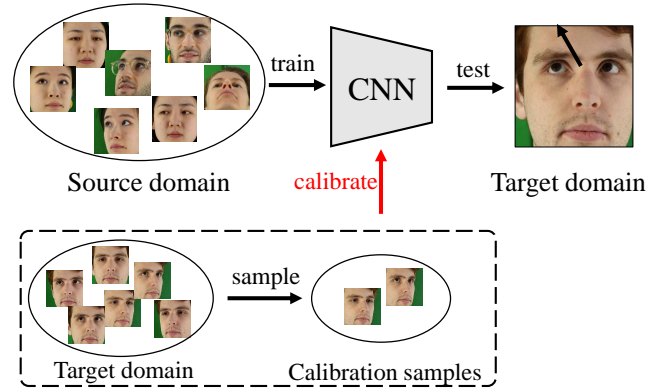


Fig. 10. Personal calibration in deep learning. The method usually samples a few images from the target domain as calibration samples. The calibration samples and training set are jointly used to improve the performance in target domain.

model to generate a pictorial image, where the pictorial image eliminates the appearance variance. They first generate pictorial images from original images using CNN and use another CNN to estimate gaze directions from pictorial image.

Eye Movement Pattern. Common eye movements, such as fixation, saccade and smooth pursuits, are independent of viewing contents and subjects. Wang *et al.* propose to incorporate the generic eye movement pattern in dynamic gaze estimation [86]. They recover the eye movement pattern from videos and use a CNN to estimate gaze from static images.

Two eye asymmetry Property. Cheng *et al.* discover the 'two eye asymmetry' property that the appearances of two eyes are different while the gaze directions of two eyes are approximately the same [44]. Based on this observation, Cheng *et al.* propose to treat two eyes asymmetrically in the CNN. They design an asymmetry regression network for adaptively weighting two eyes.

Gaze data distribution. The basic assumption of regression models is independent identically distributed, however, gaze data is not *i.i.d.* Xiong *et al.* discuss the problem [114] and design a mixed-effect model to consider person-specific information.

Inter-subject bias. Chen *et al.* observe the inter-subject bias in most datasets [131, 132]. They make the assumption that there exists a subject-dependent bias that cannot be estimated from images. Thus, they propose a gaze decomposition method. They decompose the gaze into the subject-dependent bias and the subject-independent gaze estimated from images. During test, they use some image samples to calibrate the subject-dependent bias.

3.3 Calibration

It is non-trivial to learn an accurate and universal gaze estimation model. Conventional 3D eye model recovery methods usually build a unified gaze model including subject-specific parameters such as eyeball radius [28]. They perform a personal calibration to estimate these subject-specific parameters. In the field of deep learning-based gaze estimation, personal calibration is also explored to improve person-specific performance. Fig. 10 shows a common pipeline of personal calibration in deep learning.

3.3.1 Calibration via Domain Adaption

The calibration problem can be considered as domain adaption problems, where the training set is the source domain and the test set is the target domain. The test set usually contains unseen

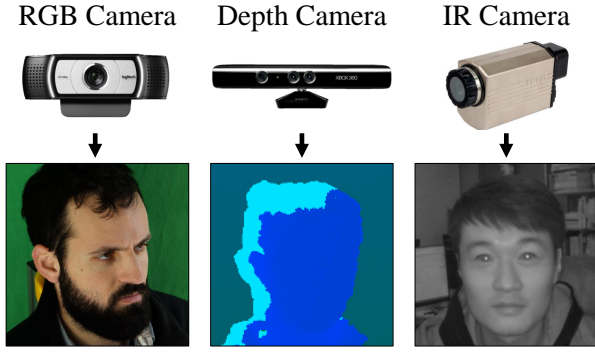



Fig. 11. Different cameras and their captured images.

subjects or unseen environment. Researchers aim to improve the performance in the target domain using calibration samples.

The common approach of domain adaption is to fine-tune the model in the target domain [42, 133, 134]. This is simple but effective. Krafka *et al.* replace the fully-connected layer with an SVM and fine-tune the SVM layer to predict the gaze location [42]. Zhang *et al.* split the CNN into three parts: the encoder, the feature extractor, and the decoder [133]. They fine-tune the encoder and decoder in each target domain. Zhang *et al.* also learn a third-order polynomial mapping function between the estimated and ground-truth of 2D gaze locations [135]. Some studies introduce person-specific feature for gaze estimation [117, 136]. They learn the person-specific feature during fine-tuning. Linden *et al.* introduce user embedding for recording personal information. They obtain user embedding of the unseen subjects by fine-tuning using calibration samples [136]. Chen *et al.* [131, 132] observe the different gaze distributions of subjects. They use the calibration samples to estimate the bias between the estimated gaze and the ground-truth of different subjects. They use bias to refine the estimates. Yu *et al.* generate additional calibration samples through the synthesis of gaze-redirectioned eye images from calibration samples [46]. The generated samples are also directly used for training. These methods all need labeled samples for supervised calibration.

Unsupervised calibration methods attract much attention recently. These methods use unlabeled calibration samples to improve performance. Wang *et al.* propose an adversarial method for aligning features. They build a discriminator to judge the source of images from the extracted feature. The feature extractor has to confuse the discriminator, *i.e.*, the generated feature should be domain-invariant. Guo *et al.* [137] use source samples to form a locally linear representation of each target domain prediction in gaze space. The same linear relationships are applied in the feature space to generate the feature representation of target samples. Meanwhile, they minimize the difference between the generated feature and extracted feature of target sample for alignment. Cheng *et al.* [72] propose a domain generalization method. They improve the cross-dataset performance without knowing the target dataset or touching any new samples. They propose a self-adversarial framework to remove the gaze-irrelevant features in face images. Cui *et al.* define a new adaption problem [138]: adaptation from adults to children. They use the conventional domain adaption method, geodesic flow kernel [139], to transfer the feature in the adult domain into the children domain. Bao *et al.* [140] estimate the point-of-regard by aligning the predicted gaze distribution with known gaze distribution.

Some well-known strategies in universal tasks are proved effective for gaze estimation. Meta learning and metric learn-



| Platforms | Computer | HMD Devices | Mobile Devices |
|-------------------------|------------|-------------|----------------|
| Gaze Zone | Medium | Large | Small |
| User-camera distance | ~60 cm | Near-eye | <30 cm |
| Camera types | RGB | IR | RGB |
| Computational Resources | Sufficient | Limited | Limited |

Fig. 12. Different platforms and their characteristics.

ing show great potentials in personalized gaze estimation. They usually require few-shot annotated samples for calibration. Park *et al.* propose a meta learning-based calibration approach [47]. They train a highly adaptable gaze estimation network through meta learning. The network can be converted into a person-specific network once training with target person samples. Liu *et al.* propose a differential CNN based on metric learning [141]. The network predicts gaze difference between two eye images. During test stage, it estimates the differences between inputs and calibration images, and takes the average results as the estimation.

Contrastive learning and mean teacher [142] perform well in unsupervised domain adaption. They are usually used for cross-dataset task in gaze estimation. Liu *et al.* propose an outlier-guided collaborative learning for unsupervised cross-dataset tasks [112]. They create a group of teacher-student networks where teacher networks are pre-trained in source domain. They design the outlier-guided loss which requires the outputs of teacher and student networks to be consistent. Bao *et al.* also propose a mean teacher architecture for unsupervised cross-dataset task [111]. They perform data augmentation w.r.t. rotation in target domains and require the rotation consistency in gaze estimation. Wang *et al.* [143] propose a contrastive learning for cross-dataset gaze estimation. They propose a contrastive loss function to encourage close feature distance for the samples with close gaze directions.

3.3.2 Calibration via User-unaware Data Collection

It is difficult to acquire enough samples for calibration in practical applications. Collecting calibration samples in a user-unaware manner is an alternative solution [144–146]. Salvalaio *et al.* implicitly collect calibration data when users are using computers. They collect data when the user is clicking a mouse, this is based on the assumption that users are gazing at the position of the cursor when clicking the mouse [146]. They use online learning to fine-tune their model with the calibration samples. Some studies investigate the relation between the gaze points and the saliency maps [125, 126]. Chang *et al.* utilize saliency information to adapt the gaze estimation algorithm to a new user without explicit calibration [144]. They transform the saliency map into a differentiable loss map that can be used to optimize the CNN models. Wang *et al.* introduce a stochastic calibration procedure. They minimize the difference between the probability distribution of predicted gaze and ground truth [145].

3.4 Devices and Platforms

3.4.1 Camera

The majority of gaze estimation systems use a single RGB camera to capture eye images, while some studies use different

TABLE 1
Summary of gaze estimation methods.

| Perspectives | | Methods | | | | | | | |
|--------------|-------------------------------|---------|-------|-----------------------------|--|---|--|-----------------------------------|---|
| | | 2015 | 2016 | 2017 | 2018 | 2019 | 2020 | 2021 | 2022 |
| Feature | Eye image | [17] | — | [49, 138, 147] | [45, 54, 55, 119, 148] | [47, 53, 59, 121, 123, 141, 149] | [56, 57, 60–62, 90, 91] | [58, 78, 118] | [132, 140, 150, 151] |
| | Facial image | — | [42] | [50, 76, 124] | [79, 82, 127, 133] | [53, 83, 114, 117, 122, 135, 136, 144, 146, 152] | [44, 56, 57, 66, 68, 71, 73–75, 80, 89, 93, 131, 137, 153, 154] | [14, 78, 110, 112, 155] | [34, 69, 72, 111, 132, 143, 156–163] |
| | Video | — | — | — | [79] | [43, 84–86] | [67] | [110] | — |
| Model | Supervised CNN | [17] | [42] | [49, 50, 76, 124, 138, 147] | [45, 54, 79, 82, 119, 127, 133, 148] | [43, 47, 53, 84–86, 121–123, 135, 141, 144, 146, 149, 152] | [44, 56, 57, 60, 61, 66–68, 71, 73–75, 89–91, 93, 131, 153, 154] | [14, 58, 78, 155] | [34, 69, 72, 132, 150, 156, 158, 159, 161] |
| | Semi-/Self-/Un-Supervised CNN | — | — | — | [55] | [59, 83, 114, 117, 136] | [44, 62, 137] | [110, 112, 118, 140] | [69, 72, 111, 140, 143, 151, 157, 160, 162, 163] |
| | Multi-task CNN | — | — | [76, 124] | [119, 127] | [43, 121–123] | [67] | — | [151] |
| | Recurrent CNN | — | — | — | [79] | [43, 84] | [56, 67] | — | — |
| | CNN with Priors | — | — | [76] | [45, 55] | [86, 114] | [131] | — | [132] |
| Calibration | Domain Adaption | — | [42] | [138] | [133] | [46, 47, 117, 135, 136, 141] | [131, 134, 137] | [112, 140] | [111, 132, 140, 143, 151, 157, 160] |
| | User-unaware Data Collection | — | [145] | — | — | [144, 146] | — | — | — |
| Camera | Single camera | [17] | [42] | [49, 50, 76, 138] | [45, 54, 55, 79, 82, 119, 127, 133, 148] | [43, 47, 53, 83–86, 117, 123, 135, 141, 144, 146, 149, 152] | [44, 56, 57, 60, 61, 66–68, 71, 73–75, 89–91, 93, 131, 153, 154] | [58, 78, 110, 112, 118, 140, 155] | [34, 69, 72, 111, 132, 140, 143, 150, 151, 156–162] |
| | Multi cameras | — | — | [147] | — | [121] | — | — | [163] |
| | IR Camera | — | — | — | — | [123, 149] | [61] | — | — |
| | RGBD Camera | — | — | — | — | [122] | — | — | — |
| | Near-eye Camera | — | — | [147] | — | [123, 149] | — | — | — |
| Platform | Computer | [17] | — | [49, 50, 76, 138] | [45, 54, 55, 79, 82, 119, 127, 133] | [43, 47, 53, 83–86, 121, 122, 135, 141, 144, 146] | [44, 56, 60, 61, 66–68, 71, 73–75, 89–91, 93, 131, 153] | [58, 78, 110, 112, 118, 140, 155] | [34, 69, 72, 111, 132, 140, 143, 150, 151, 156, 157, 159–163] |
| | Mobile Device | — | [42] | — | [133] | [117, 152] | [57, 154] | — | — |
| | HMD Device | — | — | [147] | [148] | [123, 149] | — | — | — |

camera settings, *e.g.*, using multiple cameras to capture multi-view images [121, 147, 164], using infrared (IR) cameras to handle low illumination condition [123, 149], and using RGBD cameras to provide the depth information [122]. Different cameras and their captured images are shown in Fig. 11.

