# UniGaze: Towards Universal Gaze Estimation via Large-scale Pre-Training

Jiawei Qin[1], Xucong Zhang[2], Yusuke Sugano[1]

[1] Institute of Industrial Science, The University of Tokyo, Komaba 4-6-1, Tokyo, Japan

[2] Computer Vision Lab, Delft University of Technology, Mekelweg 5, Delft, Netherlands

{jqin, sugano}@iis.u-tokyo.ac.jp

xucong.zhang@tudelft.nl

Figure 1. Leveraging self-supervised pre-training on large-scale facial data, the proposed UniGaze demonstrates strong generalization to unseen in-the-wild face images under diverse conditions, including facial appearance, lighting conditions, variant head poses, and face resolutions. We draw the estimated gaze direction with green arrows. More examples are shown in the supplementary materials.

## Abstract

*Despite decades of research on data collection and model architectures, current gaze estimation models encounter significant challenges in generalizing across diverse data domains. Recent advances in self-supervised pre-training have shown remarkable performances in generalization across various vision tasks. However, their effectiveness in gaze estimation remains unexplored. We propose UniGaze, for the first time, leveraging large-scale in-the-wild facial datasets for gaze estimation through self-supervised pre-training. Through systematic investigation, we clarify critical factors that are essential for effective pre-training in gaze estimation. Our experiments reveal that self-supervised approaches designed for semantic tasks fail when applied to gaze estimation, while our carefully designed pre-training pipeline consistently improves cross-domain performance. Through comprehensive experiments of challenging cross-dataset evaluation and novel proto-cols including leave-one-dataset-out and joint-dataset settings, we demonstrate that UniGaze significantly improves generalization across multiple data domains while minimizing reliance on costly labeled data. Source code and model are available at* https://github.com/ut-vision/UniGaze.

## 1. Introduction

Gaze estimation is a key task in computer vision with broad applications in human-computer interaction [5, 59], virtual reality [58], and behavioral analysis [26, 38, 65]. The methodology in unconstrained environments has been actively researched in past decades [2, 14, 24, 56, 57, 71, 75, 88]. However, achieving robust and accurate gaze estimation in unseen test environments remains a fundamental challenge. Although current models can achieve high accuracy when tailored to specific datasets or individ-

uals [60], they consistently show significant performance degradation in novel environments. This challenge is evidenced by the persistent difficulty in training one single model to achieve high accuracy across various gaze estimation datasets [22, 39, 45, 91, 92]. Despite decades of research on data collection and model architecture, the problem of generalization across variations in head pose, identity, and lighting conditions remains largely unsolved.

Recent advances in self-supervised pre-training have shown remarkable potential for improving model generalization to arbitrary data domains. Following scaling laws [36], larger models trained with extensive data and compute resources can be more sample-efficient, which has been observed with significant performance improvements across numerous tasks such as image classification [19, 67, 69], segmentation [43, 48, 80, 94], human-centric tasks [40], gaze following [50], and face-focused applications [8, 23, 70, 77, 96]. However, their effectiveness in gaze estimation remains unexplored. Unlike these typical tasks, gaze estimation is a geometric regression task to map facial appearance to precise 3D vectors. Our experiments show that existing pre-trained models [10, 11, 96] fail to improve gaze estimation, as they are optimized for semantic understanding rather than the fine-grained facial structure crucial for gaze direction. This raises an essential question: Can large-scale pre-training benefit the fine-grained geometric nature of gaze estimation? Our research offers an answer: Yes, but *only when the pre-training is specifically tailored to the unique constraints of gaze estimation*.

In this work, we present a novel approach toward **Uni**versal **Gaze** estimation, dubbed UniGaze, exploring the potential of large-scale self-supervised pre-training for appearance-based gaze estimation. We employ Masked Autoencoder (MAE) [29] on the Vision Transformer (ViT) architecture [17], using diverse in-the-wild face image datasets. Through systematic experimentation, we discovered that pre-training for gaze estimation requires three essential components that differ from typical pre-training approaches: (1) pre-training on normalized facial images maintaining the spatial configuration required by downstream gaze models [90], (2) ensuring diverse yet balanced head pose distributions to learn robust facial representations across viewing angles, and (3) incorporating sufficient identity diversity to generalize across different facial appearances. This strategy allows the model to learn appropriate feature representations within the specific input space required by gaze estimation models, enabling effective transfer to the downstream gaze estimation task.

Through extensive experiments, we demonstrate that training our pre-trained UniGaze on gaze-specific datasets yields substantial generalization performance improvements across multiple data domains [22, 39, 45, 91, 92], sur-passing state-of-the-art domain generalization methods [3, 13, 83, 85, 86, 93]. Pre-training with general semantics [10, 11, 98], face-specific semantics [96], and vanilla MAE trained on small-scale gaze data [33] all fail to transfer effectively to gaze estimation, sometimes even performing worse than simple CNNs. In contrast, our carefully designed pre-training approach consistently improves accuracy across diverse environments. This highlights our crucial discovery that suitable datasets, normalized facial images, and varied yet balanced head pose distributions are vital for developing transferable representations for gaze estimation, offering practical guidelines for future research in this field.

Additionally, in gaze estimation, the widely used cross-dataset evaluation usually only trains the model on a single dataset, which cannot fully represent the appearance diversity and pose range in real-world environments. There also remains the possibility that some adaptation methods merely overfit some specific training/testing dataset combinations. To address this limitation and reflect the practical requirements of real-world applications, we propose two novel evaluation protocols utilizing multiple datasets for training: a *leave-one-dataset-out* setting that assesses generalization to unseen datasets and a *joint-dataset* setting that evaluates the achievable performance across multiple datasets. Our comprehensive evaluation demonstrates that UniGaze consistently achieves superior performance across diverse environments under these protocols, suggesting the effectiveness of large-scale pre-training.

In summary, our contributions are threefold: (i) We present UniGaze, a novel gaze estimation model that addresses the fundamental challenge of cross-domain generalization in appearance-based gaze estimation via large-scale pre-training. (ii) We provide empirical evidence that MAE pre-training on normalized face images can learn meaningful representations for gaze estimation tasks. (iii) We propose leave-one-dataset-out and joint-dataset evaluation protocols that offer practical benchmarks for assessing gaze estimation performance in real-world scenarios.

## 2. Related Works

### 2.1. Appearance-based Gaze Estimation

Appearance-based methods for gaze estimation have gained prominence, leveraging the ability of deep learning to learn gaze representations from gaze-labeled datasets [20, 22, 39, 45, 88, 91, 92]. However, these approaches are limited by the scarcity of comprehensive, well-labeled gaze datasets, which hinders performance in unconstrained settings. One solution to data scarcity has been synthetic data generation, where gaze images are synthesized with controllable variables such as lighting, head pose, and redirected gaze [34, 61, 64, 74, 84, 95]. Despite its utility, syn-

thetic data often lacks realism, leading to domain adaptation challenges [41, 78]. Furthermore, subtle inaccuracies in gaze labels can compromise model performance when transferred to real-world scenarios.

In addition to data-driven approaches, method-driven solutions have been explored to improve robustness and generalizability [3, 13, 53, 83, 86, 93]. For instance, gaze frontalization [82], multi-view consistency [4, 30], and contrastive learning [28, 35] have been used to learn generalized gaze representations. Clip-Gaze [86] and LG-Gaze [85] use the linguistic features extracted from the vision-language model to regularize the gaze feature learning. Alternatively, Kothari *et al*. [44] utilizes in-the-wild face datasets with *look-at-each-other* labels as weak supervision. Furthermore, domain adaptation methods [60, 76] incorporate target domain samples with limited or no labels. For example, PnP-GA+ uses the plug-and-play method to adapt the gaze estimation model to new domains with assembling model variants [52].

Network architectures in gaze estimation remain predominantly CNN-based, particularly relying on ResNet [27]. Recent work by Cheng *et al*. [12] explored ViTs [17] for gaze estimation, finding that although ViTs hold promise, they require extensive pre-training data beyond standard ImageNet [16] to perform effectively. This motivates the exploration of ViT models specifically tailored to learn diverse gaze representations through pre-training approaches.

## 2.2. Large-scale Pre-Training in Vision Models

Large-scale pre-training has become fundamental for foundational model development in computer vision. Recent studies show that pre-trained generative models enhance representation learning across various applications. For example, diffusion models [31, 63] have been successfully applied to image classification [47] and segmentation [48, 80, 94]. Another effective unsupervised pre-training approach is masked autoencoding [19, 29, 67, 69], which demonstrates strong capabilities in image recognition and robust feature extraction. It is also shown that masked image modeling benefits more for ViT than CNN family [18, 46].

MAE pre-training has been adapted to several domains such as audio [32], video [72], and microscopy [46]. Recently, Sapiens [40] leverages human-centric data for MAE pre-training, resulting in strong generalization across multiple human-related tasks. These adaptations to specific domains have made a crucial step in advancing research. Similarly, multiple self-supervised pre-training works have been proposed to achieve notable improvement in face-centric tasks, such as expression recognition, facial attribute recognition, and face alignment [8, 23, 70, 77, 96]. These studies underscore the effectiveness of pre-training masked autoen-

| Dataset | Type | # Identities | # Samples |
|---------|------|--------------|-----------|
| CelebV-Text [87] | Real | 13,179 | 666,967 |
| VFHQ [81] | Real | 10,382 | 231,809 |
| VGGFace2 [9] | Real | 9,131 | 182,603 |
| FaceSynthetics [79] | Syn. | 86,878† | 86,878 |
| SFHQ-T2I [6] | Syn. | 120,241† | 120,241 |
| FFHQ-NV [37, 61] | Syn. | 25,000 | 100,000 |
| XGaze-Dense [55, 62, 92] | Syn. | 60 | 267,160 |
| Total | - | 264,871 | 1,655,668 |

Table 1. Statistics of face datasets used to pre-train UniGaze in terms of data type to be real or synthetic (Syn.), number of identities, and number of samples. The † indicates that we assume there are no duplicated identities during the synthesis image generation.

coders on large, diverse datasets for enhancing model robustness and generalizability. In contrast, gaze estimation has received less attention in this context. A concurrent work [33] even highlighted a limitation in applying MAE to gaze estimation with ViT, noting that random masking tends to focus on global semantics while neglecting critical gaze-related information.

