

# CrossGaze: A Strong Method for 3D Gaze Estimation in the Wild

Andy Cătrună, Adrian Cosma, Emilian Rădoi

andy.eduard.catruna@upb.ro, ioan-adrian.cosma@upb.ro, emilian.radoi@upb.ro

Faculty of Automatic Control and Computer Science, POLITEHNICA Bucharest

**Abstract**—Gaze estimation, the task of predicting where an individual is looking, is a critical task with direct applications in areas such as human-computer interaction and virtual reality. Estimating the direction of looking in unconstrained environments is difficult, due to the many factors that can obscure the face and eye regions. In this work we propose CrossGaze, a strong baseline for gaze estimation, that leverages recent developments in computer vision architectures and attention-based modules. Unlike previous approaches, our method does not require a specialized architecture, utilizing already established models that we integrate in our architecture and adapt for the task of 3D gaze estimation. This approach allows for seamless updates to the architecture as any module can be replaced with more powerful feature extractors. On the Gaze360 benchmark, our model surpasses several state-of-the-art methods, achieving a mean angular error of  $9.94^\circ$ . Our proposed model serves as a strong foundation for future research and development in gaze estimation, paving the way for practical and accurate gaze prediction in real-world scenarios.

## I. INTRODUCTION

Human gaze estimation, the ability to infer the direction of a person’s gaze from visual cues, is a critical aspect in understanding a person’s intent, engagement and attention, which is highly valuable in fields such as automotive safety [22], human-robot interaction [19] and virtual reality [24]. Gaze estimation scenarios usually fall into two broad categories [21]: model-based and appearance-based. Model-based approaches employ specialized hardware and sensors such as near-infrared cameras (NIR) and are performed in constrained environments. Appearance-based methods, on the other hand, do not require specialized hardware and are meant to be used in unconstrained, real-world scenarios. While gaze estimation has been well-studied in controlled laboratory settings, where participants’ head and eye movements can be tightly controlled, gaze estimation in the wild presents unique challenges due to the uncontrolled and dynamic nature of real-world environments. In the wild, individuals exhibit a wide range of head poses, lighting conditions, occlusions and distractions, making gaze estimation a challenging problem to solve. For commercial applications, gaze estimation in unconstrained environments can be essential for gathering actionable insights into customer behaviour [2]. Coupled with other unintrusive soft-biometrics such as face and gait analysis [4], [8], gaze estimation from existing CCTV infrastructure enables a wide range of analytics which can be used to optimize customer experience and satisfaction.

In this work we present a series of simple improvements to the gaze estimation pipeline, showcasing state-of-the-art results without constructing a specialized architecture for this task. We combine the recent advancements in computer vi-

sion and image processing fields and construct a model which outperforms previous works on the Gaze360 benchmark.

We propose the CrossGaze architecture which utilizes two separate encoders, one for the face and one for the eye region. This separation allows the face encoder to focus on extracting the global information as the local eye information is obtained by the second backbone. Our model leverages the cross-attention mechanism to combine the extracted features, which enables a prediction focused on the eye region that also accounts for the global information of the whole face.

This work makes the following contributions:

- Our proposed method, CrossGaze, achieves a mean angular error of  $9.94^\circ$  on the Front  $180^\circ$  and  $7.17^\circ$  on the Front Facing subsets of the Gaze360 benchmark, surpassing several state-of-the-art methods. For CrossGaze we chose the most performant components according to our experiments, resulting in a strong baseline for gaze estimation in the wild.
- We provide an ablation study on the face encoder backbone, the procedure for incorporating eye features, and pretraining datasets. Our results demonstrate that gaze estimation can be considerably improved by employing a multi-scale feature extractor, by pretraining on large-scale face datasets, and by incorporating eye-specific information through cross-attention.

## II. RELATED WORK

Recent years have seen a surge of interest in gaze estimation from images [16], [30], [5], [6], [1], driven by advances in computer vision architectures [15], [26], [18], [17]. One of the main challenges in this task is training models that are sufficiently robust to estimate human gaze in unconstrained scenarios. Subsequently, progress in this direction is driven by the development of large and diverse datasets for gaze estimation in the wild [29], [16], that capture subjects in natural settings. Zhang et al. [29] introduced the MPIIGaze dataset, which contains 200K images of 15 subjects in various settings. MPIIGaze has some limitations, such as the release of eye patches, and not full face images, or the small amount of participants which can lead to overfitting.