Tonsen *et al.* embed multiple millimeter-sized RGB cameras into a normal glasses frame [147]. They use multi-layer perceptrons to process the eye images captured by different cameras, and concatenate the extracted feature to estimate gaze. Lian *et al.* mount three cameras at the bottom of a screen [121]. They build a multi-branch network to extract the features of each view and concatenate them to estimate 2D gaze position on the screen. Wu *et al.* collect gaze data using near-eye IR cameras [123]. They use CNN to detect the location of glints, pupil centers and corneas from IR images. Then, they build an eye model using the detected feature and estimate gaze from the gaze model. Kim *et al.* collect a large-scale dataset of near-eye IR eye images [149]. They synthesize additional IR eye images that cover large variations in face shape, gaze direction, pupil and iris etc.. Lian *et al.* use RGBD cameras to capture depth facial images [122]. They extract the depth information of eye regions and concatenate it with RGB image features to estimate gaze.

3.4.2 Platform

Eye gaze can be used to estimate human intent in various applications, *e.g.*, product design evaluation [165], marketing studies [166] and human-computer interaction [10, 167, 168]. These applications can be simply categorized into three types of platforms: computers, mobile devices and head-mounted devices. We summarize the characteristics of these platforms in Fig. 12.

The computer is the most typical platform for appearance-based gaze estimation. The cameras are usually placed below/above the computer screen [17, 45, 55, 56, 169]. Some works focus on using deeper neural networks [17, 50, 54] or extra modules [45, 55, 56] to improve gaze estimation performance, while the other studies seek to use custom devices for gaze estimation, such as multi-cameras and RGBD cameras [121, 122].

The mobile device contains front cameras but has limited computational resources. The related methods usually estimate PoG instead of gaze directions due to the difficulty of geometric calibration. Krafka *et al.* propose iTracker for mobile devices [42], which combines the facial image, two eye images and the face grid to estimate the gaze. The face grid encodes the position of the face in captured images and is proved to be effective for gaze estimation in mobile devices in many works [57, 117].

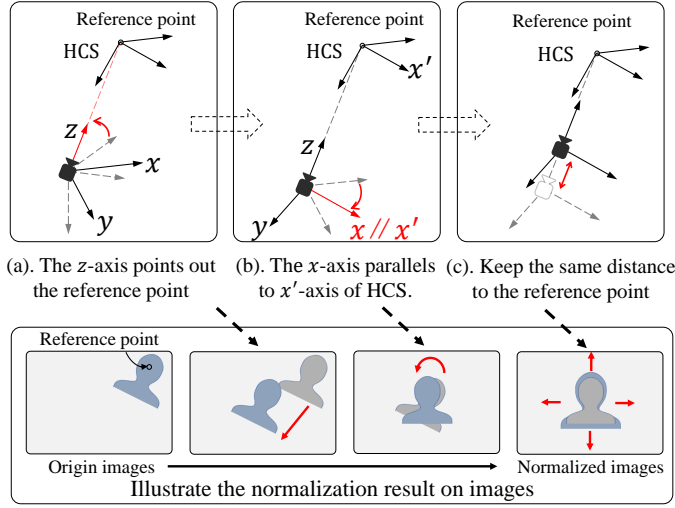


Fig. 13. A data rectification method [37]. The virtual camera is rotated so that the z -axis points at the reference point and the x -axis is parallel with the x' -axis of the head coordinate system (HCS). The bottom row illustrates the rectification result on images. Overall, the reference point is moved to the center of images, the image is rotated to straighten face and scaled to align the size of face in different images.

He *et al.* propose a more accurate and faster method based on iTracker [117]. They replace the face grid with eye corner landmark feature. Guo *et al.* propose a generalized gaze estimation method [152]. They observe the notable jittering problem in gaze point estimates and propose to use adversarial training to address this problem. Valliappan [170] evaluate the eye tracking with deep learning on smartphone. They show the algorithm can achieve competitive result compared with modern eye tracking devices.

The head-mounted device usually employs near-eye cameras to capture eye images. Tonsen *et al.* embed millimetre-sized RGB cameras into a normal glasses frame [147]. In order to compensate for the low-resolution captured images, they use multi-cameras to capture multi-view images and use a neural network to regress gaze from these images. IR cameras are also employed by head-mounted devices. Wu *et al.* collect the MagicEyes dataset using IR cameras [123]. They propose EyeNet, a neural network that solves multiple heterogeneous tasks related to eye gaze estimation for an off-axis camera setting. They use the CNN to model 3D cornea and 3D pupil and estimate the gaze from these two 3D models. Lemley *et al.* use the single near-eye image as input to the neural network and directly regress gaze [148]. Kim *et al.* follow a similar approach and collect the NVGaze dataset [149].

3.5 Summary

Tab. 1 summarizes the existing CNN-based gaze estimation methods. Note that many methods do not specify a platform [17, 56]. Thus, we categorize these methods into the platform of "computer". In general, there is an increasing trend in developing supervised or semi-/self-/un-supervised CNN structures to estimate gaze. Many recent research interests shift to different calibration approaches through domain adaptation or user-unaware data collection. The first CNN-based gaze direction estimation method is proposed by Zhang *et al.* in 2015 [17], the first CNN-based PoG estimation method is proposed by Krafka *et al.* in 2016 [42]. These two studies both provide large-scale gaze datasets, the MPIIGaze and the GazeCapture, which have been widely used for evaluating gaze estimation algorithms in later studies.

TABLE 2
Summary of face alignment methods

| Names | Years | Pub. | Links |
|----------------|-------|-------|---|
| Dlib [171] | 2014 | CVPR | https://pypi.org/project/dlib/19.6.0/ |
| MTCNN [172] | 2016 | SPL | https://github.com/kpzhang93/MTCNN_face_detection_alignment |
| DAN [173] | 2017 | CVPRW | https://github.com/MarekKowalski/DeepAlignmentNetwork |
| OpenFace [174] | 2018 | FG | https://github.com/TadasBaltrusaitis/OpenFace |
| PRN [175] | 2018 | ECCV | https://github.com/YadiraF/PRNet |
| 3DDFA_V2 [176] | 2020 | ECCV | https://github.com/cleardusk/3DDFA_V2 |

4 DATASETS AND BENCHMARKS

4.1 Data Pre-processing

4.1.1 Face and Eye Detection

Raw images often contain unnecessary information for gaze estimation, such as the background. Directly using raw images to regress gaze not only increases the computational resource but also brings nuisance factors such as changes in scenes. Therefore, face or eye detection is usually applied in raw images to prune unnecessary information. Generally, researchers first perform face alignment in raw images to obtain facial landmarks and crop face/eye images using these landmarks. Several face alignment methods have been proposed recently [177–179]. We list some typical face alignment methods in Tab. 2.

After the facial landmarks are obtained, face or eye image are cropped accordingly. There is no protocol to regulate the cropping procedure. We provide a common cropping procedure here as an example. We let $x_i \in \mathbb{R}^2$ be the x, y -coordinates of the i th facial landmark in an raw image I . The center point \bar{x} is calculated as $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$, where n is the number of facial landmarks. The face image is defined as a square region with the center \bar{x} and an width w . The w is usually set empirically. For example, [50] set w as 1.5 times of the maximum distance between the landmarks. The eye cropping is similar to face cropping, while the eye region is usually defined as a rectangle with the center set as the centroid of eye landmarks. The width of the rectangle is set based the distance between eye corners, *e.g.*, 1.2 times.

4.1.2 Data Rectification

Data rectification eliminate environment factors such as head pose and illumination. It simplifies gaze regression problem with data pre-processing methods. Sugano *et al.* propose to rectify the eye image by rotating the virtual camera to point at the same reference point in the human face [37]. They assume that the captured eye image is a plane in 3D space, the rotation of the virtual camera can be performed as a perspective transformation on the image. The whole data rectification process is shown in Figure 13. They compute the transformation matrix $M = SR$, where R is the rotation matrix and S is the scale matrix. R also indicates the rotated camera coordinate system. The z -axis z_c of the rotated camera coordinate system is defined as the line from cameras to reference points, where the reference point is usually set as the face center or eye center. It means that the rotated camera is pointing towards the reference point. The rotated x -axis x_c is defined as the x -axis of the head coordinate system so that the appearance captured by the rotated cameras is facing the front. The rotated y -axis y_c can be computed by $y_c = z_c \times x_c$, the x_c is recalculated by $x_c = y_c \times z_c$ to maintain orthogonality. As a result, the rotation matrix $R = [\frac{x_c}{\|x_c\|}, \frac{y_c}{\|y_c\|}, \frac{z_c}{\|z_c\|}]$. The S

TABLE 3
Symbol table in data post-processing

| Symbol | Meaning |
|-----------------------------------|---|
| $p \in \mathbb{R}^2$ | $p = (u, v)$, gaze targets. |
| $g \in \mathbb{R}^3$ | $g = (g_x, g_y, g_z)$, gaze directions. |
| $o \in \mathbb{R}^3$ | $o = (x_o, y_o, z_o)$, origins of gaze directions. |
| $t \in \mathbb{R}^3$ | $t = (x_t, y_t, z_t)$, targets of gaze directions. |
| $R_s \in \mathbb{R}^{3 \times 3}$ | The rotation matrix of SCS w.r.t. CCS. |
| $T_s \in \mathbb{R}^3$ | $T_s = (t_x, t_y, t_z)$, the translation matrix of SCS w.r.t. CCS. |
| $n \in \mathbb{R}^3$ | $n = (n_x, n_y, n_z)$, the normal vectors of x-y plane of SCS. |

maintains the distance between the virtual camera and the reference point, which is defined as $diag(1, 1, \frac{d_n}{d_o})$, where d_o is the original distance between the camera and the reference point, and d_n is the new distance that can be adjusted manually. They apply a perspective transformation on images with $W = C_n M C_r^{-1}$, where C_r is the intrinsic matrix of the original camera and C_n is the intrinsic matrix of the new camera. Gaze directions can also be calculated in the rotated camera coordinate system as $\hat{g} = M g$. The method eliminates the ambiguity caused by different head positions and aligns the intrinsic matrix of cameras. It also rotates the captured image to cancel the degree of freedom of roll in head rotation. Zhang *et al.* further explore the method in [180]. They argue that scaling can not change the gaze direction vector. The gaze direction is computed by $\hat{g} = R g$.