## 3. UniGaze

Our UniGaze model consists of a large-scale pre-training stage followed by task-specific training. The pre-training stage utilizes a carefully curated dataset of $1.6\,\mathrm{M}$ facial images that combines real-world and synthetic data to ensure broad coverage of head poses and facial appearances. We adopt the MAE framework to learn robust facial representations from this diverse dataset and then train it with labeled gaze data for the downstream gaze estimation task.

### 3.1. Pre-Training Datasets

Learning robust gaze representations requires extensive head pose and facial appearance variations. To achieve this, we combine two complementary data sources: real data that captures natural face appearances and synthetic data that allows us to cover extreme pose variations and diverse appearances systematically.

Table 1 summarizes our pre-training dataset composition, which comprises approximately $1.6\,\mathrm{M}$ samples. This combined dataset includes over $260\,\mathrm{k}$ unique identities, vastly exceeding the 1,474 subjects in the existing Gaze-Capture dataset [45]. Although GazeCapture [45] and ETH-XGaze [92] offer over $2.4\,\mathrm{M}$ and $1\,\mathrm{M}$ samples, respectively, our combined dataset provides a significantly higher level of diversity. In this manner, our pre-training data offers a significantly broader representation regarding facial appearances, head poses, and environmental conditions, surpassing existing gaze estimation datasets.

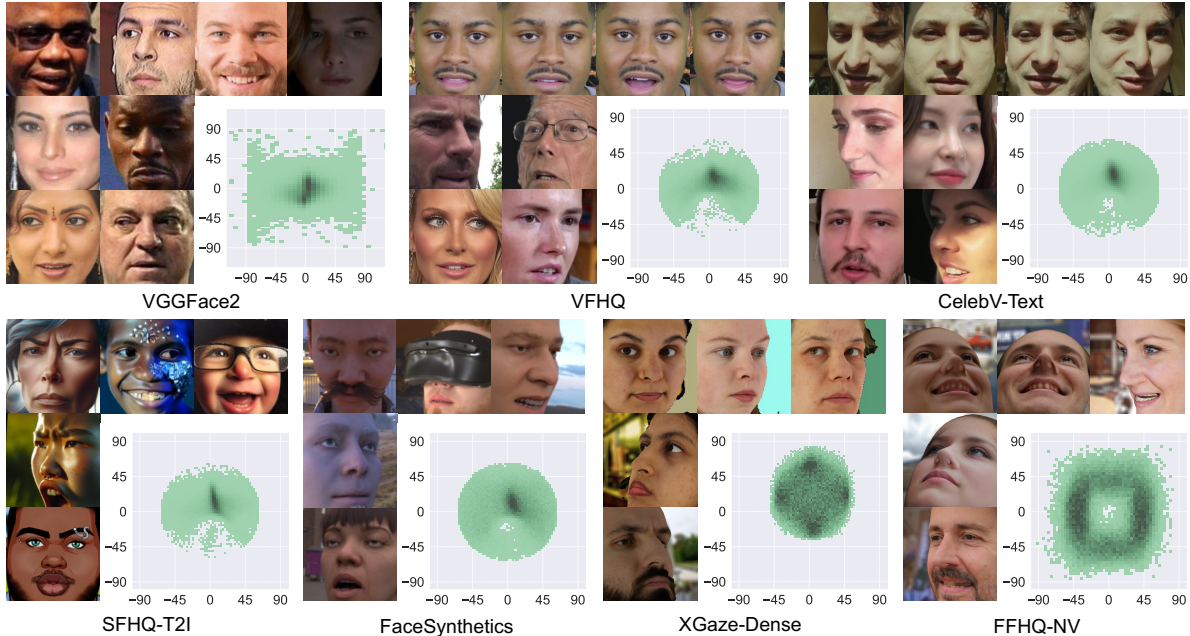**Real Datasets.** We use VGGFace2 [9] for its large number

Figure 2. Example of the normalized facial images from different datasets in the pre-training stage. We also draw their head pose distributions where the vertical axis is the pitch rotation angle and the horizontal axis is the yaw rotation angle in degrees.

of identities with diverse conditions. Additionally, we incorporate two high-quality video datasets, VFHQ [81] and CelebV-Text [87], leveraging their natural pose and gaze variations provided by video sequences. To reduce redundancy while maintaining diversity, we sub-sample every 15 frames from VFHQ, 45 frames from CelebV-Text, and select 20 images per identity from VGGFace2.

**Synthetic Datasets.** To ensure diverse facial appearances, we use two synthetic datasets: FaceSynthetics [79], generated via computer graphics, and SFHQ-T2I [6], created with diffusion models. We further employ novel-view synthesis techniques to extend the range in head poses. We reconstruct 3D facial shapes from FFHQ [37] via the single-view method [61], synthesizing FFHQ-NV. Following [62], we use Metashape [55] for multi-view 3D face reconstruction and apply novel-view rendering to get XGaze-Dense. Since XGaze-Dense is used without gaze labels, its role is equivalent to that of a generic facial dataset in our pre-training.

**Data Pre-processing** We perform facial landmark detection [7] to estimate the 3D head pose by perspective-n-point (PnP) algorithm [21]. We then apply data normalization [90] to crop face images, ensuring alignment of the input space between MAE pre-training and gaze estimation. We filter out samples with extreme head poses for VFHQ, CelebV-Text, SFHQ-T2I, and FaceSynthetics to eliminate extreme cases where the face is invisible. Precisely, we discard samples with an $L_2$ norm of pitch and yaw angles exceeding 80 degrees. Figure 2 shows example face images

and head pose distributions for each dataset.

## 3.2. Training Procedure

We follow the MAE [29] to pre-train the ViT model without any gaze labels. Briefly, it randomly masks patches of the input image and the model is trained to predict the masked contents. Given an input image $X$, it is divided into $N$ patches $\{x_i\}_{i=1}^N$. A subset of these patches is masked out, denoted by $\{x_i\}_{i \in M}$ with $M \subset \{1, \ldots, N\}$. The encoder processes the visible patches $\{x_i\}_{i \in V}$, where $V = \{1, \ldots, N\} \setminus M$, and generates a latent representation $z$. An extra decoder then takes $z$ and reconstructs the masked patches $\{\hat{x}_i\}_{i \in M}$ with the loss $\mathcal{L}_{\text{MAE}} = \frac{1}{|M|} \sum_{i \in M} \|x_i - \hat{x}_i\|^2$. This process encourages the model to capture meaningful structures and features in $z$, making it well-suited for downstream tasks [40]. During pre-training, we randomly apply flip, central crop, and color jitter, and the mask ratio is 75%. The loss is computed with the pixel value normalized within each patch, which is suggested to have better representation [29].

**Gaze Estimation Training** With labeled gaze datasets, we train the pre-trained model for the gaze estimation downstream task. We replace the decoder with a fully connected layer to predict gaze direction from the latent representation $z$. The gaze direction is represented by a 2D vector in the polar angle coordinate system. We use the loss $\mathcal{L}_1 = \|g - \hat{g}\|_1$, where $g$ and $\hat{g}$ denote the ground-truth gaze label and the prediction, respectively.

| Models \ Test | M | GC | E | G360 |
|---|---|---|---|---|
| ResNet-50 | 6.75 | 10.08 | 10.28 | 19.80 |
| GazeTR-50 | 7.09 | 10.95 | 9.47 | 21.10 |
| DINO-B [10] | 7.73 | 9.12 | 11.98 | 19.44 |
| MoCo-v3-B [11] | 7.19 | 8.88 | 9.50 | 19.09 |
| FaRL-B [96] | 7.18 | 8.56 | 9.53 | 17.28 |
| UniGaze-B | 6.21 | 7.35 | 6.64 | 12.18 |
| ViT-H | 7.68 | 9.36 | 10.40 | 20.38 |
| UniGaze-H | **5.57** | **6.56** | **6.53** | **11.19** |

(a) Results when trained on XGaze.

| Models \ Test | $X_{Test}$ | GC | E | G360 |
|---|---|---|---|---|
| ResNet-50 | 32.80 | 6.75 | 17.11 | 27.92 |
| GazeTR-50 | 29.00 | 7.06 | 19.38 | 28.22 |
| DINO-B [10] | 33.41 | 8.13 | 20.07 | 30.39 |
| MoCo-v3-B [11] | 30.99 | 6.53 | 17.56 | 27.14 |
| FaRL-B [96] | 30.43 | 6.14 | 16.55 | 25.97 |
| UniGaze-B | 30.56 | 5.68 | 17.54 | **20.85** |
| ViT-H | **28.25** | 6.87 | 16.02 | 25.30 |
| UniGaze-H | 33.11 | **4.87** | **12.66** | 21.28 |

(b) Results when trained on MPIIFaceGaze.

| Models \ Test | $X_{Test}$ | M | E | G360 |
|---|---|---|---|---|
| ResNet-50 | 26.56 | 5.84 | 13.39 | 25.33 |
| GazeTR-50 | 23.57 | 5.49 | 14.25 | 25.48 |
| DINO-B [10] | 27.56 | 6.97 | 16.75 | 26.49 |
| MoCo-v3-B [11] | 27.26 | 5.29 | 14.75 | 25.94 |
| FaRL-B [96] | 26.55 | 5.74 | 15.18 | 23.49 |
| UniGaze-B | **20.63** | 5.01 | 10.91 | 18.91 |
| ViT-H | 23.49 | 5.54 | 14.60 | 23.98 |
| UniGaze-H | 22.67 | **4.89** | **10.61** | **17.77** |

(c) Results when trained on GazeCapture.

| Models \ Test | $X_{Test}$ | M | GC | G360 |
|---|---|---|---|---|
| ResNet-50 | 37.50 | 16.88 | 16.37 | 30.03 |
| GazeTR-50 | 35.82 | 16.37 | 16.63 | 25.69 |
| DINO-B [10] | 38.21 | 17.84 | 18.45 | 33.49 |
| MoCo-v3-B [11] | 31.98 | 14.29 | 14.15 | 24.90 |
| FaRL-B [96] | 34.74 | 12.96 | 12.31 | 26.88 |
| UniGaze-B | 25.37 | **8.91** | **8.89** | 18.27 |
| ViT-H | 30.90 | 15.51 | 13.51 | 26.12 |
| UniGaze-H | **23.79** | 8.93 | 9.97 | **16.00** |

(d) Results when trained on EYEDIAP.

| Models \ Test | $X_{Test}$ | M | GC | E |
|---|---|---|---|---|
| ResNet-50 | 18.83 | 10.25 | 9.90 | 12.06 |
| GazeTR-50 | 18.04 | 11.28 | 10.82 | 13.69 |
| DINO-B [10] | 22.25 | 11.02 | 11.17 | 17.53 |
| MoCo-v3-B [11] | 17.59 | 7.50 | 8.72 | 11.91 |
| FaRL-B [96] | 20.55 | 7.62 | 7.75 | 12.43 |
| UniGaze-B | **12.22** | 6.00 | 8.11 | 8.74 |
| ViT-H | 18.63 | 8.24 | 9.43 | 11.50 |
| UniGaze-H | 16.39 | **5.43** | **6.48** | **6.97** |

(e) Results when trained on Gaze360.