Kellnhofer et al. [16] tackled some of the limitations of the MPIIGaze dataset by introducing the Gaze360 dataset, which includes a larger number of subjects and more varied scenarios. They also proposed an architecture that incorporates temporal information, using a CNN backbone to extract features from individual consecutive frames, followed by bidirectional LSTM layers to model temporal information, and an MLP that predicts the gaze direction. Zhang et al.

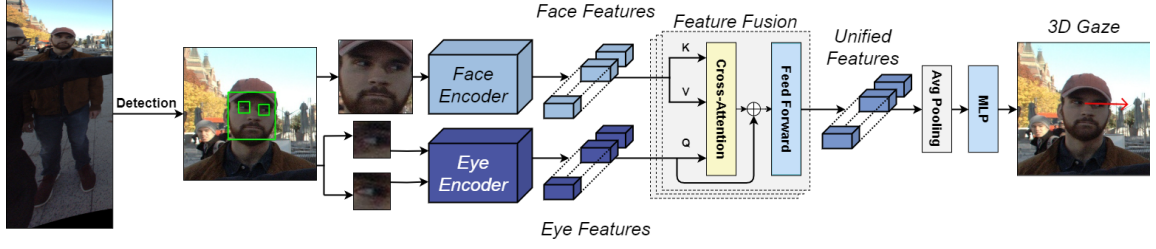


Fig. 1. A high-level overview of the CrossGaze architecture. After the face is detected with a pretrained model, we process the features using a separate encoder for the face and for the eyes. The two resulting feature maps are processed using a cross-attention module to obtain the final gaze prediction.

[30] proposed the FullFace architecture which is a custom CNN with spatial weights that takes the image as input and outputs the 3D gaze estimate. Chen et al. [5] constructed a model based on dilated convolutions that not only takes as input the face image but also the eye crops to help the model focus more on the ocular region for the final prediction.

Cheng et al. [6] proposed CA-Net, which uses a coarse-to-fine approach with a CNN that processes the full face image to estimate a basic gaze direction, and another CNN that processes the eye images in order to make a fine prediction. However, the authors constructed a custom CNN model, which is suboptimal as it cannot easily be modified based on new improvements in deep learning. In contrast, our approach enables straightforward updates to every module in the architecture and can utilize any image feature extractors.

The L2CS-Net [1] architecture utilizes a ResNet-50 [15] backbone on the input image and the extracted gaze features are fed to two branches for predicting the yaw and pitch gaze angles. Each branch includes a fully connected layer that generates continuous as well as discrete predictions. Yan et al. [27] improve on this architecture by modifying the backbone to employ strip pooling which makes the receptive field more suited to gaze estimation. They also incorporate multi-criss-cross attention to capture dependencies between the eye features.

Fang et al. [12] propose a 3D gaze estimation approach which employs a module for head pose detection and eye features extraction to obtain the final prediction. This information is then utilized in conjunction with a depth map to also detect the gaze target. In contrast, our method does not require any additional labels for head pose or depth as it relies solely on the image input to output the 3D gaze.

### III. METHOD

**CrossGaze Architecture.** A high level overview of the proposed architecture is shown in Figure 1. Our model takes as input an in-the-wild image and uses a face detection model to obtain the bounding box predictions of the face and eye regions. The face image is fed into an encoder that captures global gaze features, while the eye images are given as input to the eye encoder which extracts local features. To obtain an eye-informed gaze prediction, the extracted features of the 3 images are combined through a cross-attention module. The resulting combination of features is averaged and processed by a linear layer that outputs the 3D gaze direction.