Illumination also influences the appearance of the human eye. To handle this, researchers usually take gray-scale images rather than RGB images as input and apply histogram equalization in the gray-scale images to enhance the image.

4.2 Data Post-processing

Various applications require different forms of gaze estimates. For example, in a real-world interaction task, it requires 3D gaze direction to estimate the human intent [7, 181], while it requires 2D PoG for the screen-based interaction [10, 182]. In this section, we introduce how to convert different forms of gaze estimates by post-processing. We list the symbols in Tab. 3 and illustrate the symbols in Fig. 14. We denote the PoG as 2D gaze and the gaze direction as 3D gaze in this section.

4.2.1 2D/3D Gaze Conversion

The 2D gaze estimation algorithm usually estimates gaze targets on a computer screen [42, 86, 144, 152, 183], while the 3D gaze estimation algorithm estimates gaze directions in 3D space [43, 49, 50, 56, 114]. We first introduce how to convert between the 2D gaze and the 3D gaze.

Given a 2D gaze target $p = (u, v)$ on the screen, our goal is to compute the corresponding 3D gaze direction $g = (g_x, g_y, g_z)$. The processing pipeline is that we first compute the 3D gaze target t and 3D gaze origin o in the camera coordinate system (CCS). The gaze direction can be computed as

$$g = \frac{t - o}{\|t - o\|}. \quad (1)$$

To derive the 3D gaze target t , we obtain the pose $\{R_s, T_s\}$ of screen coordinate system (SCS) w.r.t. CCS by geometric calibration, where R_s is the rotation matrix and T_s is the translation matrix. The t is computed as $t = R_s[u, v, 0]^T + T_s$, where the additional 0 is the z -axis coordinate of p in SCS. The 3D gaze

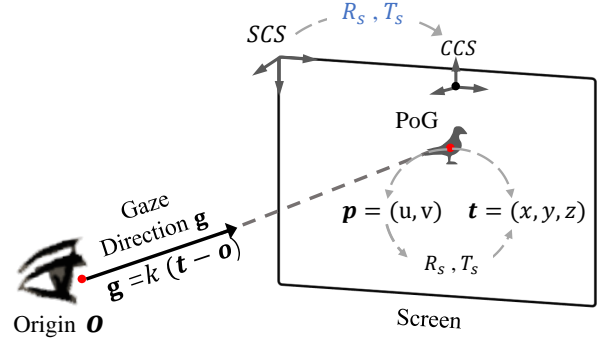


Fig. 14. We illustrate the relation between gaze directions and PoG. Gaze directions are originated from an origin o and intersect with the screen at the PoG. The PoG is usually denoted as a 2D coordinate p . It can be converted to 3D coordinate t in CCS with screen pose $\{R_s, T_s\}$. Gaze directions can be computed with $k(t - o)$, where k is a scale factor.

origin o is usually defined as the face center or the eye center. It can be estimated by landmark detection algorithms or stereo measurement methods.

On the other hand, given a 3D gaze direction g , we aim to compute the corresponding 2D target point p on the screen. Note that, we also need to acquire the screen pose $\{R_s, T_s\}$ as well as the origin point o as mentioned previously. We first compute the intersection of gaze direction and screen, i.e., 3D gaze target t , in CCS, and then we convert the 3D gaze target to the 2D gaze target using the pose $\{R_s, T_s\}$.

To deduce the equation of screen plane, we compute $n = R_s[:, 2] = (n_x, n_y, n_z)$, where n is the normal vector of screen plane. $T_s = [t_x, t_y, t_z]^T$ also represents a point on the screen plane. Therefore, the equation of the screen plane is

$$n_x x + n_y y + n_z z = n_x t_x + n_y t_y + n_z t_z. \quad (2)$$

Given a gaze direction g and the origin point o , we can write the equation of the line of sight as

$$\frac{x - x_o}{g_x} = \frac{y - y_o}{g_y} = \frac{z - z_o}{g_z}. \quad (3)$$

By solving Eq. (2) and Eq. (3), we obtain the intersection t , and $(u, v, z) = R_s^{-1}(t - T_s)$, where z usually equals to 0 and $p = (u, v)$ is the coordinate of 2D target point in metre.

4.2.2 Gaze Origin Conversion

Conventional gaze estimation methods usually estimate gaze directions w.r.t. each eye. They define the origin of gaze directions as each eye center [45, 49, 59, 141]. Recently, more attention has been paid to gaze estimation using the face images and they estimate gaze direction w.r.t. the whole face. They define the gaze vector starting from the face center to the gaze target [47, 50, 54, 56]. Here we introduce a gaze origin conversion method to bridge the gap between these two types of gaze estimates.

We first compute the pose $\{R_s, T_s\}$ of SCS and the origin o of the predicted gaze direction g through calibration. Then we can write Eq. (2) and Eq. (3) based on these parameters. The 3D gaze target point t can be calculated by solving the equation of Eq. (2) and Eq. (3). Next, we obtain the new origin o_n of the gaze direction through 3D landmark detection. The new gaze direction can be computed by

$$g_{new} = \frac{t - o_n}{\|t - o_n\|}. \quad (4)$$

TABLE 4
Summary of common gaze estimation datasets.

| Datasets | Subjects | Total | Annotations | | | Brief Introduction | Links |
|--|----------|---------------|-------------|---------|---------|---|---|
| | | | Full face | 2D Gaze | 3D Gaze | | |
| Columbia [184], 2013, (Columbia University) | 58 | 6K images | ✓ | × | ✓ | Collected in laboratory; 5 head pose and 21 gaze directions per head pose. | https://cs.columbia.edu/CAVE/databases/columbia_gaze |
| UTMultiview [37], 2014, (The University of Tokyo; Microsoft Research Asia) | 50 | 1.1M images | × | ✓ | ✓ | Collected in laboratory; Fixed head pose; Multiview eye images; Synthesis eye images. | https://ut-vision.org/datasets |
| EyeDiap [185], 2014, (Idiap Research Institute) | 16 | 94 videos | ✓ | ✓ | ✓ | Collected in laboratory; Free head poses; Additional depth videos. | https://idiap.ch/dataset/eyediap |
| MPIIGaze [49], 2015, (Max Planck Institute) | 15 | 213K images | × | ✓ | ✓ | Collected by laptops in daily life; Free head pose and illumination. | https://mpi-inf.mpg.de/mpiigaze |
| GazeCapture [42], 2016, (University of Georgia; MIT; Max Planck Institute) | 1,474 | 2.4M images | ✓ | ✓ | × | Collected by mobile devices in daily life; Variable lighting condition and head motion. | https://gazecapture.csail.mit.edu |
| MPIIFaceGaze [50], 2017, (Max Planck Institute) | 15 | ~ 45K images | ✓ | ✓ | ✓ | Collected by laptops in daily life; Free head pose and illumination. | footnote ¹ |
| InvisibleEye [147], 2017, (Max Planck Institute; Osaka University) | 17 | 280K Images | × | ✓ | × | Collected in laboratory; Multiple near-eye camera; Low resolution cameras. | https://mpi-inf.mpg.de/invisibleeye |
| TabletGaze [186], 2017, (Rice University) | 51 | 816 videos | ✓ | ✓ | × | Collected by tablets in laboratory; Four postures to hold the tablets; Free head pose. | https://sh.rice.edu/cognitive-engagement/tabletgaze |
| RT-Gene [54], 2018, (Imperial College London) | 15 | 123K images | ✓ | × | ✓ | Collected in laboratory; Free head pose; Annotated with mobile eye-tracker; Use GAN to remove the eye-tracker in face images. | https://github.com/Tobias-Fischer/rt_gene |
| Gaze360 [43], 2019, (MIT; Toyota Research Institute) | 238 | 172K images | ✓ | × | ✓ | Collected in indoor and outdoor environments; A wide range of head poses and distances between subjects and cameras. | https://gaze360.csail.mit.edu |
| NVGaze [149], 2019, (NVIDIA; UNC) | 30 | 4.5M images | × | ✓ | × | Collected in laboratory; Near-eye Images; Infrared illumination. | https://sites.google.com/nvidia.com/nvgaze |
| ShanghaiTechGaze [121], 2019, (ShanghaiTech University; UESTC) | 137 | 224K images | ✓ | ✓ | × | Collected in laboratory; Free head poses; Multiview gaze dataset. | https://github.com/dongzeliang/multi-view-gaze |
| ETH-XGaze [66], 2020, (ETH Zurich; Google) | 110 | 1.1M images | ✓ | ✓ | ✓ | Collected in laboratory; High-resolution images; Extreme head pose; 16 illumination conditions. | https://ait.ethz.ch/projects/2020/ETH-XGaze |
| EVE [67], 2020, (ETH Zurich) | 54 | ~ 4.2K videos | ✓ | ✓ | ✓ | Collected in laboratory; Free head pose; Free view; Annotated with desktop eye tracker; Pupil size annotation. | https://ait.ethz.ch/projects/2020/EVE/ |

¹ <https://mpi-inf.mpg.de/departments/computer-vision-and-machine-learning/research/gaze-based-human-computer-interaction/its-written-all-over-your-face-full-face-appearance-based-gaze-estimation>

TABLE 5

Benchmark of **within-dataset evaluation**. We use the provided source codes or re-implement ([†]) the methods for comparison. The underlines indicate the top three best performances. Note that the methods in the last row are proposed for **point of gaze estimation**, we convert the result using the post-processing method in Sec. 4.2.

| Methods | Pub. | MPIIGaze [49] | EyeDiap [185] | UT [37] | MPIIFaceGaze [50] | EyeDiap [185] | Gaze360 [43] | RT-Genie [54] | ETH-XGaze [66] |
|-------------------------------|---------|------------------|------------------|--------------|----------------------|------------------|-----------------|------------------|-------------------|
| Mnist [†] [17] | CVPR15 | 6.27° | 7.60° | <u>6.34°</u> | 6.39° | 7.37° | N/A | N/A | N/A |
| GazeNet [†] [49] | TPAMI17 | 5.70° | 7.13° | <u>6.44°</u> | 5.76° | 6.79° | N/A | N/A | N/A |
| Dilated-Net [†] [53] | ACCV19 | 4.39° | 6.57° | N/A | 4.42° | 6.19° | 13.73° | 8.38° | N/A |
| Gaze360 [43] | ICCV19 | <u>4.07°</u> | <u>5.58°</u> | N/A | <u>4.06°</u> | 5.36° | <u>11.04°</u> | <u>7.06°</u> | <u>4.46°</u> |
| RT-Genie [54] | ECCV18 | 4.61° | <u>6.30°</u> | N/A | 4.66° | 6.02° | 12.26° | 8.60° | N/A |
| FullFace [50] | CVPRW17 | 4.96° | 6.76° | N/A | 4.93° | 6.53° | 14.99° | 10.00° | <u>7.38°</u> |
| RCNN [†] [79] | BMVC18 | N/A | N/A | N/A | 4.10° | <u>5.31°</u> | 11.23° | 10.30° | N/A |
| CA-Net [56] | AAAI20 | <u>4.27°</u> | <u>5.63°</u> | N/A | 4.27° | <u>5.27°</u> | <u>11.20°</u> | <u>8.27°</u> | N/A |
| GazeTR-Pure [34] | ICPR22 | N/A | N/A | N/A | 4.74° | 5.72° | 13.58° | 8.06° | N/A |
| GazeTR-Hybrid [34] | ICPR22 | N/A | N/A | N/A | 4.00° | 5.17° | 10.62° | 6.55° | N/A |
| Itracker [†] [42] | CVPR16 | 7.25° | 7.50° | N/A | 7.33° | 7.13° | N/A | N/A | N/A |
| AFF-Net [57] | ICPR20 | <u>3.69°</u> | 6.75° | N/A | <u>3.73°</u> | 6.41° | N/A | N/A | N/A |