Table 2. Cross-dataset evaluation of different models trained on one dataset and tested on multiple unseen datasets. Each subtable corresponds to a specific training dataset, with columns representing the testing datasets ($X_{Test}$: XGaze Test, **M**: MPIIFaceGaze, **GC**: GazeCapture, **E**: EYEDIAP, **G360**: Gaze360). Results demonstrate the generalization ability of each model, with UniGaze consistently outperforming other baselines in most settings, showcasing the effectiveness of our pre-training approach.

## 3.3. Implementation Details

We use the Adam optimizer [42] with a base learning rate of $1.5 \times 10^{-4}$ and a weight decay of $0.05$. We set a batch size of 4,096 and pre-train the ViT-Huge model for 300 epochs, which requires approximately 120 hours on four NVIDIA H100 GPUs. During gaze estimation training, we do not apply any image augmentation. We use the Adam optimizer [42] with a learning rate of $1 \times 10^{-4}$ and a weight decay of $1 \times 10^{-6}$, and the one-cycle learning rate [68]. More details and examples can be found in supplementary materials.

## 4. Experiments

### 4.1. Experimental Settings

**Gaze Datasets** We conduct experiments on multiple gaze estimation datasets following the common way of utilization in recent works. **MPIIFaceGaze** [89] contains 15 subjects with nearly frontal head poses. In experiments requiring splitting, the first ten subjects are used for training and the remaining five for testing. **ETH-XGaze** [92] comprises over $750\,\text{k}$ publicly available gaze-labeled images of 80 subjects. We refer to the "XGaze" as the 80 subjects, while "XGaze Train/Test" indicates 60/20-subject split. When training, we randomly select three out of the 18 cameras to reduce redundancy without losing effectiveness, and we

utilize all cameras for testing. **EYEDIAP** [22] includes 16 subjects and two sessions with screen (CS) and 3D floating object (FT) targets. We use both sessions and split the data into training and test sets by subjects, with an 8/8 split. The data is pre-processed using the pipeline by Park *et al.* [60]. **Gaze360** [39] consists of indoor and outdoor images of 238 subjects with wide ranges of head poses and gaze directions. We use the training and test split defined in the original paper. **GazeCapture** [45] contains around 1400 subjects collected through crowd-sourcing. We use the training and test split defined in the original GazeCapture paper and pre-process with the pipeline by Park *et al.* [60]. When training, we sample every 15 frames to reduce redundancy, and we use all samples for testing.

**Baseline Architectures** We draw upon existing gaze estimation research and consider baselines of convolutional neural networks, ViTs, and hybrid models. **ResNet** [27] models are lightweight yet powerful CNNs, which dominate the backbone and baseline in most of the gaze estimation works [13, 83, 92, 93]. We use ResNet-50 in our experiments. **GazeTR-50 (Hybrid)** [12] is a hybrid network where the image features extracted from the ResNet-50 are fed into the transformer. We include **DINO-B** [10], **MoCo-v3-B** [11], and **FaRL-B** [96] as representative pre-trained ViT-Base models: DINO and MoCo-v3 have rich

| Models | X→M | X→E$_{CS}$ | G360→M | G360→E$_{CS}$ |
|---|---|---|---|---|
| ResNet-18[†] [93] | 8.02 | 9.11 | 8.04 | 9.20 |
| PureGaze[†] [13] | 7.08 | 7.48 | 9.28 | 9.32 |
| Gaze-Consistent[†] [83] | 6.50 | 7.44 | 7.55 | 9.03 |
| AGG[†] [3] | <u>5.91</u> | 6.75 | 7.87 | 7.93 |
| CLIP-Gaze[†] [86] | 6.41 | 7.51 | 6.89 | 7.06 |
| LG-Gaze[†] [85] | 6.45 | 7.22 | 6.83 | 6.86 |
| Gaze-BAR[†] [93] | 6.35 | 6.72 | 6.96 | 8.79 |
| ViT-H | 7.68 | 8.58 | 8.24 | 7.79 |
| UniGaze-B | 6.21 | <u>5.08</u> | <u>6.00</u> | <u>6.63</u> |
| UniGaze-H | **5.57** | **4.65** | **5.43** | **5.35** |

Table 3. Domain generalization compared with SOTA methods in the cross-dataset setting. Since most of the methods in the table are based on ResNet-18, we also show the ViT-H baseline that only pre-trained on ImageNet for fair comparison with UniGaze-H.

| Methods \ Dataset | X | M | GC | E | G360 |
|---|---|---|---|---|---|
| Abdelrahman *et al.* [1] | - | 3.92 | - | - | - |
| Guan *et al.* [25] | - | - | - | - | 9.81 |
| Shi *et al.* [66] | - | **3.61** | - | 4.78 | - |
| 3DGazeNet [73] | 4.2 | 4.0 | 3.1 | - | 9.6 |
| UniGaze-H | **3.96** | 4.07 | **3.01** | **4.34** | **9.44** |

Table 4. Within-dataset evaluation compared with SOTA methods. Note we here follow the train-test split from individual previous works.

| Models \ Test | X$_{Test}$ | M | GC | E | G360 |
|---|---|---|---|---|---|
| ResNet-50 | 16.31 | 4.94 | 6.71 | 7.89 | 18.69 |
| GazeTR-50 | 15.47 | 4.94 | 7.22 | 8.26 | 19.75 |
| ViT-L | 19.32 | 5.09 | 6.33 | 8.98 | 22.68 |
| FaRL-B | 19.08 | 5.09 | 6.08 | 8.18 | 18.28 |
| UniGaze-B | 11.78 | 4.73 | 5.86 | 6.31 | 12.41 |
| UniGaze-L | **10.93** | 4.64 | 5.79 | 6.56 | 12.44 |
| UniGaze-H | 11.29 | **4.51** | **5.47** | **5.88** | **12.37** |

Table 5. Results of leave-one-dataset-out evaluation. By taking the five gaze estimation datasets, we train the model on four datasets and test on the remaining one respectively. We show the results of each test dataset in the column. Except for baselines, we evaluate three versions of UniGaze with backbones of ViT-Base (UniGaze-B), ViT-Large (UniGaze-L), and ViT-Huge (UniGaze-H).

general semantic representations, while FaRL is specialized for face analysis tasks and pre-trained on 20 M LAION-Face samples. We also include the ImageNet [16] pre-trained **ViTs** [17] for comparison. We use the Base (B), Large (L), and Huge (H) variants in our experiments.

## 4.2. Cross-Dataset Evaluation

We first assess the generalization of our UniGaze with the commonly used cross-dataset evaluation setting, where models are trained on one dataset and tested on an unseen dataset. Unlike previous studies that primarily report results for models trained on ETH-XGaze and/or Gaze360 [13, 51, 54, 83, 93], we evaluate models trained individually on five different gaze datasets and measure their cross-dataset performance on the remaining datasets. In Tab. 2, each subtable corresponds to a training dataset, with each column representing a different test dataset. We compare several baseline models with our UniGaze. Note the 60 subjects from XGaze-Dense used during pre-training (Sec. 3.1) are excluded from the XGaze Test set.

Although, commonly, larger models tend to achieve better performance in typical computer vision tasks, the large ViT-based models (ViT-H, DINO-B, MoCo-v3-B) do not consistently surpass the ResNet-50 on gaze estimation as revealed in Tab. 2. This discrepancy indicates that simply

increasing the model size without curated pre-training data is not effective for the downstream gaze estimation task. Especially, the FaRL-B pre-trained on large face-centric data still cannot effectively handle gaze information across diverse datasets.

By contrast, the proposed UniGaze-H, pre-trained on diverse facial datasets, demonstrates superior performance across nearly all settings. For completeness of fairness, we also show the UniGaze-B based on the ViT-B backbone, which still outperforms similar model size models DINO-B and FaRL-B. Note UniGaze-B outperforms UniGaze-H in some cases, which suggests again that our combination of large-model with our curated large-scale pre-training data effectively enhances ViT's ability to learn gaze-specific representations, rather than only increasing the size of the model. This answers our question that MAE pre-training with diverse facial data can strengthen the model's ability to capture fine-grained gaze geometry.
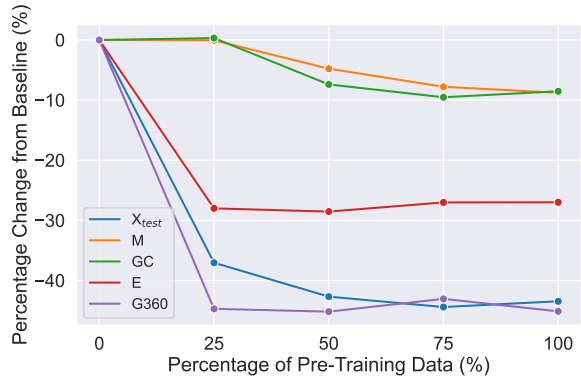
Besides, the results emphasize the importance of label range in training data. While pre-training improves generalization across domains, it alone is insufficient when the training gaze data has a narrow label range. For example, the model trained on MPIIFaceGaze (Tab. 2b) exhibits high errors on datasets such as XGaze Test (33.11 degrees) and Gaze360 (21.28 degrees), which have larger ranges of head poses and gaze directions than MPIIFaceGaze. Thus, while MAE pre-training improves ViT's capacity to learn robust representations, diverse label coverage in the gaze estimation training stage remains crucial for achieving consistent performance in gaze estimation across domains.