We consider the face image  $I$  of size  $H \times W \times 3$  to be the input for the Face Encoder  $Enc$  and the images  $I', I''$  of size  $H' \times W' \times 3$  corresponding to the left and right eyes as the input for the Eye Encoder  $Enc'$ . The Face Encoder obtains global gaze features  $F \in \mathbb{R}^{n \times d}$  ( $n$  - number of features,  $d$  - feature dimension) that take into account both the visual information around the eye area and outside of it. On the other hand, the Eye Encoder obtains rich local features  $F' \in \mathbb{R}^{n \times d}$  as it processes images in which the eyes region has a higher resolution than in the face image. The global and local features are merged by a module which employs cross-attention. The fused features are averaged and processed by an MLP which predicts the 3D gaze. This is formulated as:

$$\begin{aligned} F, F' &= Enc(I), Enc'(I', I'') \\ F'' &= CrossAttention(F, F') \\ \hat{g} &= MLP(AvgPool(F'')) \end{aligned} \quad (1)$$

where  $F''$  are the unified features,  $AvgPool$  stands for Average Pooling and  $MLP$  for Multi-Layer Perceptron.  $CrossAttention$  is computed as:

$$CrossAttention = Softmax(Q'K^T/\sqrt{d})V \quad (2)$$

where the Queries ( $Q' \in \mathbb{R}^{n \times d}$ ) are obtained from the sequence of local features  $F'$  and the sequences of Keys ( $K \in \mathbb{R}^{n \times d}$ ) and Values ( $V \in \mathbb{R}^{n \times d}$ ) are obtained from the sequence of global features  $F$ .

**Implementation Details.** While many existing methods [1], [16] use polar coordinates as outputs and labels for training, we found that there is no improvement in comparison to the use of 3D vectors. Consequently, we use normalized 3D vectors  $(x, y, z)$  as both outputs and labels.

For training the CrossGaze model we utilize the cosine loss which represents the complement of the cosine similarity between the predicted gaze vector  $\hat{g}$  and the ground truth  $g$ :

$$L_{Cosine} = 1 - \frac{\hat{g}_i \cdot g_i}{\|\hat{g}_i\| \cdot \|g_i\|} \quad (3)$$

We adapt the RandAugment algorithm [9] as training-time augmentation to increase the robustness of the gaze estimation models. For simplicity and to preserve the gaze information, we remove structural augmentations in the form of rotation-based and shear-based transformations. Additionally, we incorporate into the pool of image transformations the cutout augmentation [11], which removes a small patch of the image. Cutout can act as a form of dropout in cases

where the removed patch corresponds to the eyes region, further regularizing the gaze estimation model.

The models are trained using the AdamW [20] optimizer with a batch size of 256 for 200 epochs. The learning rate employs a step schedule that has a starting value of 0.0001, a step size of 10 epochs and a multiplicative factor of decay equal to 0.95. For training the models a single NVIDIA A100 with 40GB of VRAM was utilised. The training of the CrossGaze model takes approximately 10 hours.

#### IV. EXPERIMENTS AND RESULTS

We chose to conduct our experiments on the Gaze360 dataset as it is a large-scale dataset containing in-the-wild scenarios. Multiple works [16], [14] demonstrated that models trained on Gaze360 generalize well to other gaze benchmarks, obtaining better results than pretraining on different datasets. Consequently, we conduct all our experiments on the Gaze360 dataset which contains 129K training images and 26K testing images. We utilize 2 subsets of Gaze360: Front 180° and Front Facing. Front 180° limits the gaze angles to below 90° while the Front Facing subset limits them to below 20°. This is a similar methodology to other works in gaze estimation [1], [7], [16], [12], as the gaze cannot be inferred from images in which the eyes are not visible. Furthermore, all ablation experiments are conducted on the Front 180° subset which contains approximately 85k training images and 16k testing images.

The most widely used evaluation metric in the field of gaze estimation is the angular error, which is calculated as the angle in degrees between the predicted 3D gaze vector and the ground truth gaze vector. In line with most works on gaze estimation [1], [16], [12] we also utilize it as the main evaluation metric. The angular error is formulated as:

$$\text{AngularError} = \arccos \frac{\hat{g}_i \cdot g_i}{\|\hat{g}_i\| \cdot \|g_i\|} \quad (4)$$

##### A. Evaluation in the Wild

We present a comparative analysis of our architecture with other state-of-the-art models in appearance-based 3D gaze estimation. **FullFace** [30] employs convolutional modules with spatial weights to predict the 2D and 3D gaze from the input image. **Dilated-Net** [5] leverages dilated convolutions to capture minor changes in the gaze information. **RT-Genie** [13] utilizes an ensemble of 4 VGG-16 [23] models on the eye images along with a head pose prediction model to obtain the gaze estimation. **CA-Net** [6] utilizes a CNN on the face image to obtain a coarse prediction and separate encoders for the eye images to refine the prediction. **Gaze360 LSTM** [16] uses a CNN on a window of 7 frames and the extracted features are temporally aggregated with a bidirectional LSTM. **L2CS-Net** [1] uses a ResNet-50 as backbone and employs 2 different prediction heads for the yaw and pitch angles. **SPMCCA-Net** [27] integrates strip pooling and criss-cross attention to the ResNet backbone. **DAM** [12] predicts the head pose and extracts the features of the eye crops to make the gaze prediction, however it benefits from additional annotations.