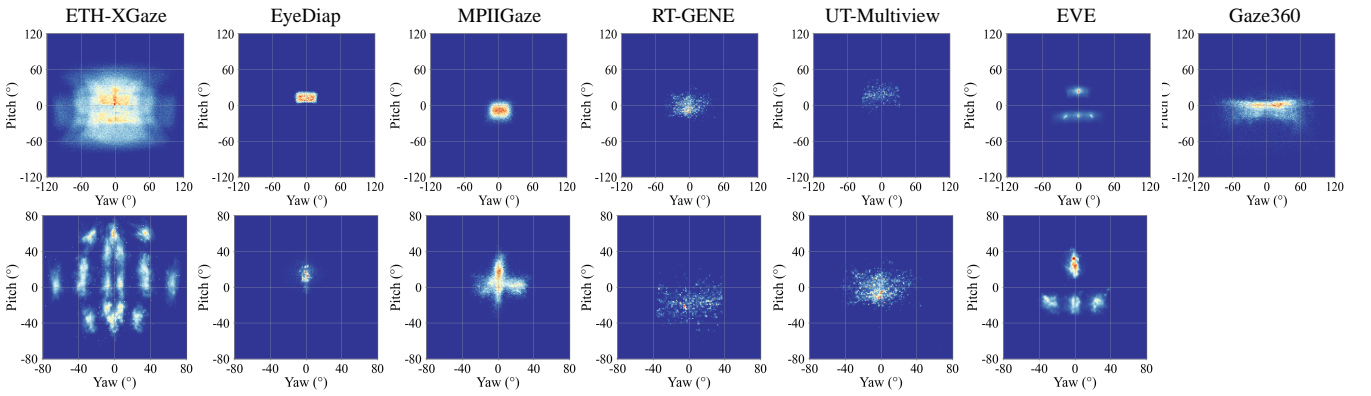


Fig. 15. Distribution of head pose and gaze in different datasets. The first row is the distribution of gaze and the second row show head distribution.

4.3 Evaluation Metrics

Two types of metric are used for performance evaluation: **angular error** and the **Euclidean distance**. The angular error measures the accuracy of 3D gaze estimation methods [47, 49, 56]. Assuming the actual gaze direction is $\mathbf{g} \in \mathbb{R}^3$ and the estimated gaze direction is $\hat{\mathbf{g}} \in \mathbb{R}^3$, the **angular error** can be computed as:

$$\mathcal{L}_{\text{angular}} = \frac{\mathbf{g} \cdot \hat{\mathbf{g}}}{\|\mathbf{g}\| \|\hat{\mathbf{g}}\|} \quad (5)$$

The Euclidean distance has been used for measuring the accuracy of **2D gaze estimation** methods in [42, 144, 183]. We denote the actual gaze position as $\mathbf{p} \in \mathbb{R}^2$ and the estimated gaze position as $\hat{\mathbf{p}} \in \mathbb{R}^2$. We can compute the Euclidean distance as

$$\mathcal{L}_{\text{Euclidean}} = \|\mathbf{p} - \hat{\mathbf{p}}\|_2, \quad (6)$$

Two kinds of evaluation protocols are commonly used for deep-learning based gaze estimation methods, including within-dataset and cross-dataset evaluation. The within-dataset evaluation assesses the model performance on the unseen subjects from the same dataset. The dataset is divided into training and test set according to subjects. There is no overlap in subjects between the training and test set. Note that, most of the gaze datasets provide within-dataset evaluation protocol. They divide the data into training and test set in advance. The **cross-dataset evaluation** assesses the model performance on the unseen environment. The model is **trained on one dataset and tested on another dataset**.

4.4 Public Datasets

We try our best to summarize all the public datasets on gaze estimation, as shown in Tab. 4. The gaze and head pose distribution of these datasets are shown in Fig. 15. Note that, the Gaze360 dataset do not provide the head information. We also discuss three typical datasets that are widely used in gaze estimation studies.

Zhang *et al.* proposed the MPIIGaze [17] dataset. It is the most popular dataset for appearance-based gaze estimation methods. It contains a total of 213,659 images collected from 15 subjects. The images are collected in daily life over several months and there is no constraint for the head pose. MPIIGaze dataset provides both 2D and 3D gaze annotation. It also provides a standard evaluation set, which contains 15 subjects and 3,000 images for each subject. The 3000 images are consisted of 1,500 left-eye images and 1,500 right-eye images. The author further extends the original datasets in [49, 50]. They supply the corresponding face images in [50] and manual landmark annotations in [49].

EyeDiap [185] dataset consists of 94 video clips from 16 participants. It is collected in laboratory environments and has three visual target sessions: continuous moving targets, discrete moving targets, and floating ball. For each subject, they recorded a total of six sessions containing two head movements: static head pose and free head movement. Two cameras are used for data collection: an RGBD camera and an HD camera. The disadvantage of this dataset is the lack of illumination variation.

GazeCapture [42] dataset is collected through crowdsourcing. It contains a total of 2,445,504 images from 1,474 participants. All images are collected using mobile phones or tablets. Each

TABLE 6

Benchmark of cross-domain gaze estimation. ‘Source-free’ indicates whether the method requires source images during domain adaption. ‘Target num’ presents the number of images used for domain adaption. \mathcal{D}_E , \mathcal{D}_G , \mathcal{D}_M , \mathcal{D}_D denotes ETH-XGaze [66], Gaze360 [43], MPIIGaze [50] and EyeDiap [185] datasets. The second-row methods use unannotated images while the third-row methods use annotated images.

| Methods | Pub. | Year | Source-free | Target num | $\mathcal{D}_E \rightarrow \mathcal{D}_M$ | $\mathcal{D}_E \rightarrow \mathcal{D}_D$ | $\mathcal{D}_G \rightarrow \mathcal{D}_M$ | $\mathcal{D}_G \rightarrow \mathcal{D}_D$ |
|------------------|-------|------|-------------|------------|---|---|---|---|
| FullFace [50] | CVPRW | 2017 | | | 11.13° | 14.42° | 12.35° | 30.15° |
| CA-Net [56] | AAAI | 2020 | | | N/A | N/A | <u>27.13°</u> | 31.41° |
| PureGaze [72] | AAAI | 2022 | | | 9.28° | 9.32° | 7.08° | 7.48° |
| RAT [111] | CVPR | 2022 | | | 7.40° | 6.91° | 7.69° | 7.08° |
| PnP-GA [112] | ICCV | 2021 | × | 10 | 6.00° | 6.17° | <u>5.74°</u> | 7.04° |
| CSA [143] | CVPR | 2022 | ✓ | unreport | <u>5.37°</u> | 6.77° | 7.30° | 7.73° |
| CRGA-100 [143] | CVPR | 2022 | × | 100 | 5.68° | <u>5.72°</u> | 6.09° | <u>6.68°</u> |
| CRGA [143] | CVPR | 2022 | × | unreport | <u>5.48°</u> | <u>5.66°</u> | 5.89° | <u>6.49°</u> |
| RUDA [111] | CVPR | 2022 | × | 100 | 5.78° | 5.10° | 6.88° | 6.73° |
| PureGaze-FT [72] | AAAI | 2022 | ✓ | ~ 50 | <u>5.20°</u> | 7.36° | <u>5.30°</u> | <u>6.42°</u> |

TABLE 7
Benchmark of 2D gaze estimation (cm).

| Methods | Pub. | MPIIGaze [50] | EyeDiap [185] | GazeCapture [42] Tablet | Phone |
|------------------------------|---------|---------------|---------------|----------------------------|-------------|
| Itracker [†] [42] | CVPR16 | 7.67 | 10.13 | 2.81 | 1.86 |
| AFF-Net [57] | ICPR20 | <u>4.21</u> | 9.25 | <u>2.30</u> | <u>1.62</u> |
| SAGE [117] | ICCVW19 | N/A | N/A | 2.72 | 1.78 |
| TAT [152] | ICCVW19 | N/A | N/A | <u>2.66</u> | <u>1.77</u> |
| EFE [187] | CVPRW23 | <u>3.89</u> | N/A | <u>2.48</u> | <u>1.61</u> |
| Mnist [†] [17] | CVPR15 | 7.29 | 9.06 | N/A | N/A |
| GazeNet [†] [49] | TPAMI17 | 6.62 | 8.51 | N/A | N/A |
| DilatedNet [†] [53] | ACCV19 | 5.07 | 7.36 | N/A | N/A |
| Gaze360 [43] | ICCV19 | <u>4.66</u> | <u>6.37</u> | N/A | N/A |
| RT-Gene [54] | ECCV18 | 5.36 | <u>7.19</u> | N/A | N/A |
| FullFace [50] | CVPRW17 | 5.65 | 7.70 | N/A | N/A |
| CA-Net [56] | AAAI20 | 4.90 | <u>6.30</u> | N/A | N/A |

participant is required to gaze at a circle shown on the devices without any constraint on their head movement. As a result, the GazeCapture dataset covers various lighting conditions and head motions. The GazeCapture dataset does not provide 3D coordinates of targets. It is usually used for the evaluation of unconstrained 2D gaze point estimation methods.

In addition to the dataset mentioned above, there are several datasets being proposed recently. For example, in 2018, Fischer *et al.* proposed RT-Gene dataset [54]. This dataset provides accurate 3D gaze data since they collect gaze with a dedicated eye tracking device. In 2019, Kellnhofe *et al.* proposed the Gaze360 dataset [43]. The dataset consists of 238 subjects of indoor and outdoor environments with 3D gaze across a wide range of head poses and distances. In 2020, Zhang *et al.* propose the ETH-XGaze dataset [66]. This dataset provides high-resolution images that cover extreme head poses. It also contains 16 illumination conditions for exploring the effects of illumination.

4.5 Benchmarks

We build benchmarks for 2D PoG and 3D gaze estimation in this section. We re-implemented the typical gaze estimation methods as annotated with [†] or report the performance from their manuscripts for comparison. Note that, 2D PoG estimation methods and 3D gaze estimation methods are not comparable since they estimate different forms of gaze. We follow Sec. 4.2.1 to convert the estimation results. We convert 2D PoG into 3D gaze and vice versa. Besides, there are two different gaze definitions in 3D gaze estimation methods. Conventional methods define the origin of

gaze direction as eye centers. Recent methods estimate gaze from face images where the gaze origin is defined as face centers. The difference between the two definitions is minor but makes the direct comparison unfair. We also convert the two definitions with post-processing methods following Sec. 4.2.2. We respectively conduct benchmarks for 2D PoG and 3D gaze estimation. The 3D gaze estimation also are divided into within-dataset and cross-dataset evaluation. We mark the top three performance in all benchmarks with underlines.