## 4.3. Comparison with Domain Generalization

We further assemble the experiment results from our UniGaze and compare them with the current SOTA domain generalization methods [3, 13, 83, 86, 93]. In Tab. 3, we pick the typical cross-dataset setting, where the training dataset is either ETH-XGaze or Gaze360, and test on MPI-

(a) The results when trained on ETH-XGaze.



(b) The results of the leave-one-dataset-out evaluation.

Figure 3. Effect of MAE pre-training data size on gaze estimation performance. The horizontal axis is the percentage of the pre-training data and the vertical axis is the percentage of the error reduction from the 0% baseline.

| Models \ Test | $X_{test}$ | $M_{test}$ | $GC_{test}$ | $E_{test}$ | $G360_{test}$ |
|---|---|---|---|---|---|
| ResNet-50 | 5.04 | 5.88 | 3.59 | 6.04 | 10.55 |
| GazeTR-50 | 4.63 | 5.82 | 3.52 | 6.06 | 10.39 |
| ViT-L | 5.24 | 5.42 | 3.66 | 6.30 | 10.54 |
| UniGaze-B | 4.75 | 5.40 | 3.37 | 5.52 | 9.64 |
| UniGaze-L | 4.67 | 5.12 | **3.17** | 5.33 | 9.29 |
| UniGaze-H | **4.46** | **5.08** | 3.20 | **5.16** | **9.07** |

Table 6. Results of the joint-dataset evaluation protocol. We train the model on the gathered training splits of all five gaze estimation datasets. We test the model on the test set of each dataset individually, shown in each column.

IFaceGaze and screen targets (CS) subset of EYEDIAP.

Across all datasets, our UniGaze-H model consistently achieves the lowest error, surpassing SOTA methods by significant margins. Notably, while the ViT-H model pre-trained on ImageNet alone does not outperform domain generalization methods, our UniGaze-H model pre-trained on large-scale data demonstrates significant improvements. Again, this highlights that MAE pre-training on large-scale data greatly enhances ViT's ability for domain generalization, validating the effectiveness of our approach in improving gaze estimation across diverse domains. Moreover, even UniGaze-B, despite having a model size comparable to other SOTA methods, still achieves substantially better performance.

## 4.4. Within-Dataset Evaluation

To align with other SOTA methods, the splits for MPI-IFaceGaze, ETH-XGaze, and EYEDIAP used in this section differ from the subject splits used elsewhere in the paper. For GazeCapture and Gaze360, it is the same as defined in Sec. 4.1. For MPIIFaceGaze, we use the leave-

one-subject-out protocol [1, 66, 73]. For ETH-XGaze, we follow the train-test split provided in [73]. For EYEDIAP, we apply the four-fold validation scheme from [12, 14]. The results in Tab. 4 demonstrate that UniGaze-H outperforms SOTA on most datasets, except MPIIFaceGaze, indicating that the proposed pre-training enhances generalizability even within the same dataset.

## 4.5. Leave-One-Dataset-Out Evaluation

We conduct a *leave-one-dataset-out* evaluation to further assess the generalizability of UniGaze. Using five gaze estimation datasets, we train the model on a combination of four datasets and test it on the held-out dataset. In this manner, we assess the upper-bound performance that can be achieved in each domain by maximizing the use of existing datasets. We compare the baseline models of ResNet-50, GazeTR-50, and ViT-H that are pre-trained on ImageNet. We present the three versions of UniGaze with different backbones of ViT-Base, ViT-Large, and ViT-Huge.

We show the results in Tab. 5 with each column showing the results on each test dataset. Across test datasets, the proposed UniGaze with three backbone sizes consistently surpasses all baselines, achieving substantial error reductions. By comparing the UniGaze with three sizes of backbones, we can see that the trend clearly shows the larger model achieves better performances in general. Notably, large-margin improvements from UniGaze happen in the XGaze and Gaze360 test sets, highlighting that UniGaze significantly enhances ViT's generalization ability across diverse domains by fostering robust representations of facial appearance and image quality.

## 4.6. Joint-Dataset Evaluation

Building upon the leave-one-dataset-out evaluation, we conduct a new evaluation protocol *joint-dataset* evaluation. In this setting, we gather the training sets of all five datasets,

including XGaze Train, MPIIFaceGaze Train, GazeCapture Train, EYEDIAP Train, and Gaze360 Train to form a large and comprehensive training set. The model trained on this large training set is expected to be the most powerful gaze estimator that we can acquire. We test the trained model on the test sets of each dataset. It has another meaning of generalization, in which one model performs gaze estimation across multiple data domains.

Table 6 shows the evaluation results where each column represents a specific test dataset. As expected, our UniGaze-H model improves over other methods, achieving the lowest error on most test datasets compared to all baselines and smaller backbone versions. Since the model is trained on data that covers the domains of all test datasets, the resulting errors are generally lower and saturated compared to other settings. It aligns with our purpose of training one model for all domains with the lowest error.

### 4.7. Ablation Studies

#### 4.7.1. Effect of Pre-Training Data Size

The large-scale pre-training with MAE is the most critical factor for our UniGaze. To analyze the impact of MAE pre-training data size, we vary the amount of data used for the pre-training stage. With the full pre-training data of around $1.6\,\mathrm{M}$ images, we randomly sample subsets from each component dataset (Sec. 3.1) to create 25%, 50%, and 75% subsets of the full data. We use the UniGaze-L as the backbone, and the 0% refers to the ImageNet pre-trained ViT-L.

We present two experiment settings in Fig. 3, the XGaze training (left) and the leave-one-dataset-out setting (right). Beginning with the 0% baseline, we illustrate the percentage reduction in error (vertical axis) as the pre-training dataset size (horizontal axis) increases, allowing us to capture and compare relative performance improvements across various pre-training levels.

Overall, the results indicate that the pre-trained model on a larger subset consistently achieves lower error across all test domains. There are notable performance improvements with 25% and 50% of the data, and the improvement gap becomes smaller when reaching the 75% subset to the full data, which is expected due to a sufficient amount of data diversity for the training. These results confirm that increasing the amount of data in MAE pre-training strengthens ViT's representation learning, leading to improved accuracy and generalization across diverse gaze estimation tasks.

#### 4.7.2. Effect of Pre-Training Data Attributes

To better understand the performance gap from FaRL-B [96], we examine the impact of different data in Tab. 7. To assess the effect of the synthetic data, we first limit the pre-training dataset to the real datasets: CelebV-Text, VFHQ, and VGGFace2 (*Real*). To evaluate the impact of pre-processing specifically for gaze estimation, we also assess

| Models \ Test | $\mathbf{X_{test}}$ | **M** | **GC** | **E** | **G360** |
|---|---|---|---|---|---|
| FaRL-B [96] | 19.08 | 5.09 | 6.08 | 8.18 | 18.28 |
| UniGaze-B (*Real, limited poses*) | 16.70 | 5.29 | 6.87 | 6.57 | 13.78 |
| UniGaze-B (*Real, w/o norm.*) | 14.56 | 4.95 | 6.27 | 6.93 | 14.48 |
| UniGaze-B (*Real*) | 11.95 | 4.86 | 6.14 | **6.26** | 12.71 |
| UniGaze-B | **11.78** | **4.73** | **5.86** | 6.31 | **12.41** |

Table 7. Comparison of different formats of the input data in the leave-one-dataset-out evaluation setting. Each column shows the results on each test dataset. For fair compassion with FaRL-B [96], we use the ViT-B backbone (UniGaze-B). *Real* stands for the combination of our real datasets CelebV-Text, VFHQ, and VGGFace2.

the performance when pre-training is conducted using the loose face bounding boxes directly, without applying any data normalization [90] (*Real, w/o norm.*). To evaluate the effect of wide head pose range, we limit head pose variability by filtering out samples with a pitch-yaw $L_2$ norm exceeding 10 degrees, reducing the dataset size to about 20% (*Real, limited poses*). Note that the different amounts of data could also affect the model's performance. We use the leave-one-dataset-out evaluation protocol, given the trade-off between the task difficulty and simplification.

The results demonstrate that the UniGaze-B trained on full data performs best compared to other baselines in almost all settings. Models without normalization (*Real, w/o norm.*) show substantially degraded performance despite using identical data. Models with limited pose ranges (*Real, limited poses*) also perform significantly worse than those trained on a more diverse range. Integrating real and synthetic data yields the best results, which can be attributed to combining naturalistic appearance variations from real data with controlled pose variations from synthetic data.

## 5. Conclusion

In this paper, we provide the first extensive study of the self-supervised large-scale pre-training on gaze estimation. Through our extensive experimentation, we have established key principles for effective pre-training in gaze estimation. With careful data curation for the MAE pre-training, the proposed UniGaze achieves distinguished performances for cross-dataset, leave-one-dataset-out, and joint-dataset evaluations compared to current SOTAs. Interestingly, we show the importance of the pre-training data selection and pre-processing for the final performance, rather than simply gathering a large amount of data. Looking forward, we believe there remains significant potential for further refinement of pre-training data selection based on the principles identified in this work. Another potential improvement option could be using a large-size face image input for the high-resolution scenarios.