Table I displays the comparison between the proposed CrossGaze model and the other 3D gaze estimation architectures on the Front 180° and Front Facing subsets of the Gaze360 benchmark. For our model, we show the average and standard deviation of the mean angular error for 3 different runs. The CrossGaze architecture is composed of an Inception ResNet [25] face encoder, a ResNet-18 [15] eye encoder, and a cross-attention module that combines the extracted features of both models. The pretrained version of CrossGaze on VggFace2 [3] achieves state-of-the-art mean angular error on the Front Facing testing subset, demonstrating its capability for gaze estimation in the wild. On Front 180°, CrossGaze is competitive with DAM [12] which utilizes additional head pose annotations during training. The randomly initialized version of our architecture manages to be on par with the pretrained versions of other models, showcasing its capability to generalize with less data.

TABLE I  
COMPARISON OF GAZE ESTIMATION MODELS ON SUBSETS OF THE GAZE360 DATASET. MODELS HIGHLIGHTED WITH \* EMPLOY INITIALIZATION FROM PRETRAINED WEIGHTS WHILE THOSE WITH \*\* ALSO UTILIZE HEAD POSE ANNOTATIONS. TABLE ADAPTED FROM [1].

Model	Front 180°	Front Facing
FullFace [30]	14.99°	N/A
Dilated-Net [5]	13.73°	N/A
RT-Genie (4 ensemble) [13]	12.26°	N/A
CA-Net [6]	11.20°	N/A
Gaze360 LSTM* [16]	11.04°	N/A
L2CS-Net ( $\beta = 1$ )* [1]	10.41°	9.02°
SPMCCA-Net ( $\beta = 2$ )* [27]	10.13°	8.40°
DAM** [12]	9.6°	9.2°
<b>CrossGaze (Random Init)</b>	10.65° ± 0.03	7.84° ± 0.27
<b>CrossGaze (Pretrained)</b>	<b>9.94° ± 0.06</b>	<b>7.17° ± 0.04</b>

##### B. Ablation Studies

In this section, we present ablation studies for each critical component of CrossGaze: the face encoder backbone, pretraining dataset, and methods for incorporating eye information in the computation. Our CrossGaze architecture is constructed from scratch in a progressive manner, starting from the face encoder backbone and incrementally adding improvements such as the pretraining dataset and the eye feature integration. For all comparisons we show the mean and standard deviation of the performance of each setting computed for 3 different runs.

**Multi-scale face features aids gaze estimation.** To obtain a strong gaze estimation architecture we experiment with multiple image processing backbones. These models extract the features of the face image which are then passed to a gaze estimation head consisting of an MLP that outputs the 3D vector. We experiment with CNN backbones such as EfficientNet [26], ConvNeXt [18], Inception ResNet [25] and an attention-based backbone (i.e. Swin Transformer [17]).

TABLE II

PERFORMANCE OF EACH ARCHITECTURE ON THE FRONT 180° SUBSET.  
THE INCEPTION RESNET OBTAINS THE LOWEST ERROR,  
DEMONSTRATING ITS SUITABILITY FOR GAZE ESTIMATION IN THE WILD.

Model Init.	Model	Mean Angular Error
Random Init	EfficientNet	$14.1^\circ \pm 0.32$
	Swin Transformer	$12.03^\circ \pm 0.3$
	ConvNeXt	$12.79^\circ \pm 0.05$
	Inception ResNet	<b><math>10.91^\circ \pm 0.09</math></b>
Pretrained ImageNet	EfficientNet	$12.03^\circ \pm 0.3$
	Swin Transformer	$10.87^\circ \pm 0.05$
	ConvNeXt	$11.11^\circ \pm 0.08$
	Inception ResNet	<b><math>10.82^\circ \pm 0.08</math></b>

The results of this experiment are shown in Table II. In both initialization scenarios, the Inception ResNet architecture obtains the lowest mean angular error with a value of  $10.91^\circ$  for random initialization and  $10.82^\circ$  for ImageNet pretrained weights. These results demonstrate the suitability of the Inception ResNet architecture for processing the face: the capability to process the input at multiple scales in every layer helps the model focus both on the eyes region and on the periocular area for the gaze prediction.