Within-dataset evaluation. We first show the comparison of within-dataset evaluation in Tab. 5. The second row contains methods estimating 3D gaze from eye images where the gaze origin is eye center. The methods in the third row estimate 3D gaze from face images. They define the gaze origin as face centers. The last row contains the methods which estimate 2D PoG from face images. Evaluation datasets contain two categories based on data pre-processing process. We obtain eye images from MPIIGaze [49], EyeDiap [185] and UT [37], and evaluate the method which define eye centers as gaze origin in the three datasets. The result is shown in the third column of Tab. 5. We obtain face and eye images from MPIIFaceGaze [50], EyeDiap [185], Gaze360 [43], RT-Gene [54] and ETH-XGaze datasets [66]. We evaluate the method which defines face centers as gaze origin in these datasets. The result is shown in the fourth column of Tab. 5.

Conventional approaches typically estimate gaze using eye images. The Mnist [17] and GazeNet [49] methods employ eye images and head pose vector as input for gaze estimation. Recent methods, *i.e.*, the third-row methods, focus on estimating gaze from facial images. Despite incurring higher computational costs, methods reliant on facial images outperform those centered on eye images. Notably, face image-based methods also usually maintain an acceptable inference speed exceeding 20 frames per second.

Among face image-based methods, GazeTR-Hybrid [34], CA-Net [56] and Gaze360 [43] have better performance. Gaze360 employs ResNet18 for feature extraction while GazeTR-Hybrid adopts a mixed architecture of ResNet18 and transformers. Pre-training significantly enhances the performance of the two methods. In contrast, CA-Net leverages features from both facial and eye images, It requires no pre-training but has a complex network. Regarding datasets, Gaze360 [43] and RT-Gene [54] are collected with large user-camera distances. Most of methods demonstrate significant errors in the two datasets due to low-resolution images. Other datasets are collected with a small user-camera distance or with high-resolution cameras. Appearance-based gaze estimation methods usually achieve approximately 5° in these environments.

Cross-Dataset Evaluation. We conduct four tasks including $\mathcal{D}_E \rightarrow \mathcal{D}_M$, $\mathcal{D}_E \rightarrow \mathcal{D}_D$, $\mathcal{D}_G \rightarrow \mathcal{D}_M$ and $\mathcal{D}_G \rightarrow \mathcal{D}_D$, where \mathcal{D}_E , \mathcal{D}_G , \mathcal{D}_M and \mathcal{D}_D represents ETH-XGaze [66], Gaze360 [43], MPIIFaceGaze [50] and EyeDiap [185] datasets. ETH-XGaze and Gaze360 are used as training set since they have large gaze and head pose ranges. The result is shown in Tab. 6. Unsupervised domain adaption methods are usually proposed to solve the cross-dataset problem. These methods require target images for domain adaption. We summarize the number of required target images. The source-free column indicates whether the method requires source images during domain adaption.

The methods in the second row train models on source datasets without adaption. PureGaze [72] and RAT [111] integrate specific algorithms to enhance model generalization. Their models can be directly applied into multiple domains and achieve reasonable performance. The third row shows the performance of unsupervised domain adaption methods. CRGA [143] and RUDA [111] have better performance while PnP-GA [112] has lower requirement. Compared with the second-row methods, these methods leverage target images to improve the model performance within specific domains. This approach yields a dedicated model for each domain, outperforming PureGaze and RAT. Notably, CSA [143] stands out as a source-free method that dispenses with the need for a source dataset during adaptation. This trend is noteworthy for its implications in privacy protection. PureGaze-FT [72] samples 5 images per person for fine-tuning. Although the method achieves good performance with 50 images, it requires annotated images while previous methods only require unannotated images.

2D PoG estimation. We conduct experiment for 2D PoG estimation. We use MPIIGaze [49], EyeDiap [185] and GazeCapture [42] for evaluation sets and Euclidean distance for evaluation metric. MPIIGaze and EyeDiap datasets collect 2D PoG in screen. The two datasets both provide calibrated screen pose, where we can convert gaze directions to 2D PoG. GazeCapture dataset collects 2D PoG in mobile devices. We count the result based on the types of devices, e.g., tablets and phones. The second row in Tab. 7 shows the result of PoG estimation methods. AFF-Net [57] and EFE [187] shows the best performance than other compared methods. The third and fourth rows show the converted results. Compared methods are designed for gaze direction estimation and we convert the result into PoG. The converted result shows good accuracy in EyeDiap dataset while AFF-Net also shows the best performance in MPIIGaze dataset.

5 CONCLUSIONS AND FUTURE DIRECTIONS

In this survey, we present a comprehensive overview of deep learning-based gaze estimation methods. Unlike the conventional gaze estimation methods that requires dedicated devices, the deep learning-based approaches regress the gaze from the eye appearance captured by web cameras. This makes it easy to implement the algorithm in real world applications. We introduce the gaze estimation method from four perspectives: deep feature extraction, deep neural network architecture design, personal calibration as well as device and platform. We summarize the public datasets on appearance-based gaze estimation and provide benchmarks to compare of the state-of-the-art algorithms. This survey can serve as a guideline for future gaze estimation research.

We further suggest several future directions of deep learning-based gaze estimation. 1) Extracting more robust gaze features. The perfect gaze estimation method should be accurate under

all different subjects, head poses and environments. Therefore, an environment-invariant gaze feature is crucial. 2) Improving performance with fast and simple calibration. There is a trade-off between the system performance and calibration time. The longer calibration time leads to more accurate estimates. How to achieve satisfactory performance with fast calibration procedure is a promising direction. 3) Interpreting learned features. Deep learning approaches often serve as a black-box tool for gaze estimation. Interpretation of the learned features in these methods will bring insight for the deep learning-based gaze estimation.

REFERENCES

- [1] M. K. Eckstein, B. Guerra-Carrillo, A. T. Miller Singley, and S. A. Bunge, "Beyond eye gaze: What else can eyetracking reveal about cognition and cognitive development?" *Developmental Cognitive Neuroscience*, vol. 25, pp. 69–91, 2017.
- [2] G. E. Raptis, C. Katsini, M. Belk, C. Fidas, G. Samaras, and N. Avouris, "Using eye gaze data and visual activities to infer human cognitive styles: Method and feasibility studies," in *Proceedings of the Conference on User Modeling, Adaptation and Personalization*, 2017, p. 164–173.
- [3] M. Meißner and J. Oll, "The promise of eye-tracking methodology in organizational research: A taxonomy, review, and future avenues," *Organizational Research Methods*, vol. 22, no. 2, pp. 590–617, 2019.
- [4] J. Kerr-Gaffney, A. Harrison, and K. Tchanturia, "Eye-tracking research in eating disorders: A systematic review," *International Journal of Eating Disorders*, vol. 52, no. 1, pp. 3–27, 2019.
- [5] A. Palazzi, D. Abati, S. Calderara, F. Solera, and R. Cucchiara, "Predicting the driver's focus of attention: The dr(eye)ve project," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 41, no. 7, pp. 1720–1733, July 2019.
- [6] V. Sitzmann, A. Serrano, A. Pavel, M. Agrawala, D. Gutierrez, B. Masia, and G. Wetzstein, "Saliency in vr: How do people explore virtual environments?" *IEEE Transactions on Visualization and Computer Graphics (TVCG)*, vol. 24, no. 4, pp. 1633–1642, 2018.
- [7] H. Wang, J. Pi, T. Qin, S. Shen, and B. E. Shi, "Slam-based localization of 3d gaze using a mobile eye tracker," in *The ACM Symposium on Eye Tracking Research & Applications (ETRA)*, 2018, pp. 1–5.
- [8] C. Katsini, Y. Abdrabou, G. Raptis, M. Khamis, and F. Alt, "The role of eye gaze in security and privacy applications: Survey and future hci research directions," in *Conference on Human Factors in Computing Systems (CHI)*, 04 2020.
- [9] M. Khamis, F. Alt, and A. Bulling, "The past, present, and future of gaze-enabled handheld mobile devices: Survey and lessons learned," in *International Conference on Human-Computer Interaction with Mobile Devices and Services*, 2018.
- [10] H. Wang, X. Dong, Z. Chen, and B. E. Shi, "Hybrid gaze/eeg brain computer interface for robot arm control on a pick and place task," in *Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*. IEEE, 2015, pp. 1476–1479.
- [11] B. Lai, M. Liu, F. Ryan, and J. Rehg, "In the eye of transformer: Global-local correlation for egocentric gaze estimation," *The British Machine Vision Conference (BMVC)*, 2022.
- [12] L. Jiang, Y. Li, S. Li, M. Xu, S. Lei, Y. Guo, and B. Huang, "Does text attract attention on e-commerce images: A novel saliency prediction dataset and method," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 2088–2097.
- [13] Y. Fang, J. Tang, W. Shen, W. Shen, X. Gu, L. Song, and G. Zhai, "Dual attention guided gaze target detection in the wild," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 11 390–11 399.
- [14] Y. Li, W. Shen, Z. Gao, Y. Zhu, G. Zhai, and G. Guo, "Looking here or there? gaze following in 360-degree images," in *International Conference on Computer Vision (ICCV)*, October 2021, pp. 3742–3751.
- [15] D. Tu, X. Min, H. Duan, G. Guo, G. Zhai, and W. Shen, "End-to-end human-gaze-target detection with transformers," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 2202–2210.
- [16] L. R. Young and D. Sheena, "Survey of eye movement recording methods," *Behavior research methods & instrumentation*, vol. 7, no. 5, pp. 397–429, 1975.
- [17] X. Zhang, Y. Sugano, M. Fritz, and A. Bulling, "Appearance-based gaze estimation in the wild," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [18] T. Eggert, "Eye movement recordings: methods," *Neuro-Ophthalmology*, vol. 40, pp. 15–34, 2007.