# References

[1] Ahmed A Abdelrahman, Thorsten Hempel, Aly Khalifa, Ayoub Al-Hamadi, and Laslo Dinges. L2cs-net: Fine-grained gaze estimation in unconstrained environments. In *2023 8th International Conference on Frontiers of Signal Processing (ICFSP)*, pages 98–102. IEEE, 2023. 6, 7, 1

[2] Shumeet Baluja and Dean Pomerleau. Non-intrusive gaze tracking using artificial neural networks. *Proc. NIPS*, 6, 1993. 1

[3] Yiwei Bao and Feng Lu. From feature to gaze: A generalizable replacement of linear layer for gaze estimation. In *Proc. CVPR*, pages 1409–1418, 2024. 2, 3, 6

[4] Yiwei Bao and Feng Lu. Unsupervised gaze representation learning from multi-view face images. In *Proc. CVPR*, pages 1419–1428, 2024. 3

[5] Anna Belardinelli. Gaze-based intention estimation: principles, methodologies, and applications in hri. *ACM Transactions on Human-Robot Interaction*, 13(3):1–30, 2024. 1

[6] David Beniaguev. Synthetic faces high quality - text 2 image (sfhq-t2i) dataset, 2024. 3, 4

[7] Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In *Proc. ICCV*, 2017. 4

[8] Zhixi Cai, Shreya Ghosh, Kalin Stefanov, Abhinav Dhall, Jianfei Cai, Hamid Rezatofighi, Reza Haffari, and Munawar Hayat. Marlin: Masked autoencoder for facial video representation learning. In *Proc. CVPR*, pages 1493–1504, 2023. 2, 3

[9] Qiong Cao, Li Shen, Weidi Xie, Omkar M Parkhi, and Andrew Zisserman. Vggface2: A dataset for recognising faces across pose and age. In *Proc. FG*, pages 67–74. IEEE, 2018. 3

[10] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proc. ICCV*, pages 9650–9660, 2021. 2, 5

[11] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *Proc. ICCV*, pages 9640–9649, 2021. 2, 5

[12] Yihua Cheng and Feng Lu. Gaze estimation using transformer. In *Proc. ICPR*, pages 3341–3347. IEEE, 2022. 3, 5, 7

[13] Yihua Cheng, Yiwei Bao, and Feng Lu. Puregaze: Purifying gaze feature for generalizable gaze estimation. In *Proc. AAAI*, pages 436–443, 2022. 2, 3, 5, 6, 4

[14] Yihua Cheng, Haofei Wang, Yiwei Bao, and Feng Lu. Appearance-based gaze estimation with deep learning: A review and benchmark. *PAMI*, 2024. 1, 7

[15] Eunji Chong, Yongxin Wang, Nataniel Ruiz, and James M Rehg. Detecting attended visual targets in video. In *Proc. CVPR*, pages 5396–5406, 2020. 4, 7

[16] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proc. CVPR*, pages 248–255. Ieee, 2009. 3, 6

[17] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *Proc. ICLR*, 2021. 2, 3, 6

[18] Yuxin Fang, Li Dong, Hangbo Bao, Xinggang Wang, and Furu Wei. Corrupted image modeling for self-supervised visual pre-training. *Proc. ICLR*, 2023. 3

[19] Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning at scale. In *Proc. CVPR*, pages 19358–19369, 2023. 2, 3

[20] Tobias Fischer, Hyung Jin Chang, and Yiannis Demiris. Rt-gene: Real-time eye gaze estimation in natural environments. In *Proc. ECCV*, pages 334–352, 2018. 2

[21] Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981. 4

[22] Kenneth Alberto Funes Mora, Florent Monay, and Jean-Marc Odobez. Eyediap: A database for the development and evaluation of gaze estimation algorithms from rgb and rgb-d cameras. In *Proc. ETRA*, pages 255–258, 2014. 2, 5

[23] Zheng Gao and Ioannis Patras. Self-supervised facial representation learning with facial region awareness. In *Proc. CVPR*, pages 2081–2092, 2024. 2, 3

[24] Shreya Ghosh, Abhinav Dhall, Munawar Hayat, Jarrod Knibbe, and Qiang Ji. Automatic gaze analysis: A survey of deep learning based approaches. *PAMI*, 46(1):61–84, 2023. 1

[25] Yiran Guan, Zhuoguang Chen, Wenzheng Zeng, Zhiguo Cao, and Yang Xiao. End-to-end video gaze estimation via capturing head-face-eye spatial-temporal interaction context. *IEEE Signal Processing Letters*, 30:1687–1691, 2023. 6

[26] Katarzyna Harezlak and Pawel Kasprowski. Application of eye tracking in medicine: A survey, research issues and challenges. *Computerized Medical Imaging and Graphics*, 65: 176–190, 2018. 1

[27] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. CVPR*, pages 770–778, 2016. 3, 5

[28] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proc. CVPR*, pages 9729–9738, 2020. 3

[29] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proc. CVPR*, pages 16000–16009, 2022. 2, 3, 4

[30] Yoichiro Hisadome, Tianyi Wu, Jiawei Qin, and Yusuke Sugano. Rotation-constrained cross-view feature fusion for multi-view appearance-based gaze estimation. In *Proc. WACV*, pages 5985–5994, 2024. 3

[31] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Proc. NIPS*, 33:6840–6851, 2020. 3

[32] Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metze, and Christoph Feichtenhofer. Masked autoencoders that listen. *Proc. NIPS*, 35:28708–28720, 2022. 3

[33] Yangzhou Jiang, Yinxin Lin, Yaoming Wang, Teng Li, Bilian Ke, and Bingbing Ni. Learning unsupervised gaze representation via eye mask driven information bottleneck. *arXiv preprint arXiv:2407.00315*, 2024. 2, 3

[34] Shiwei Jin, Zhen Wang, Lei Wang, Ning Bi, and Truong Nguyen. Redirtrans: Latent-to-latent translation for gaze and head redirection. In *Proc. CVPR*, pages 5547–5556, 2023. 2

[35] Swati Jindal and Roberto Manduchi. Contrastive representation learning for gaze estimation. In *Annual Conference on Neural Information Processing Systems*, pages 37–49. PMLR, 2023. 3

[36] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020. 2

[37] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *Proc. CVPR*, 2020. 3, 4

[38] Fengfeng Ke, Ruohan Liu, Zlatko Sokolikj, Ibrahim Dahlstrom-Hakki, and Maya Israel. Using eye-tracking in education: review of empirical research and technology. *Educational technology research and development*, pages 1–36, 2024. 1

[39] Petr Kellnhofer, Adria Recasens, Simon Stent, Wojciech Matusik, and Antonio Torralba. Gaze360: Physically unconstrained gaze estimation in the wild. In *Proc. ICCV*, pages 6912–6921, 2019. 2, 5

[40] Rawal Khirodkar, Timur Bagautdinov, Julieta Martinez, Su Zhaoen, Austin James, Peter Selednik, Stuart Anderson, and Shunsuke Saito. Sapiens: Foundation for human vision models. In *Proc. ECCV*, pages 206–228. Springer, 2025. 2, 3, 4

[41] Joohwan Kim, Michael Stengel, Alexander Majercik, Shalini De Mello, David Dunn, Samuli Laine, Morgan McGuire, and David Luebke. Nvgaze: An anatomically-informed dataset for low-latency, near-eye gaze estimation. In *Proceedings of the 2019 CHI conference on human factors in computing systems*, pages 1–12, 2019. 3

[42] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proc. ICLR*, 2015. 5, 4

[43] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proc. ICCV*, pages 4015–4026, 2023. 2

[44] Rakshit Kothari, Shalini De Mello, Umar Iqbal, Wonmin Byeon, Seonwook Park, and Jan Kautz. Weakly-supervised physically unconstrained gaze estimation. In *Proc. CVPR*, pages 9980–9989, 2021. 3

[45] Kyle Krafka, Aditya Khosla, Petr Kellnhofer, Harini Kannan, Suchendra Bhandarkar, Wojciech Matusik, and Antonio Torralba. Eye tracking for everyone. In *Proc. CVPR*, pages 2176–2184, 2016. 2, 3, 5

[46] Oren Kraus, Kian Kenyon-Dean, Saber Saberian, Maryam Fallah, Peter McLean, Jess Leung, Vasudev Sharma, Ayla Khan, Jia Balakrishnan, Safiye Celik, et al. Masked autoencoders for microscopy are scalable learners of cellular biology. In *CVPR*, pages 11757–11768, 2024. 3

[47] Alexander C Li, Mihir Prabhudesai, Shivam Duggal, Ellis Brown, and Deepak Pathak. Your diffusion model is secretly a zero-shot classifier. In *Proc. ICCV*, pages 2206–2217, 2023. 3

[48] Daiqing Li, Huan Ling, Amlan Kar, David Acuna, Seung Wook Kim, Karsten Kreis, Antonio Torralba, and Sanja Fidler. Dreamteacher: Pretraining image backbones with deep generative models. In *Proc. ICCV*, pages 16698–16708, 2023. 2, 3

[49] Junxuan Li, Chen Cao, Gabriel Schwartz, Rawal Khirodkar, Christian Richardt, Tomas Simon, Yaser Sheikh, and Shunsuke Saito. Uravatar: Universal relightable gaussian codec avatars. *arXiv preprint arXiv:2410.24223*, 2024. 4, 6

[50] Zhi-Yi Lin, Jouh Yeong Chew, Jan van Gemert, and Xucong Zhang. Gazehta: End-to-end gaze target detection with head-target association. *arXiv preprint arXiv:2404.10718*, 2024. 2

[51] Ruicong Liu, Yiwei Bao, Mingjie Xu, Haofei Wang, Yunfei Liu, and Feng Lu. Jitter does matter: Adapting gaze estimation to new domains. *arXiv preprint arXiv:2210.02082*, 2022. 6

[52] Ruicong Liu, Yunfei Liu, Haofei Wang, and Feng Lu. Pnpga+: Plug-and-play domain adaptation for gaze estimation using model variants. *PAMI*, 2024. 3

[53] Ruicong Liu, Haofei Wang, and Feng Lu. From gaze jitter to domain adaptation: Generalizing gaze estimation by manipulating high-frequency components. *IJCV*, pages 1–16, 2024. 3

[54] Yunfei Liu, Ruicong Liu, Haofei Wang, and Feng Lu. Generalizing gaze estimation with outlier-guided collaborative adaptation. In *Proc. ICCV*, pages 3835–3844, 2021. 6

[55] Agisoft LLC. Agisoft metashape. https://www.agisoft.com/, 2024. 3, 4

[56] Feng Lu, Yusuke Sugano, Takahiro Okabe, and Yoichi Sato. Head pose-free appearance-based gaze sensing via eye image synthesis. In *Proc. ICPR*, pages 1008–1011. IEEE, 2012. 1

[57] Feng Lu, Takahiro Okabe, Yusuke Sugano, and Yoichi Sato. Learning gaze biases with head motion for head pose-free gaze estimation. *Image and Vision Computing*, 32(3):169–179, 2014. 1

[58] Mathias N Lystbæk, Peter Rosenberg, Ken Pfeuffer, Jens Emil Grønbæk, and Hans Gellersen. Gaze-hand alignment: Combining eye gaze and mid-air pointing for interacting with menus in augmented reality. *Proceedings of the ACM on Human-Computer Interaction*, 6(ETRA):1–18, 2022. 1

[59] Päivi Majaranta and Andreas Bulling. Eye tracking and eye-based human–computer interaction. In *Advances in physiological computing*, pages 39–65. Springer, 2014. 1

[60] Seonwook Park, Shalini De Mello, Pavlo Molchanov, Umar Iqbal, Otmar Hilliges, and Jan Kautz. Few-shot adaptive gaze estimation. In *Proc. ICCV*, pages 9368–9377, 2019. 2, 3, 5