**Face pretraining enhances gaze feature extraction.** We analyze the impact of different pretraining datasets on the task of gaze estimation. For this, we only use the top performing face model from the previous experiments, an Inception ResNet, pretrained on the following datasets: ImageNet [10], Casia-WebFace [28], and VGGFace2 [3]. We initialize the Inception ResNet with the pretrained weights of each dataset and train the entire network for gaze estimation.

Table III shows the results of the experiment, which indicate that the models pretrained on face images obtain better results compared to those pretrained on ImageNet. Face datasets have a closer distribution to the data in the Gaze360 dataset, which naturally translates in improved performance. Additionally, the model pretrained on the larger dataset (VGGFace2) performs better than the model pretrained on Casia-WebFace due to the increased amount of training data. While Casia-WebFace contains approximately 500K images, VGGFace2 contains 3.3M images, resulting in an improvement of  $0.24^\circ$  in mean angular error.

TABLE III

RESULTS OF PRETRAINING AN INCEPTION RESNET ON DIFFERENT DATASETS AND TRANSFER LEARNING TO GAZE ESTIMATION.

Pretraining Dataset	Mean Angular Error
ImageNet [10]	$10.82^\circ \pm 0.08$
CASIA-WebFace [28]	$10.26^\circ \pm 0.05$
VggFace2 [3]	<b><math>10.02^\circ \pm 0.07</math></b>

**Integrating eye features enriches the gaze information.** To enable the model to focus on the eye region for the gaze prediction task we insert an additional encoder that

separately takes as input both eye images. The extracted face and eye features are combined to obtain the 3D gaze estimation. We experiment with 2 different fusing strategies: the first consists of combining the face and eye features through a fully connected network while the second involves aggregating the information with a cross-attention module.

Table IV displays the results of the eye features integration for an Inception ResNet face encoder and a ResNet-18 eye encoder. We utilize the ResNet-18 as the secondary encoder because the eye images are of a lower resolution (64x64 after rescaling). In the case of both random initialization and pretrained weights initialization, the use of additional eye features brings an improvement in gaze estimation performance. Furthermore, the cross-attention mechanism manages to better aggregate the information, as it obtains a smaller error in both scenarios compared to the fully connected network. These results motivated the design of the CrossGaze architecture, shown in Figure 1, to also employ the secondary eye encoder and the cross-attention module.

TABLE IV

RESULTS FOR COMBINING FEATURES. CROSS-ATTENDING TO EYE FEATURES IMPROVES GAZE ESTIMATION PERFORMANCE.

Model Init.	Eye features combination	Mean Angular Err.
Random Init	No eye features	$10.91^\circ \pm 0.09$
	FCN	$10.75^\circ \pm 0.02$
	Cross-Attention	<b><math>10.65^\circ \pm 0.03</math></b>
Pretrained VggFace2	No eye features	$10.02^\circ \pm 0.07$
	FCN	$10.01^\circ \pm 0.05$
	Cross-Attention	<b><math>9.94^\circ \pm 0.06</math></b>

## V. CONCLUSION

This work proposes CrossGaze, a CNN-based architecture that leverages cross-attention, designed for the task of 3D gaze estimation from images in the wild. On the Gaze360 benchmark, our architecture outperforms several state-of-the-art methods. The model processes both the full face image as well as the eye images to make an eye-informed prediction. To combine the global features of the face with the local images of the eyes, we employ a cross-attention module that captures the most relevant characteristics of the gaze.

We conduct an ablation study, starting from scratch and incrementally adding improvements to the architecture. We start by choosing a suitable face backbone with multi-scale processing at every layer, which helps in extracting relevant gaze features. For this backbone, we show that pretraining on large scale datasets of face images improves generalization, as opposed to other general pretraining datasets. Our results demonstrate that incorporating a separate eye encoder and combining global and local gaze features through cross-attention further improves the performance.