- [19] F. Lu, Y. Sugano, T. Okabe, and Y. Sato, "Adaptive linear regression for appearance-based gaze estimation," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 36, no. 10, pp. 2033–2046, 2014.
- [20] F. Martinez, A. Carbone, and E. Pissaloux, "Gaze estimation using local features and non-linear regression," in *International Conference on Image Processing (ICIP)*. IEEE, 2012, pp. 1961–1964.
- [21] Kar-Han Tan, D. J. Kriegman, and N. Ahuja, "Appearance-based eye gaze estimation," in *The IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2002, pp. 191–195.
- [22] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Conference on Neural Information Processing Systems (NeurIPS)*, 2017.
- [23] O. Mowrer, T. C. Ruch, and N. Miller, "The corneo-retinal potential difference as the basis of the galvanometric method of recording eye movements," *American Journal of Physiology-Legacy Content*, vol. 114, no. 2, pp. 423–428, 1936.
- [24] E. Schott, "Über die registrierung des nystagmus und anderer augenbewegungen verm itteltes des saitengalvanometers," *Deut Arch fur klin Med*, vol. 140, pp. 79–90, 1922.
- [25] C. Morimoto and M. Mimica, "Eye gaze tracking techniques for interactive applications," *Computer Vision and Image Understanding*, vol. 98, no. 1, pp. 4–24, 2005.
- [26] D. M. Stampe, "Heuristic filtering and reliable calibration methods for video-based pupil-tracking systems," *Behavior Research Methods, Instruments, & Computers*, vol. 25, no. 2, pp. 137–142, 1993.
- [27] J. Qiang and X. Yang, "Real-time eye, gaze, and face pose tracking for monitoring driver vigilance," *Real-Time Imaging*, vol. 8, no. 5, pp. 357–377, 2002.
- [28] E. D. Guestrin and M. Eizenman, "General theory of remote gaze estimation using the pupil center and corneal reflections," *IEEE Transactions on Biomedical Engineering*, vol. 53, no. 6, pp. 1124–1133, 2006.
- [29] Z. Zhu and Q. Ji, "Novel eye gaze tracking techniques under natural head movement," *IEEE Transactions on Biomedical Engineering*, vol. 54, no. 12, pp. 2246–2260, 2007.
- [30] R. Valenti, N. Sebe, and T. Gevers, "Combining head pose and eye location information for gaze estimation," *IEEE Transactions on Image Processing (TIP)*, vol. 21, no. 2, pp. 802–815, 2012.
- [31] K. A. Funes Mora and J. Odobez, "Geometric generative gaze estimation (g3e) for remote rgb-d cameras," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [32] S. Baluja and D. Pomerleau, "Non-intrusive gaze tracking using artificial neural networks," in *Conference on Neural Information Processing Systems (NeurIPS)*, 1994.
- [33] O. Williams, A. Blake, and R. Cipolla, "Sparse and semi-supervised visual mapping with the s^3gp ," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2006.
- [34] Y. Cheng and F. Lu, "Gaze estimation using transformer," *International Conference on Pattern Recognition (ICPR)*, 2022.
- [35] Y. Sugano, Y. Matsushita, and Y. Sato, "Appearance-based gaze estimation using visual saliency," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 35, no. 2, pp. 329–341, 2013.
- [36] K. A. Funes Mora and J. Odobez, "Person independent 3d gaze estimation from remote rgb-d cameras," in *International Conference on Image Processing (ICIP)*, 2013, pp. 2787–2791.
- [37] Y. Sugano, Y. Matsushita, and Y. Sato, "Learning-by-synthesis for appearance-based 3d gaze estimation," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [38] F. Lu and X. Chen, "Person-independent eye gaze prediction from eye images using patch-based features," *Neurocomputing*, vol. 182, pp. 10–17, 2016.
- [39] Y. Sugano, Y. Matsushita, Y. Sato, and H. Koike, "An incremental learning method for unconstrained gaze estimation," in *The European Conference on Computer Vision (ECCV)*, 2008, p. 656–667.
- [40] F. Lu, T. Okabe, Y. Sugano, and Y. Sato, "Learning gaze biases with head motion for head pose-free gaze estimation," *Image and Vision Computing*, vol. 32, no. 3, pp. 169–179, 2014.
- [41] F. Lu, Y. Sugano, T. Okabe, and Y. Sato, "Gaze estimation from eye appearance: A head pose-free method via eye image synthesis," *IEEE Transactions on Image Processing (TIP)*, vol. 24, no. 11, pp. 3680–3693, Nov 2015.
- [42] K. Krafka, A. Khosla, P. Kellnhofer, H. Kannan, S. Bhandarkar, W. Matusik, and A. Torralba, "Eye tracking for everyone," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [43] P. Kellnhofer, A. Recasens, S. Stent, W. Matusik, and A. Torralba, "Gaze360: Physically unconstrained gaze estimation in the wild," in *International Conference on Computer Vision (ICCV)*, 2019.
- [44] Y. Cheng, X. Zhang, F. Lu, and Y. Sato, "Gaze estimation by exploring two-eye asymmetry," *IEEE Transactions on Image Processing (TIP)*, vol. 29, pp. 5259–5272, 2020.
- [45] S. Park, A. Spurr, and O. Hilliges, "Deep pictorial gaze estimation," in *The European Conference on Computer Vision (ECCV)*, 2018.
- [46] Y. Yu, G. Liu, and J.-M. Odobez, "Improving few-shot user-specific gaze adaptation via gaze redirection synthesis," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [47] S. Park, S. D. Mello, P. Molchanov, U. Iqbal, O. Hilliges, and J. Kautz, "Few-shot adaptive gaze estimation," in *International Conference on Computer Vision (ICCV)*, 2019.
- [48] M. Xu, H. Wang, Y. Liu, and F. Lu, "Vulnerability of appearance-based gaze estimation," *arXiv preprint arXiv:2103.13134*, 2021.
- [49] X. Zhang, Y. Sugano, M. Fritz, and A. Bulling, "Mpiigaze: Real-world dataset and deep appearance-based gaze estimation," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 41, no. 1, pp. 162–175, Jan 2019.
- [50] —, "It's written all over your face: Full-face appearance-based gaze estimation," in *The IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2017, pp. 2299–2308.
- [51] L.-Q. Xu, D. Machin, and P. Sheppard, "A novel approach to real-time non-intrusive gaze finding," in *The British Machine Vision Conference (BMVC)*, 1998.
- [52] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [53] Z. Chen and B. E. Shi, "Appearance-based gaze estimation using dilated-convolutions," in *Asian Conference on Computer Vision (ACCV)*, C. Jawahar, H. Li, G. Mori, and K. Schindler, Eds., 2019, pp. 309–324.
- [54] T. Fischer, H. Jin Chang, and Y. Demiris, "Rt-gene: Real-time eye gaze estimation in natural environments," in *The European Conference on Computer Vision (ECCV)*, 2018.
- [55] Y. Cheng, F. Lu, and X. Zhang, "Appearance-based gaze estimation via evaluation-guided asymmetric regression," in *The European Conference on Computer Vision (ECCV)*, 2018.
- [56] Y. Cheng, S. Huang, F. Wang, C. Qian, and F. Lu, "A coarse-to-fine adaptive network for appearance-based gaze estimation," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2020.
- [57] Y. Bao, Y. Cheng, Y. Liu, and F. Lu, "Adaptive feature fusion network for gaze tracking in mobile tablets," in *International Conference on Pattern Recognition (ICPR)*, 2020.
- [58] P. Biswas *et al.*, "Appearance-based gaze estimation using attention and difference mechanism," in *The IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2021, pp. 3143–3152.
- [59] K. Wang, R. Zhao, H. Su, and Q. Ji, "Generalizing eye tracking with bayesian adversarial learning," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [60] J.-H. Kim and J.-W. Jeong, "Gaze estimation in the dark with generative adversarial networks," in *The ACM Symposium on Eye Tracking Research & Applications (ETRA)*, 2020, pp. 1–3.
- [61] A. Rangesh, B. Zhang, and M. M. Trivedi, "Driver gaze estimation in the real world: Overcoming the eyeglass challenge," in *The IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2020, pp. 1054–1059.
- [62] Y. Yu and J.-M. Odobez, "Unsupervised representation learning for gaze estimation," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [63] H. Yamazoe, A. Utsumi, T. Yonezawa, and S. Abe, "Remote gaze estimation with a single camera based on facial-feature tracking without special calibration actions," in *The ACM Symposium on Eye Tracking Research & Applications (ETRA)*, 2008, p. 245–250.
- [64] J. Chen and Q. Ji, "3d gaze estimation with a single camera without ir illumination," in *International Conference on Pattern Recognition (ICPR)*. IEEE, 2008, pp. 1–4.
- [65] L. A. Jeni and J. F. Cohn, "Person-independent 3d gaze estimation using face frontalization," in *The IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2016, pp. 87–95.
- [66] X. Zhang, S. Park, T. Beeler, D. Bradley, S. Tang, and O. Hilliges, "Ethxgaze: A large scale dataset for gaze estimation under extreme head pose and gaze variation," in *The European Conference on Computer Vision (ECCV)*, 2020.
- [67] S. Park, E. Aksan, X. Zhang, and O. Hilliges, "Towards end-to-end video-based eye-tracking," in *The European Conference on Computer Vision (ECCV)*. Springer, 2020, pp. 747–763.
- [68] A. Mishra and H.-T. Lin, "360-degree gaze estimation in the wild using multiple zoom scales," *arXiv preprint arXiv:2009.06924*, 2020.
- [69] M. Zhang, Y. Liu, and F. Lu, "Gazeonce: Real-time multi-person gaze