[61] Jiawei Qin, Takuru Shimoyama, and Yusuke Sugano. Learning-by-novel-view-synthesis for full-face appearance-based 3d gaze estimation. In *Proc. CVPRW*, pages 4981–4991, 2022. 2, 3, 4

[62] Jiawei Qin, Takuru Shimoyama, Xucong Zhang, and Yusuke Sugano. Domain-adaptive full-face gaze estimation via novel-view-synthesis and feature disentanglement. *arXiv preprint arXiv:2305.16140*, 2023. 3, 4

[63] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proc. CVPR*, pages 10684–10695, 2022. 3

[64] Alessandro Ruzzi, Xiangwei Shi, Xi Wang, Gengyan Li, Shalini De Mello, Hyung Jin Chang, Xucong Zhang, and Otmar Hilliges. Gazenerf: 3d-aware gaze redirection with neural radiance fields. In *Proc. CVPR*, pages 9676–9685, 2023. 2

[65] Pavan Kumar Sharma and Pranamesh Chakraborty. A review of driver gaze estimation and application in gaze behavior understanding. *Engineering Applications of Artificial Intelligence*, 133:108117, 2024. 1

[66] Yichen Shi, Feifei Zhang, Wenming Yang, Guijin Wang, and Nan Su. Agent-guided gaze estimation network by two-eye asymmetry exploration. In *2024 IEEE International Conference on Image Processing (ICIP)*, pages 2320–2326. IEEE, 2024. 6, 7, 1

[67] Mannat Singh, Quentin Duval, Kalyan Vasudev Alwala, Haoqi Fan, Vaibhav Aggarwal, Aaron Adcock, Armand Joulin, Piotr Dollár, Christoph Feichtenhofer, Ross Girshick, et al. The effectiveness of mae pre-pretraining for billion-scale pretraining. In *Proc. ICCV*, pages 5484–5494, 2023. 2, 3

[68] Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of neural networks using large learning rates. In *Artificial intelligence and machine learning for multi-domain operations applications*, pages 369–386. SPIE, 2019. 5, 4

[69] Siddharth Srivastava and Gaurav Sharma. Omnivec: Learning robust representations with cross modal sharing. In *Proc. WACV*, pages 1236–1248, 2024. 2, 3

[70] Haomiao Sun, Mingjie He, Tianheng Lian, Hu Han, and Shiguang Shan. Face-mllm: A large face perception model. *arXiv preprint arXiv:2410.20717*, 2024. 2, 3

[71] Kar-Han Tan, David J Kriegman, and Narendra Ahuja. Appearance-based eye gaze estimation. In *Sixth IEEE Workshop on Applications of Computer Vision, 2002.(WACV 2002). Proceedings.*, pages 191–195. IEEE, 2002. 1

[72] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. *Proc. NIPS*, 35: 10078–10093, 2022. 3

[73] Evangelos Ververas, Polydefkis Gkagkos, Jiankang Deng, Michail Christos Doukas, Jia Guo, and Stefanos Zafeiriou. 3dgazenet: Generalizing 3d gaze estimation with weak-supervision from synthetic views. In *ECCV*, pages 387–404. Springer, 2025. 6, 7, 1, 3, 4

[74] Hengfei Wang, Zhongqun Zhang, Yihua Cheng, and Hyung Jin Chang. High-fidelity eye animatable neural radiance fields for human face. *arXiv preprint arXiv:2308.00773*, 2023. 2

[75] Kang Wang, Rui Zhao, Hui Su, and Qiang Ji. Generalizing eye tracking with bayesian adversarial learning. In *Proc. CVPR*, pages 11907–11916, 2019. 1

[76] Yaoming Wang, Yangzhou Jiang, Jin Li, Bingbing Ni, Wenrui Dai, Chenglin Li, Hongkai Xiong, and Teng Li. Contrastive regression for domain adaptation on gaze estimation. In *Proc. CVPR*, pages 19376–19385, 2022. 3

[77] Yue Wang, Jinlong Peng, Jiangning Zhang, Ran Yi, Liang Liu, Yabiao Wang, and Chengjie Wang. Toward high quality facial representation learning. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 5048–5058, 2023. 2, 3

[78] Erroll Wood, Tadas Baltrusaitis, Xucong Zhang, Yusuke Sugano, Peter Robinson, and Andreas Bulling. Rendering of eyes for eye-shape registration and gaze estimation. In *Proc. ICCV*, pages 3756–3764, 2015. 3

[79] Erroll Wood, Tadas Baltrušaitis, Charlie Hewitt, Sebastian Dziadzio, Thomas J Cashman, and Jamie Shotton. Fake it till you make it: face analysis in the wild using synthetic data alone. In *Proc. ICCV*, pages 3681–3691, 2021. 3, 4

[80] Weilai Xiang, Hongyu Yang, Di Huang, and Yunhong Wang. Denoising diffusion autoencoders are unified self-supervised learners. In *Proc. ICCV*, pages 15802–15812, 2023. 2, 3

[81] Liangbin Xie, Xintao Wang, Honglun Zhang, Chao Dong, and Ying Shan. Vfhq: A high-quality dataset and benchmark for video face super-resolution. In *Proc. CVPR*, pages 657–666, 2022. 3, 4

[82] Mingjie Xu and Feng Lu. Gaze from origin: Learning for generalized gaze estimation by embedding the gaze frontalization process. In *Proc. AAAI*, pages 6333–6341, 2024. 3

[83] Mingjie Xu, Haofei Wang, and Feng Lu. Learning a generalized gaze estimator from gaze-consistent feature. In *Proc. AAAI*, pages 3027–3035, 2023. 2, 3, 5, 6

[84] Pengwei Yin, Jiawu Dai, Jingjing Wang, Di Xie, and Shiliang Pu. Nerf-gaze: A head-eye redirection parametric model for gaze estimation. *arXiv preprint arXiv:2212.14710*, 2022. 2

[85] Pengwei Yin, Jingjing Wang, Guanzhong Zeng, Di Xie, and Jiang Zhu. Lg-gaze: Learning geometry-aware continuous prompts for language-guided gaze estimation. In *Proc. ECCV*, 2024. 2, 3, 6

[86] Pengwei Yin, Guanzhong Zeng, Jingjing Wang, and Di Xie. Clip-gaze: Towards general gaze estimation via visual-linguistic model. In *Proc. AAAI*, pages 6729–6737, 2024. 2, 3, 6

[87] Jianhui Yu, Hao Zhu, Liming Jiang, Chen Change Loy, Weidong Cai, and Wayne Wu. Celebv-text: A large-scale facial text-video dataset. In *Proc. CVPR*, pages 14805–14814, 2023. 3, 4, 2

[88] Xucong Zhang, Yusuke Sugano, Mario Fritz, and Andreas Bulling. Appearance-based gaze estimation in the wild. In *Proc. CVPR*, pages 4511–4520, 2015. 1, 2

11

[89] Xucong Zhang, Yusuke Sugano, Mario Fritz, and Andreas Bulling. It's written all over your face: Full-face appearance-based gaze estimation. In *Proc. CVPRW*, pages 51–60, 2017. 5

[90] Xucong Zhang, Yusuke Sugano, and Andreas Bulling. Revisiting data normalization for appearance-based gaze estimation. In *Proc. ETRA*, 2018. 2, 4, 8

[91] Xucong Zhang, Yusuke Sugano, Mario Fritz, and Andreas Bulling. Mpiigaze: Real-world dataset and deep appearance-based gaze estimation. *IEEE TPAMI*, 41(1):162–175, 2019. 2

[92] Xucong Zhang, Seonwook Park, Thabo Beeler, Derek Bradley, Siyu Tang, and Otmar Hilliges. Eth-xgaze: A large scale dataset for gaze estimation under extreme head pose and gaze variation. In *Proc. ECCV*, pages 365–381. Springer, 2020. 2, 3, 5

[93] Ruijie Zhao, Pinyan Tang, and Sihui Luo. Improving domain generalization on gaze estimation via branch-out auxiliary regularization. *arXiv preprint arXiv:2405.01439*, 2024. 2, 3, 5, 6, 4

[94] Wenliang Zhao, Yongming Rao, Zuyan Liu, Benlin Liu, Jie Zhou, and Jiwen Lu. Unleashing text-to-image diffusion models for visual perception. In *Proc. ICCV*, pages 5729–5739, 2023. 2, 3

[95] Yufeng Zheng, Seonwook Park, Xucong Zhang, Shalini De Mello, and Otmar Hilliges. Self-learning transformations for improving gaze and head redirection. *Proc. NIPS*, 33: 13127–13138, 2020. 2

[96] Yinglin Zheng, Hao Yang, Ting Zhang, Jianmin Bao, Dongdong Chen, Yangyu Huang, Lu Yuan, Dong Chen, Ming Zeng, and Fang Wen. General facial representation learning in a visual-linguistic manner. In *Proc. CVPR*, pages 18697–18709, 2022. 2, 3, 5, 8

[97] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE TPAMI*, 2017. 4

[98] Yang Zhou, Zichong Chen, and Hui Huang. Deformable one-shot face stylization via dino semantic guidance. In *Proc. CVPR*, pages 7787–7796, 2024. 2

| Training Data \ Test | $X_{test}$ | $M_{test}$ | $GC_{test}$ | $E_{test}$ | $G360_{test}$ |
|---|---|---|---|---|---|
| **ResNet-50** | | | | | |
| *same-domain* | 5.25 | 5.11 | 3.49 | 8.51 | 11.87 |
| *leave-one-dataset-out* | 16.31 (↑210.7%) | 6.23 (↑21.9%) | 6.35 (↑82.0%) | 8.25 (↓3.1%) | 20.38 (↑71.7%) |
| *joint-dataset* | **5.04** (↓4.0%) | 5.88 (↑15.1%) | 3.59 (↑2.9%) | **6.04** (↓29.0%) | **10.55** (↓11.1%) |
| **UniGaze-H** | | | | | |
| *same-domain* | 4.62 | 5.19 | 3.01 | 6.11 | 9.44 |
| *leave-one-dataset-out* | 11.29 (↑144.4%) | 5.22 (↑0.6%) | 5.13 (↑70.4%) | 6.14 (↑0.5%) | 13.12 (↑39.0%) |
| *joint-dataset* | **4.46** (↓3.5%) | **5.08** (↓2.1%) | 3.20 (↑6.3%) | **5.16** (↓15.6%) | **9.07** (↓3.9%) |

Table 1. Comparison of different training data configurations for gaze estimation. Each column represents a specific test dataset: XGaze Test, MPIIFaceGaze Test, GazeCapture Test, EYEDIAP Test, and Gaze360 Test. Each row corresponds to a training configuration: *Same-domain* means training on the same domain as the test set, *leave-one-dataset-out* means training on the remaining four datasets other than the test set, and *joint-dataset* means training on the aggregated Train split of all five datasets. The percentages in parentheses indicate the reduction or increment compared to the *same-domain* results, where lower errors indicate better performance. For the *leave-one-dataset-out* configuration, the errors reported here are on the Test splits, while the main paper reports errors on the entire dataset.