Our work represents a step forward towards gaze estimation in realistic environments and has the potential to enable practical applications including human-computer interaction, driver assistance systems, and assistive technology for individuals with disabilities.

## REFERENCES

- [1] A. A. Abdelrahman, T. Hempel, A. Khalifa, and A. Al-Hamadi. L2cs-net: Fine-grained gaze estimation in unconstrained environments. *arXiv preprint arXiv:2203.03339*, 2022.
- [2] C. Bermejo, D. Chatzopoulos, and P. Hui. Eyeshopper: Estimating shoppers' gaze using cctv cameras. In *Proceedings of the 28th ACM International Conference on Multimedia*, MM '20, page 2765–2774, New York, NY, USA, 2020. Association for Computing Machinery.
- [3] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman. Vggface2: A dataset for recognising faces across pose and age. In *2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018)*, pages 67–74. IEEE, 2018.
- [4] A. Catruna, A. Cosma, and I. E. Radoi. From face to gait: Weakly-supervised learning of gender information from walking patterns. In *2021 16th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2021)*, pages 1–5. IEEE, 2021.
- [5] Z. Chen and B. E. Shi. Appearance-based gaze estimation using dilated-convolutions. In *Asian Conference on Computer Vision*, pages 309–324. Springer, 2018.
- [6] 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*, volume 34, pages 10623–10630, 2020.
- [7] Y. Cheng, H. Wang, Y. Bao, and F. Lu. Appearance-based gaze estimation with deep learning: A review and benchmark. *arXiv preprint arXiv:2104.12668*, 2021.
- [8] A. Cosma and E. Radoi. Learning gait representations with noisy multi-task learning. *Sensors*, 22(18), 2022.
- [9] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 702–703, 2020.
- [10] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [11] T. DeVries and G. W. Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017.
- [12] 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11390–11399, 2021.
- [13] T. Fischer, H. J. Chang, and Y. Demiris. Rt-gene: Real-time eye gaze estimation in natural environments. In *Proceedings of the European conference on computer vision (ECCV)*, pages 334–352, 2018.
- [14] S. Ghosh, A. Dhall, M. Hayat, J. Knibbe, and Q. Ji. Automatic gaze analysis: A survey of deep learning based approaches. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(1):61–84, 2023.
- [15] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [16] P. Kellnhofer, A. Recasens, S. Stent, W. Matusik, and A. Torralba. Gaze360: Physically unconstrained gaze estimation in the wild. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6912–6921, 2019.
- [17] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021.
- [18] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11976–11986, 2022.
- [19] M. Lombardi, E. Maiettini, D. De Tommaso, A. Wykowska, and L. Natale. Toward an attentive robotic architecture: Learning-based mutual gaze estimation in human–robot interaction. *Frontiers in Robotics and AI*, 9, 2022.
- [20] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [21] P. Pathirana, S. Senarath, D. Meedeniya, and S. Jayarathna. Eye gaze estimation: A survey on deep learning-based approaches. *Expert Systems with Applications*, 199:1–16, 03 2022.
- [22] S. M. Shah, Z. Sun, K. Zaman, A. Hussain, M. Shoaib, and L. Pei. A driver gaze estimation method based on deep learning. *Sensors*, 22(10), 2022.
- [23] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [24] A. M. Soccini. Gaze estimation based on head movements in virtual reality applications using deep learning. In *2017 IEEE Virtual Reality (VR)*, pages 413–414, 2017.
- [25] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31, 2017.
- [26] M. Tan and Q. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- [27] C. Yan, W. Pan, C. Xu, S. Dai, and X. Li. Gaze estimation via strip pooling and multi-criss-cross attention networks. *Applied Sciences*, 13(10):5901, 2023.
- [28] D. Yi, Z. Lei, S. Liao, and S. Z. Li. Learning face representation from scratch. *arXiv preprint arXiv:1411.7923*, 2014.
- [29] X. Zhang, Y. Sugano, M. Fritz, and A. Bulling. Appearance-based gaze estimation in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4511–4520, 2015.
- [30] X. Zhang, Y. Sugano, M. Fritz, and A. Bulling. It's written all over your face: Full-face appearance-based gaze estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 51–60, 2017.