- estimation,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 4197–4206.
- [70] R. Ogusu and T. Yamanaka, “Lpm: Learnable pooling module for efficient full-face gaze estimation,” in *IEEE International Conference on Automatic Face Gesture Recognition*, May 2019, pp. 1–5.
- [71] X. Zhang, Y. Sugano, A. Bulling, and O. Hilliges, “Learning-based region selection for end-to-end gaze estimation,” in *The British Machine Vision Conference (BMVC)*, 2020.
- [72] Y. Cheng, Y. Bao, and F. Lu, “Puregaze: Purifying gaze feature for generalizable gaze estimation,” *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2022.
- [73] Z. Yu, X. Huang, X. Zhang, H. Shen, Q. Li, W. Deng, J. Tang, Y. Yang, and J. Ye, “A multi-modal approach for driver gaze prediction to remove identity bias,” in *Proceedings of the ACM International Conference on Multimodal Interaction (ICMI)*, 2020, pp. 768–776.
- [74] Y. Zhang, X. Yang, and Z. Ma, “Driver’s gaze zone estimation method: A four-channel convolutional neural network model,” in *The International Conference on Big-data Service and Intelligent Computation*, 2020, pp. 20–24.
- [75] Z. Wang, J. Zhao, C. Lu, F. Yang, H. Huang, Y. Guo *et al.*, “Learning to detect head movement in unconstrained remote gaze estimation in the wild,” in *The IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2020, pp. 3443–3452.
- [76] H. Deng and W. Zhu, “Monocular free-head 3d gaze tracking with deep learning and geometry constraints,” in *International Conference on Computer Vision (ICCV)*, 2017.
- [77] K. Cho, B. Van Merriënboer, D. Bahdanau, and Y. Bengio, “On the properties of neural machine translation: Encoder-decoder approaches,” *arXiv preprint*, 2014.
- [78] X. Cai, B. Chen, J. Zeng, J. Zhang, Y. Sun, X. Wang, Z. Ji, X. Liu, X. Chen, and S. Shan, “Gaze estimation with an ensemble of four architectures,” *arXiv preprint arXiv:2107.01980*, 2021.
- [79] C. Palmero, J. Selva, M. A. Bagheri, and S. Escalera, “Recurrent cnn for 3d gaze estimation using appearance and shape cues,” in *The British Machine Vision Conference (BMVC)*, 2018.
- [80] P. A. Dias, D. Malafronte, H. Medeiros, and F. Odone, “Gaze estimation for assisted living environments,” in *The IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2020.
- [81] P. Her, L. Manderle, P. A. Dias, H. Medeiros, and F. Odone, “Uncertainty-aware gaze tracking for assisted living environments,” *IEEE Transactions on Image Processing (TIP)*, 2023.
- [82] S. Jyoti and A. Dhall, “Automatic eye gaze estimation using geometric & texture-based networks,” in *International Conference on Pattern Recognition (ICPR)*, Aug 2018, pp. 2474–2479.
- [83] N. Dubey, S. Ghosh, and A. Dhall, “Unsupervised learning of eye gaze representation from the web,” in *International Joint Conference on Neural Networks*, July 2019, pp. 1–7.
- [84] X. Zhou, J. Lin, J. Jiang, and S. Chen, “Learning a 3d gaze estimator with improved itracker combined with bidirectional lstm,” in *IEEE International Conference on Multimedia and Expo*, July 2019, pp. 850–855.
- [85] Z. Wang, J. Chai, and S. Xia, “Realtime and accurate 3d eye gaze capture with dcnn-based iris and pupil segmentation,” *IEEE Transactions on Visualization and Computer Graphics (TVCG)*, pp. 1–1, 2019.
- [86] K. Wang, H. Su, and Q. Ji, “Neuro-inspired eye tracking with eye movement dynamics,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [87] K. P. Murphy and S. Russell, “Dynamic bayesian networks: representation, inference and learning,” 2002.
- [88] C. Farabet, C. Couprie, L. Najman, and Y. LeCun, “Learning hierarchical features for scene labeling,” *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 35, no. 8, pp. 1915–1929, 2012.
- [89] S. Liu, D. Liu, and H. Wu, “Gaze estimation with multi-scale channel and spatial attention,” in *The International Conference on Computing and Pattern Recognition*, 2020, pp. 303–309.
- [90] B. Mahanama, Y. Jayawardana, and S. Jayarathna, “Gaze-net: appearance-based gaze estimation using capsule networks,” in *The Augmented Human International Conference*, 2020, pp. 1–4.
- [91] J. Lemley, A. Kar, A. Drimbarean, and P. Corcoran, “Convolutional neural network implementation for eye-gaze estimation on low-quality consumer imaging systems,” *IEEE Transactions on Consumer Electronics*, vol. 65, no. 2, pp. 179–187, 2019.
- [92] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [93] Y. Zhuang, Y. Zhang, and H. Zhao, “Appearance-based gaze estimation using separable convolution neural networks,” in *The IEEE Advanced Information Technology, Electronic and Automation Control Conference (IAEAC)*, vol. 5. IEEE, 2021, pp. 609–612.
- [94] A. Bublea and C. D. Căleanu, “Deep learning based eye gaze tracking for automotive applications: An auto-keras approach,” in *The International Symposium on Electronics and Telecommunications (ISETC)*. IEEE, 2020, pp. 1–4.
- [95] K. Ruhland, S. Andrist, J. Badler, C. Peters, N. Badler, M. Gleicher, B. Mutlu, and R. McDonnell, “Look me in the eyes: A survey of eye and gaze animation for virtual agents and artificial systems,” in *Eurographics*, Apr. 2014, pp. 69–91.
- [96] L. Swirski and N. Dodgson, “Rendering synthetic ground truth images for eye tracker evaluation,” in *The ACM Symposium on Eye Tracking Research & Applications (ETRA)*, ser. ETRA ’14, 2014, p. 219–222.
- [97] S. Porta, B. Bossavit, R. Cabeza, A. Larumbe-Bergera, G. Garde, and A. Villanueva, “U2eyes: A binocular dataset for eye tracking and gaze estimation,” in *International Conference on Computer Vision Workshops (ICCVW)*, Oct 2019, pp. 3660–3664.
- [98] Y. Furukawa and J. Ponce, “Accurate, dense, and robust multiview stereo,” *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 32, no. 8, pp. 1362–1376, 2009.
- [99] E. Wood, T. Baltrušaitis, X. Zhang, Y. Sugano, P. Robinson, and A. Bulling, “Rendering of eyes for eye-shape registration and gaze estimation,” in *International Conference on Computer Vision (ICCV)*, 2015.
- [100] E. Wood, T. Baltrušaitis, L.-P. Morency, P. Robinson, and A. Bulling, “Learning an appearance-based gaze estimator from one million synthesised images,” in *The ACM Symposium on Eye Tracking Research & Applications (ETRA)*, 2016, p. 131–138.
- [101] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb, “Learning from simulated and unsupervised images through adversarial training,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [102] K. Wang, R. Zhao, and Q. Ji, “A hierarchical generative model for eye image synthesis and eye gaze estimation,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [103] S. Jindal and X. E. Wang, “Cuda-gr: Controllable unsupervised domain adaptation for gaze redirection,” *arXiv preprint arXiv:2106.10852*, 2021.
- [104] Y. Yu, G. Liu, and J.-M. Odobez, “Improving few-shot user-specific gaze adaptation via gaze redirection synthesis,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 11 937–11 946.
- [105] Z. He, A. Spurr, X. Zhang, and O. Hilliges, “Photo-realistic monocular gaze redirection using generative adversarial networks,” in *International Conference on Computer Vision (ICCV)*, 2019.
- [106] Y. Zheng, S. Park, X. Zhang, S. De Mello, and O. Hilliges, “Self-learning transformations for improving gaze and head redirection,” *Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [107] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, “Nerf: Representing scenes as neural radiance fields for view synthesis,” in *The European Conference on Computer Vision (ECCV)*, 2020.
- [108] A. Ruzzi, X. Shi, X. Wang, G. Li, S. De Mello, H. J. Chang, X. Zhang, and O. Hilliges, “Gazenerf: 3d-aware gaze redirection with neural radiance fields,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [109] P. Yin, J. Dai, J. Wang, D. Xie, and S. Pu, “Nerf-gaze: A head-eye redirection parametric model for gaze estimation,” *arXiv preprint arXiv:2212.14710*, 2022.
- [110] R. Kothari, S. De Mello, U. Iqbal, W. Byeon, S. Park, and J. Kautz, “Weakly-supervised physically unconstrained gaze estimation,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 9980–9989.
- [111] Y. Bao, Y. Liu, H. Wang, and F. Lu, “Generalizing gaze estimation with rotation consistency,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 4207–4216.
- [112] Y. Liu, R. Liu, H. Wang, and F. Lu, “Generalizing gaze estimation with outlier-guided collaborative adaptation,” in *International Conference on Computer Vision (ICCV)*, 2021, pp. 3835–3844.
- [113] M. J. Marin-Jimenez, V. Kalogeiton, P. Medina-Suarez, and A. Zisserman, “Lao-net: Revisiting people looking at each other in videos,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [114] Y. Xiong, H. J. Kim, and V. Singh, “Mixed effects neural networks (menets) with applications to gaze estimation,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [115] M. J. Beal, “Variational algorithms for approximate bayesian inference,”

- Ph.D. dissertation, UCL (University College London), 2003.
- [116] H. Robbins and S. Monro, "A stochastic approximation method," *The annals of mathematical statistics*, pp. 400–407, 1951.
 - [117] J. He, K. Pham, N. Valliappan, P. Xu, C. Roberts, D. Lagun, and V. Navalpakkam, "On-device few-shot personalization for real-time gaze estimation," in *International Conference on Computer Vision Workshops (ICCVW)*, Oct 2019, pp. 1149–1158.
 - [118] Y. Sun, J. Zeng, S. Shan, and X. Chen, "Cross-encoder for unsupervised gaze representation learning," in *International Conference on Computer Vision (ICCV)*, October 2021, pp. 3702–3711.
 - [119] Y. Yu, G. Liu, and J.-M. Odobez, "Deep multitask gaze estimation with a constrained landmark-gaze model," in *The European Conference on Computer Vision Workshops (ECCVW)*, 2018.
 - [120] S. Ruder, "An overview of multi-task learning in deep neural networks," *arXiv preprint*, 2017.
 - [121] D. Lian, L. Hu, W. Luo, Y. Xu, L. Duan, J. Yu, and S. Gao, "Multiview multitask gaze estimation with deep convolutional neural networks," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 10, pp. 3010–3023, Oct 2019.
 - [122] D. Lian, Z. Zhang, W. Luo, L. Hu, M. Wu, Z. Li, J. Yu, and S. Gao, "RgbD based gaze estimation via multi-task cnn," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, vol. 33, 2019, pp. 2488–2495.
 - [123] Z. Wu, S. Rajendran, T. V. As, V. Badrinarayanan, and A. Rabinovich, "Eyenet: A multi-task deep network for off-axis eye gaze estimation," in *International Conference on Computer Vision Workshops (ICCVW)*, Oct 2019, pp. 3683–3687.
 - [124] A. Recasens, C. Vondrick, A. Khosla, and A. Torralba, "Following gaze in video," in *International Conference on Computer Vision (ICCV)*, 2017.
 - [125] S. S. Kruthiventi, V. Gudisa, J. H. Dholakiya, and R. Venkatesh Babu, "Saliency unified: A deep architecture for simultaneous eye fixation prediction and salient object segmentation," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 5781–5790.
 - [126] W. Wang, J. Shen, X. Dong, A. Borji, and R. Yang, "Inferring salient objects from human fixations," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 42, no. 8, pp. 1913–1927, 2020.
 - [127] E. Chong, N. Ruiz, Y. Wang, Y. Zhang, A. Rozga, and J. M. Rehg, "Connecting gaze, scene, and attention: Generalized attention estimation via joint modeling of gaze and scene saliency," in *The European Conference on Computer Vision (ECCV)*, 2018.
 - [128] A. Graves, S. Fernández, and J. Schmidhuber, "Bidirectional lstm networks for improved phoneme classification and recognition," in *Proceedings of International Conference on Artificial Neural Networks*, ser. ICANN'05. Springer-Verlag, 2005, p. 799–804.
 - [129] Y. Xu, Y. Dong, J. Wu, Z. Sun, Z. Shi, J. Yu, and S. Gao, "Gaze prediction in dynamic 360° immersive videos," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
 - [130] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
 - [131] Z. Chen and B. Shi, "Offset calibration for appearance-based gaze estimation via gaze decomposition," in *The IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2020.
 - [132] Z. Chen and B. E. Shi, "Towards high performance low complexity calibration in appearance based gaze estimation," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 45, no. 1, pp. 1174–1188, 2022.
 - [133] X. Zhang, M. X. Huang, Y. Sugano, and A. Bulling, "Training person-specific gaze estimators from user interactions with multiple devices," in *Conference on Human Factors in Computing Systems (CHI)*, 2018.
 - [134] Y. Li, Y. Zhan, and Z. Yang, "Evaluation of appearance-based eye tracking calibration data selection," in *The IEEE International Conference on Artificial Intelligence and Computer Applications (ICAICA)*. IEEE, 2020, pp. 222–224.
 - [135] X. Zhang, Y. Sugano, and A. Bulling, "Evaluation of appearance-based methods and implications for gaze-based applications," in *Conference on Human Factors in Computing Systems (CHI)*, 2019.
 - [136] E. Lindén, J. Sjöstrand, and A. Proutiere, "Learning to personalize in appearance-based gaze tracking," in *International Conference on Computer Vision Workshops (ICCVW)*, Oct 2019, pp. 1140–1148.
 - [137] Z. Guo, Z. Yuan, C. Zhang, W. Chi, Y. Ling, and S. Zhang, "Domain adaptation gaze estimation by embedding with prediction consistency," in *Asian Conference on Computer Vision (ACCV)*, 2020.
 - [138] W. Cui, J. Cui, and H. Zha, "Specialized gaze estimation for children by convolutional neural network and domain adaptation," in *International Conference on Image Processing (ICIP)*, Sep. 2017, pp. 3305–3309.
 - [139] B. Gong, Y. Shi, F. Sha, and K. Grauman, "Geodesic flow kernel for unsupervised domain adaptation," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 2066–2073.
 - [140] J. Bao, B. Liu, and J. Yu, "An individual-difference-aware model for cross-person gaze estimation," *IEEE Transactions on Image Processing (TIP)*, 2022.
 - [141] G. Liu, Y. Yu, K. A. Funes Mora, and J. Odobez, "A differential approach for gaze estimation," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, pp. 1–1, 2019.
 - [142] A. Tarvainen and H. Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," in *Conference on Neural Information Processing Systems (NeurIPS)*, vol. 30, 2017.
 - [143] Y. Wang, Y. Jiang, J. Li, B. Ni, W. Dai, C. Li, H. Xiong, and T. Li, "Contrastive regression for domain adaptation on gaze estimation," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 19 376–19 385.
 - [144] Z. Chang, M. D. Martino, Q. Qiu, S. Espinosa, and G. Sapiro, "Salgaze: Personalizing gaze estimation using visual saliency," in *International Conference on Computer Vision Workshops (ICCVW)*, Oct 2019, pp. 1169–1178.
 - [145] K. Wang, S. Wang, and Q. Ji, "Deep eye fixation map learning for calibration-free eye gaze tracking," in *The ACM Symposium on Eye Tracking Research & Applications (ETRA)*, 2016, p. 47–55.
 - [146] B. Klein Salvalaio and G. de Oliveira Ramos, "Self-adaptive appearance-based eye-tracking with online transfer learning," in *Brazilian Conference on Intelligent Systems*, Oct 2019, pp. 383–388.
 - [147] M. Tonsen, J. Steil, Y. Sugano, and A. Bulling, "Invisibleeye: Mobile eye tracking using multiple low-resolution cameras and learning-based gaze estimation," *ACM Transactions on Interactive, Mobile, Wearable and Ubiquitous Technologies (IMWUT)*, vol. 1, no. 3, Sep. 2017.
 - [148] J. Lemley, A. Kar, and P. Corcoran, "Eye tracking in augmented spaces: A deep learning approach," in *IEEE Games, Entertainment, Media Conference*, Aug 2018, pp. 1–6.
 - [149] J. Kim, M. Stengel, A. Majercik, S. De Mello, D. Dunn, S. Laine, M. McGuire, and D. Luebke, "Nvgaze: An anatomically-informed dataset for low-latency, near-eye gaze estimation," in *Conference on Human Factors in Computing Systems (CHI)*, 2019.
 - [150] J. Li, Z. Chen, Y. Zhong, H.-K. Lam, J. Han, G. Ouyang, X. Li, and H. Liu, "Appearance-based gaze estimation for asd diagnosis," *IEEE Transactions on Cybernetics (TC)*, vol. 52, no. 7, pp. 6504–6517, 2022.
 - [151] S. Ghosh, M. Hayat, A. Dhall, and J. Knibbe, "Mtgls: Multi-task gaze estimation with limited supervision," in *The IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2022, pp. 3223–3234.
 - [152] T. Guo, Y. Liu, H. Zhang, X. Liu, Y. Kwak, B. I. Yoo, J.-J. Han, and C. Choi, "A generalized and robust method towards practical gaze estimation on smart phone," in *International Conference on Computer Vision Workshops (ICCVW)*, Oct 2019, pp. 1131–1139.
 - [153] Z. Zhao, S. Li, and T. Kosaki, "Estimating a driver's gaze point by a remote spherical camera," in *The IEEE International Conference on Mechatronics and Automation (ICMA)*. IEEE, 2020, pp. 599–604.
 - [154] Y. Xia and B. Liang, "Gaze estimation based on deep learning method," in *The International Conference on Computer Science and Application Engineering*, 2020, pp. 1–6.
 - [155] C.-S. Chen, H.-T. Lin *et al.*, "360-degree gaze estimation in the wild using multiple zoom scales," *The British Machine Vision Conference (BMVC)*, 2020.
 - [156] S. Nonaka, S. Nobuhara, and K. Nishino, "Dynamic 3d gaze from afar: Deep gaze estimation from temporal eye-head-body coordination," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 2192–2201.
 - [157] I. Lee, J.-S. Yun, H. H. Kim, Y. Na, and S. B. Yoo, "Latentgaze: Cross-domain gaze estimation through gaze-aware analytic latent code manipulation," in *The British Machine Vision Conference (BMVC)*, 2022, pp. 3379–3395.
 - [158] I. Kasahara, S. Stent, and H. S. Park, "Look both ways: Self-supervising driver gaze estimation and road scene saliency," in *The European Conference on Computer Vision (ECCV)*. Springer, 2022, pp. 126–142.
 - [159] J.-S. Yun, Y. Na, H. H. Kim, H.-I. Kim, and S. B. Yoo, "Haze-net: High-frequency attentive super-resolved gaze estimation in low-resolution face images," in *Asian Conference on Computer Vision (ACCV)*, 2022, pp. 3361–3378.
 - [160] A. Farkhondeh, C. Palmero, S. Scardapane, and S. Escalera, "Towards self-supervised gaze estimation," *The British Machine Vision Conference (BMVC)*, 2022.