# Supplementary Materials

In this supplementary material, we first provide an analysis of the effect of combining multiple domains. Then, we include additional ablations to investigate the effects of color-jitter augmentation and pixel normalization during the pre-training. Finally, we present qualitative results, highlighting images captured under diverse and challenging conditions.

## A. Analysis on Combining Multiple Domains

We analyze the effect of different training data configurations on gaze estimation performance. Specifically, we compare three configurations: training on the same domain (*same-domain*), training on multiple domains excluding the testing domain (*leave-one-dataset-out*), and training on multiple domains including the testing domain (*joint-dataset*).

Table 1 shows the comparison of gaze errors for these configurations. Each column corresponds to a specific test dataset: XGaze Test, MPIIFaceGaze Test, GazeCapture Test, EYEDIAP Test, and Gaze360 Test, while each row represents a training configuration. This *same-domain* setting is different from the *within-dataset* in the main paper. We use the splits defined in Sec. 4.1 of the main paper. Especially, please note that for MPIIFaceGaze dataset, we train the model on the first 10 subjects and test on the remaining five subjects, different from the typical leave-one-subject-out protocol [1, 66, 73].

The percentages in parentheses indicate the reduction or increment compared to the *same-domain* results, where lower errors indicate better performance. Note that, for the *leave-one-dataset-out* configuration, errors on the entire left-out dataset are reported in the main paper, but here we present errors on the Test split to align with the other configurations that require dataset splits.
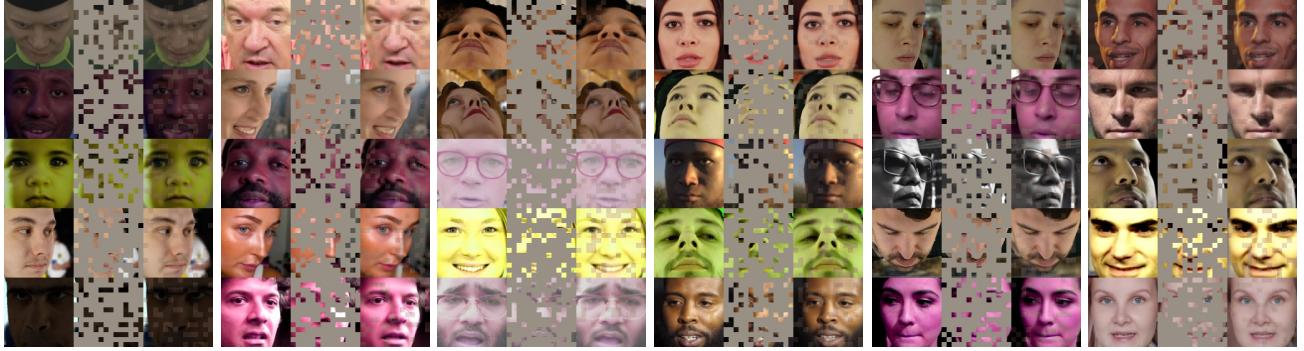
**Same-domain** In general, training and testing on the same domain (*same-domain*) yields the best results, even though datasets combined from multiple domains have the potential to be more diverse. This emphasizes the persistent challenge of achieving optimal performance when using data from different domains. The exception observed for $E_{test}$ with the ResNet-50 backbone may be attributed to the limited number of samples in the EYEDIAP Train split.
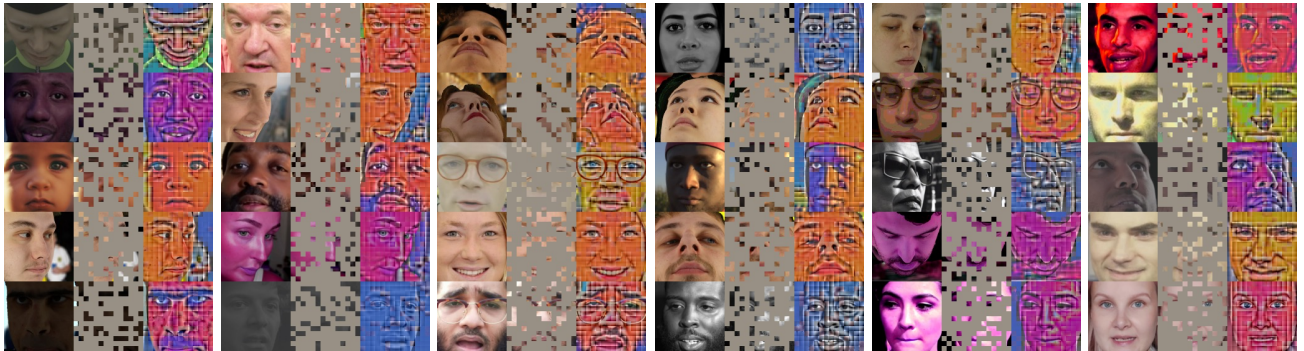
**Leave-one-dataset-out** In the *leave-one-dataset-out* configuration, we observe varying tendencies across different test datasets. Some datasets achieve errors comparable to the *same-domain* results, while others remain challenging. For instance, for $M_{test}$ and $E_{test}$, which are relatively less complex, the remaining four datasets provide sufficient information to achieve good performance. In contrast, for $X_{test}$, $GC_{test}$, and $G360_{test}$, the remaining four datasets fail to fully capture the critical factors required for optimal performance. This variation highlights the strong dependence of performance on the attributes of the training data.

Importantly, our UniGaze-H demonstrates a smaller performance gap compared to ResNet-50 in most cases, with the only exception being EYEDIAP, where the difference is marginal. This suggests that UniGaze-H is better equipped to learn gaze representations from out-of-domain data with less overfitting, underscoring its enhanced generalization capability.

**Joint-dataset** Overall, the *joint-dataset* configuration demonstrates significant promise, creating a single model robust across multiple test domains. For UniGaze-H, the

(a) MAE reconstruction examples without pixel normalization.


(b) MAE reconstruction examples with pixel normalization (Proposed).

Figure 1. Examples comparison of the pixel normalization during the MAE pre-training. The left, middle, and right columns show the original image, masked input, and the reconstructed image, respectively.

| Real | Syn. | NV. | M | GC | E | G360 |
|------|------|-----|------|------|------|-------|
| ✓ | | | 6.79 | 7.81 | 6.86 | 12.93 |
| ✓ | ✓ | | 6.57 | 7.37 | **6.51** | 13.23 |
| ✓ | ✓ | ✓ | **6.21** | **7.35** | 6.64 | **12.18** |

Table 2. We ablate the pre-training facial datasets by comparing real, synthetic, and novel-rendered images. The comparison is performed on the UniGaze-B network, followed by training on XGaze. The last row represents the full-dataset setting.

only exception is $GC_{test}$, where the *joint-dataset* configuration produces a slightly higher error (3.01→3.20). Although this suggests some negative effects from the other four datasets, the effects remain marginal. While the improvement percentages for UniGaze-H are smaller compared to ResNet-50, the absolute errors are consistently lower.

## B. Additional Ablation Studies on Pre-Training

**Effect of Pre-Training Dataset Composition** Beyond the overall pre-training dataset size, the composition of the dataset also plays a critical role in learning effective gaze representations. To investigate the impact of different fa-
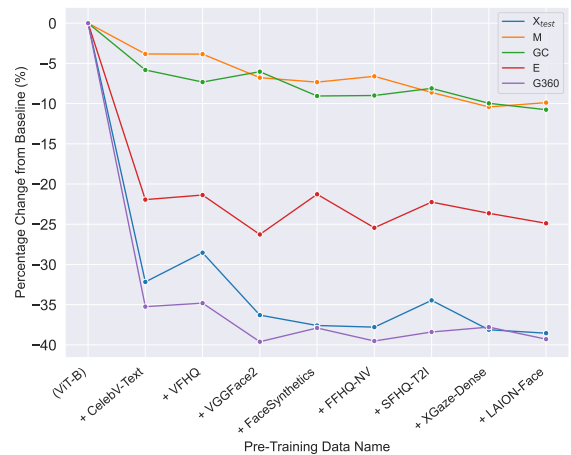

Figure 2. Effect of MAE pre-training dataset composition on downstream gaze estimation performance. The horizontal axis represents the incremental accumulation of datasets, while the vertical axis shows the percentage reduction in error relative to the first CelebV-Text dataset [87].

cial dataset components, we conduct an experiment where we incrementally accumulate datasets during the MAE pre-training stage and analyze their effect on the down-

stream gaze estimation performance. Starting with CelebV-Text [87], we progressively add datasets for pre-training and evaluate model separately on gaze estimation. Each pre-trained model is subsequently trained on gaze datasets using the same *leave-one-dataset-out* protocol. Figure 2 illustrates the error change across different test sets as more datasets are included in pre-training.

Overall, the results indicate that adding more diverse data during pre-training generally enhances gaze generalization. However, there are exceptions that adding a model can result in increased error for specific test sets. For example, adding VFHQ increases the error on the XGaze Test set from 12.65 to 13.32, while including SFHQ-T2I causes performance fluctuations across different benchmarks. This suggests that certain dataset attributes may not align well with particular test distributions, leading to suboptimal transferability. On the other hand, datasets such as VG-GFace2 and XGaze-Dense provide performance improvements on most test sets. Additionally, performance gains becomes marginal as the dataset number increases, aligning with the analysis of pre-training data size in main paper.

In conclusion, dataset diversity plays a crucial role in improving MAE pre-training for gaze estimation. A more detailed analysis of dataset attributes and their impact on gaze estimation remains an open research question, which we leave for future work. Nonetheless, our empirical results suggest that increasing data diversity in pre-training tends to improve model performance across various test domains.