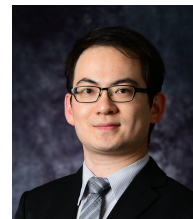
- [161] J. O. Oh, H. J. Chang, and S.-I. Choi, "Self-attention with convolution and deconvolution for efficient eye gaze estimation from a full face image," in *The IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, June 2022, pp. 4992–5000.
- [162] J. Qin, T. Shimoyama, and Y. Sugano, "Learning-by-novel-view-synthesis for full-face appearance-based 3d gaze estimation," in *The IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2021, pp. 4981–4991.
- [163] J. Gideon, S. Su, and S. Stent, "Unsupervised multi-view gaze representation learning," in *The IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, June 2022, pp. 5001–5009.
- [164] Y. Cheng and F. Lu, "Dvgaze: Dual-view gaze estimation," in *International Conference on Computer Vision (ICCV)*, October 2023, pp. 20632–20641.
- [165] S. Khalighy, G. Green, C. Scheepers, and C. Whittet, "Quantifying the qualities of aesthetics in product design using eye-tracking technology," *International Journal of Industrial Ergonomics*, vol. 49, pp. 31 – 43, 2015.
- [166] R. d. O. J. dos Santos, J. H. C. de Oliveira, J. B. Rocha, and J. d. M. E. Giraldi, "Eye tracking in neuromarketing: a research agenda for marketing studies," *International journal of psychological studies*, vol. 7, no. 1, p. 32, 2015.
- [167] X. Zhang, Y. Sugano, and A. Bulling, "Everyday eye contact detection using unsupervised gaze target discovery," in *The ACM Symposium on User Interface Software and Technology (UIST)*, 2017, p. 193–203.
- [168] Y. Sugano, X. Zhang, and A. Bulling, "Aggregaze: Collective estimation of audience attention on public displays," in *The ACM Symposium on User Interface Software and Technology (UIST)*, 2016, p. 821–831.
- [169] J. Zhang and S. Sclaroff, "Exploiting surroundedness for saliency detection: A boolean map approach," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 38, no. 5, pp. 889–902, May 2016.
- [170] N. Valliappan, N. Dai, E. Steinberg, J. He, K. Rogers, V. Ramachandran, P. Xu, M. Shojaeizadeh, L. Guo, K. Kohlhoff *et al.*, "Accelerating eye movement research via accurate and affordable smartphone eye tracking," *Nature communications*, vol. 11, no. 1, pp. 1–12, 2020.
- [171] V. Kazemi and J. Sullivan, "One millisecond face alignment with an ensemble of regression trees," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014, pp. 1867–1874.
- [172] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, "Joint face detection and alignment using multitask cascaded convolutional networks," *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1499–1503, 2016.
- [173] M. Kowalski, J. Naruniec, and T. Trzcinski, "Deep alignment network: A convolutional neural network for robust face alignment," in *The IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, July 2017.
- [174] T. Baltrusaitis, A. Zadeh, Y. C. Lim, and L. Morency, "Openface 2.0: Facial behavior analysis toolkit," in *IEEE International Conference on Automatic Face Gesture Recognition*, 2018, pp. 59–66.
- [175] Y. Feng, F. Wu, X. Shao, Y. Wang, and X. Zhou, "Joint 3d face reconstruction and dense alignment with position map regression network," in *The European Conference on Computer Vision (ECCV)*, 2018.
- [176] J. Guo, X. Zhu, Y. Yang, F. Yang, Z. Lei, and S. Z. Li, "Towards fast, accurate and stable 3d dense face alignment," in *The European Conference on Computer Vision (ECCV)*, 2020.
- [177] J. Zhang, H. Hu, and S. Feng, "Robust facial landmark detection via heatmap-offset regression," *IEEE Transactions on Image Processing (TIP)*, vol. 29, pp. 5050–5064, 2020.
- [178] P. Chandran, D. Bradley, M. Gross, and T. Beeler, "Attention-driven cropping for very high resolution facial landmark detection," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [179] P. Gao, K. Lu, J. Xue, L. Shao, and J. Lyu, "A coarse-to-fine facial landmark detection method based on self-attention mechanism," *IEEE Transactions on Multimedia*, pp. 1–1, 2020.
- [180] X. Zhang, Y. Sugano, and A. Bulling, "Revisiting data normalization for appearance-based gaze estimation," in *The ACM Symposium on Eye Tracking Research & Applications (ETRA)*, 2018.
- [181] H. Wang and B. E. Shi, "Gaze awareness improves collaboration efficiency in a collaborative assembly task," in *The ACM Symposium on Eye Tracking Research & Applications (ETRA)*, 2019, pp. 1–5.
- [182] X. Dong, H. Wang, Z. Chen, and B. E. Shi, "Hybrid brain computer interface via bayesian integration of eeg and eye gaze," in *International IEEE/EMBS Conference on Neural Engineering (NER)*. IEEE, 2015, pp. 150–153.
- [183] E. T. Wong, S. Yean, Q. Hu, B. S. Lee, J. Liu, and R. Deepu, "Gaze estimation using residual neural network," in *IEEE International*

Conference on Pervasive Computing and Communications Workshops, 2019, pp. 411–414.

- [184] B. A. Smith, Q. Yin, S. K. Feiner, and S. K. Nayar, "Gaze locking: Passive eye contact detection for human-object interaction," in *The ACM Symposium on User Interface Software and Technology (UIST)*, 2013, p. 271–280.
- [185] K. A. Funes Mora, F. Monay, and J.-M. Odobez, "Eyediap: A database for the development and evaluation of gaze estimation algorithms from rgb and rgb-d cameras," in *The ACM Symposium on Eye Tracking Research & Applications (ETRA)*, 2014.
- [186] Q. Huang, A. Veeraraghavan, and A. Sabharwal, "Tabletgaze: dataset and analysis for unconstrained appearance-based gaze estimation in mobile tablets," *Machine Vision and Applications*, vol. 28, no. 5-6, pp. 445–461, 2017.
- [187] H. Balim, S. Park, X. Wang, X. Zhang, and O. Hilliges, "Efe: End-to-end frame-to-gaze estimation," in *The IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, June 2023, pp. 2688–2697.



Yihua Cheng received the B.S. degree in computer science from Beijing University of Posts and Telecommunications in 2017, and Ph.D. degree in computer science from Beihang University in 2022. He is now a Postdoctoral Researcher with University of Birmingham, UK. His research interests include gaze estimation, hand pose estimation, object pose estimation and human-robot interaction.



Haofei Wang received the B.S. degree with distinction from Zhejiang University in 2013, and Ph.D. degree in electronic and computer engineering from the Hong Kong University of Science and Technology in 2020. He is now a Postdoctoral Researcher with the Pengcheng Laboratory, Shenzhen, China. His research interests include eye tracking, gaze estimation, human-computer interaction and mixed reality.



Yiwei Bao currently pursuing the Ph.D. degree with the State Key Laboratory of Virtual Reality Technology and Systems, School of Computer Science and Engineering, Beihang University. His research interests include computer vision and human gaze analysis.



Feng Lu received the B.S. and M.S. degrees in automation from Tsinghua University, in 2007 and 2010, respectively, and the Ph.D. degree in information science and technology from The University of Tokyo, in 2013. He is currently a Professor with the State Key Laboratory of Virtual Reality Technology and Systems, School of Computer Science and Engineering, Beihang University. His research interests include computer vision, human-computer interaction and augmented intelligence.