**Effect of Novel-View Synthesis Data in Pre-Training** To examine the effect of novel-view synthesis in pre-training data, we conduct further experiments separating these two elements. In Tab. 2, we conduct an ablation study by varying data subsets during the pre-training: real datasets (CelebV-Text, VGGFace2, and VFHQ), synthetic datasets (FaceSynthetics and SFHQ-T2I), and novel-view-rendered datasets (FFHQ-NV and XGaze-Dense). We use the UniGaze-B to conduct the experiment due to its time efficiency. After pre-training, we train on XGaze and test on the rest of the four datasets.

The results further clarify the effect of different data types on the model's generalizability. Adding synthetic data (*Real + Syn.*) reduces errors in several test domains compared to using only real data, suggesting the variability of the synthetic data contributes to generalization. Further incorporating novel-view data (*Real + Syn. + NV*) provides additional performance gains, especially in head-pose generalization, likely due to the expanded range of facial orientations. This finding supports the idea that a mix of real, synthetic, and novel-view data in MAE pre-training strengthens ViT's representation learning.

| Color-Jitter | Pixel Norm. | M | GC | E | G360 |
|:---:|:---:|:---:|:---:|:---:|:---:|
| ✗ | ✗ | 7.52 | 8.01 | 8.56 | 14.14 |
| ✓ | ✗ | 7.17 | 8.23 | 8.03 | 14.03 |
| ✗ | ✓ | 7.18 | 7.94 | 8.05 | 13.66 |
| ✓ | ✓ | **6.21** | **7.35** | **6.64** | **12.18** |

Table 3. Ablation studies on the pre-training, comparing the effect of the color-jitter augmentation and the pixel normalization. During the gaze estimation training, we train the model using XGaze and test on the other four datasets to evaluate the generalizability.

**Effect of Pixel Normalization** The patch normalization technique is applied during the MAE pre-training as suggested in [29] which is different from reconstructing the natural image, as shown in Fig. 1. We compare models pre-trained with and without patch normalization to investigate its impact.

**Effect of Color-Jitter Augmentation** Color jittering introduces randomness in brightness, contrast, saturation, and hue to simulate diverse lighting conditions, enhancing the robustness of learned features. We compare models pre-trained with and without color-jitter augmentation to investigate its impact.

**Results** We use the UniGaze-B model as the backbone and compare different pre-training settings, followed by training on XGaze and testing on the remaining four datasets. Table 3 demonstrates that both color-jitter augmentation and pixel normalization contribute to improved gaze estimation performance, highlighting their benefits for the generalization of the pre-trained model. Notably, pixel normalization consistently improves performance across all test datasets, aligning with the observations in the original MAE paper [29], which showed that pixel normalization enhances representation learning.

## C. Comparison with the SOTAs

3DGazeNet [73] collects in-the-wild face images with pseudo gaze labels and applies multi-view synthesis to obtain an augmented dataset ITWG-MV. To account for the difference in test data settings, we compare 3DGazeNet with UniGaze-H separately in Tab. 4. The results demonstrate that UniGaze-H outperforms 3DGazeNet in all domain generalization settings.

**Re-implementation** In the main paper, we compared our UniGaze-H model with state-of-the-art (SOTA) methods using their reported results. It is important to note that minor discrepancies may arise due to differences in our data pre-processing compared to prior work [13, 83, 93].

| Models | X→M | X→GC | G360→M | G360→GC |
|---|---|---|---|---|
| 3DGazeNet[†] [73] | 6.0 | 7.8 | 6.3 | 8.0 |
| UniGaze-H | **5.57** | **6.56** | **5.43** | **6.48** |

Table 4. Domain generalization compared with SOTA methods. The results marked with [†] are directly cited from previous studies [73].

| Models | X→M | X→E$_{CS}$ | G360→M | G360→E$_{CS}$ |
|---|---|---|---|---|
| ResNet-18 | 7.57 | 9.54 | 9.24 | 8.07 |
| ResNet-18[†] [93] | 8.02 | 9.11 | 8.04 | 9.20 |
| PureGaze | 6.68 | 7.62 | 8.87 | 10.53 |
| PureGaze[†] [13] | 7.08 | 7.48 | 9.28 | 9.32 |
| UniGaze-H | **5.57** | **4.65** | **5.43** | **5.35** |

Table 5. Domain generalization compared with SOTA methods and their re-implementations. The results marked with [†] are cited from previous studies [13, 93], and the rest of the results are based on our implementation.

To ensure a fair comparison, we re-implemented ResNet-18 and PureGaze [13] using our pre-processed datasets, aligning them with the reported results [13, 93]. The re-implementation results, alongside the reported values, are summarized in Tab. 5.

While minor differences exist between our re-implementation and the reported values, the improvements achieved by our UniGaze-H model remain significant, demonstrating its superior performance across all domain generalization tasks.

## D. Implementation Details

**Novel-Rendered Data Preparation** To render images from novel views, we follow the rendering approach described in [61]. To control the head pose, we randomly generate target head poses and compute the corresponding rotation matrices to apply to the 3D face models. During the rendering process, 40% of the images are assigned a random background color, while the remaining 60% use random scene images from the Places365 dataset [97] as background. Additionally, to simulate varied lighting conditions, half of the rendered images are adjusted to have lower ambient light intensity, ranging from 0.2 to 0.75.

All face images in our method are in the size of $224 \times 224$ after the data normalization process [90]. When the camera parameters are unknown, we use a camera matrix with focal length $f$ set to the image width and principal point $(c_x, c_y)$ set to half the image height and width.

**Pre-Training** We apply random color jitter augmentation with a probability of 0.5 and the following parameters: hue in the range $[-0.15, 0.15]$, saturation in $[0.8, 1.2]$, contrast in $[0.4, 1.8]$, and brightness in $[0.7, 1.3]$. We apply random grayscale with a probability of 0.05 on all images.

**Gaze Estimation Training** We use the Adam optimizer [42] with a learning rate of $1 \times 10^{-4}$ and a weight decay of $1 \times 10^{-6}$ for all experiments. For experiments with ResNet-50 and GazeTR-50, we set the batch size to 128 and decay the learning rate by 0.1 every five epochs, with a total of 12 epochs. For cross-dataset evaluation with UniGaze-H, we use a batch size of 128 and train the model for eight epochs with the one-cycle learning rate schedule [68]. For *leave-one-dataset-out* and *joint-dataset* evaluations, we set the batch size to 160 with 12 epochs.

## E. Qualitative Results

In this section, we present additional qualitative results using the UniGaze-H model trained on the aggregated datasets under the *joint-dataset* setting. We employ an off-the-shelf facial landmark detector [7] to extract landmarks and perform data normalization. Gaze estimation is conducted on the normalized images, and the results are de-normalized back to the original image for visualization. For reference, we also include the normalized faces alongside the original images.

Figure 3 and Fig. 4 showcase examples from various in-the-wild videos captured under challenging conditions, including large head poses and diverse lighting environments. Notably, we also include a synthetic example from URAvatar [49] (bottom row in Fig. 4), which generates faces with controlled viewpoints and lighting. Furthermore, Fig. 5 presents examples from the gaze-following dataset VideoAttentionTarget [15], a collection of diverse samples extracted from movies. This dataset provides annotated gaze targets, which are visualized when annotated within the image frame, as some targets may be out of frame.

These examples highlight the model's ability to predict gaze direction accurately in unseen environments, even under extreme head poses, challenging lighting conditions, and synthetic appearances.
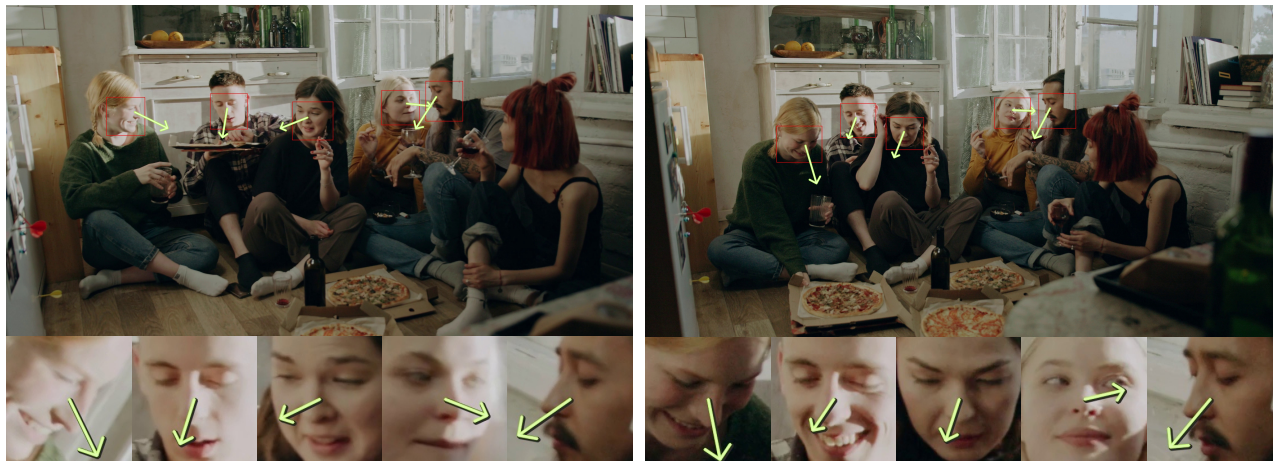
## F. Ethical Considerations

Our research involves the use of existing facial and gaze datasets. In accordance with ethical guidelines, we rely on the fact that these datasets were originally collected and published following relevant ethical and data protection standards, including obtaining consent, and we do not generate or collect additional new data. Our experimental protocols involve only image content, with no identifiable personal information or links to other personal data.

Figure 3. Qualitative results from various in-the-wild video examples. The normalized input images are displayed alongside the original image for reference.

(a) Qualitative results of in-the-wild video examples from public domain.



(b) Qualitative results of synthetic faces from URAvatar [49] .

Figure 4. Qualitative results of in-the-wild video and synthetic video. The normalized input images are displayed alongside the original image for reference.

Figure 5. Qualitative results of examples from the VideoAttentionTarget dataset [15]. Gaze targets are visualized when annotated within the image frame, as some targets may be out of frame. The normalized input images are displayed alongside the original image for reference.