# CUDA-GHR: Controllable Unsupervised Domain Adaptation for Gaze and Head Redirection

Swati Jindal, Xin Eric Wang

University of California, Santa Cruz

{swjindal, xwang366}@ucsc.edu

## Abstract

*The robustness of gaze and head pose estimation models is highly dependent on the amount of labeled data. Recently, generative modeling has shown excellent results in generating photo-realistic images, which can alleviate the need for annotations. However, adopting such generative models to new domains while maintaining their ability to provide fine-grained control over different image attributes, e.g., gaze and head pose directions, has been a challenging problem. This paper proposes CUDA-GHR, an unsupervised domain adaptation framework that enables fine-grained control over gaze and head pose directions while preserving the appearance-related factors of the person. Our framework simultaneously learns to adapt to new domains and disentangle visual attributes such as appearance, gaze direction, and head orientation by utilizing a label-rich source domain and an unlabeled target domain. Extensive experiments on the benchmarking datasets show that the proposed method can outperform state-of-the-art techniques on both quantitative and qualitative evaluations. Furthermore, we demonstrate the effectiveness of generated image-label pairs in the target domain for pretraining networks for the downstream task of gaze and head pose estimation. The source code and pre-trained models are available at* https://github.com/jswati31/cuda-ghr.

## 1. Introduction

Gaze behavior plays a pivotal role in the analysis of non-verbal cues and can provide support to various applications such as virtual reality [40, 41], human-computer interaction [34, 23], cognition [1, 43], and social sciences [21, 37]. Recent gaze estimation models rely on learning robust representations requiring a time-consuming and expensive step of collecting a large amount of training data, especially when labels are continuous. Although various methods [55, 53, 44] have been proposed to circumvent the data need, to generalize in-the-wild real-world scenarios remains
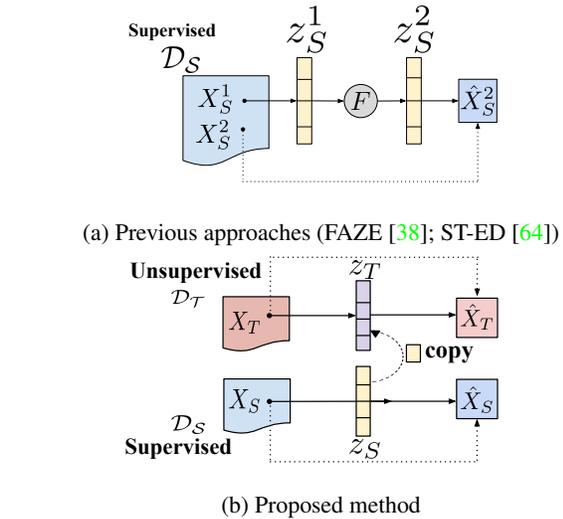


(a) Previous approaches (FAZE [38]; ST-ED [64])

(b) Proposed method

Figure 1: **Comparison of existing methods and proposed method**. In Fig (a), previous approaches [38, 64] assume conditional image-to-image translation ($X_S^1 \rightarrow X_S^2$) using a pair of labeled samples from a single domain $\mathcal{D}_S$ and use a transforming function $F$ in the latent space to ensure disentanglement. Here, $\mathcal{D}_S$ and $\mathcal{D}_T$ represent the source and target domains. In Fig (b), our method auto-encodes the images $X_S$, $X_T$ from both domains into a common disentangled space using labels only from source, and transfers latent factors via a simple copy operation.

a challenge and is an open research problem.

Different gaze redirection methods [64, 57, 25] have been explored as an alternate solution for generating more labeled training data using generative adversarial networks (GANs) [10] based frameworks. These generative methods require a pair of labeled images across both source and target domains to learn image-to-image translation; thus, these methods fail to generalize faithfully to new domains. Furthermore, various visual attributes are entangled during the generation process and cannot be manipulated independently to provide fine-grained control. Consequently, these

methods have limited applicability, as in order for the generated data to be useful on downstream tasks, the variability of these visual attributes across the generated data plays a key role in their success. Few works [51, 28] on neural image generation attempt to manipulate individual visual attributes in-the-wild real-world scenarios; however, they are constrained by the availability of simulated data with pre-defined labeled attributes. The recent work [52] proposes contrastive regression loss and utilizes unsupervised domain adaption to improve gaze estimation performance.

In this paper, we propose a novel domain adaptation framework for the task of controllable generation of eye gaze and head pose directions in the target domain while not requiring any label information in the target domain. Our method learns to render such control by disentangling explicit factors (*e.g.*, gaze and head orientations) from various implicit factors (*e.g.*, appearance, illumination, shadows, etc.) using a labeled-rich source domain and an unlabeled target domain. Both disentanglement and domain adaptation are performed jointly, thus enabling the transfer of learned knowledge from the source to the target domain. Since we use only unlabeled target-domain data to train our framework, we call it as *unsupervised domain adaptation* [65, 48].

Figure 1 illustrates the differences between the proposed method and previous approaches [38, 64]. Previous approaches use a pair of labeled samples $(X_S^1, X_S^2)$ from the source domain $\mathcal{D}_S$ to learn the conditional image-to-image translation while disentangling visual attributes using a transforming function $F$. In particular, Park *et al*. [38] provides control over only explicit factors while Zheng *et al*. [64] manipulate both explicit and implicit visual attributes. In contrast, our method can perform controllable generation without any input-output paired samples and apply auto-encoding of images $X_S$ and $X_T$ from source $\mathcal{D}_S$ and target $\mathcal{D}_T$ domains into a common disentangled latent space. Concurrently, we adapt the latent representations from the two domains, thereby allowing the transfer of learned knowledge from the labeled source to the unlabeled target domain. Unlike previous approaches, the proposed method is less constrained by label information and can be seamlessly applied to a broader set of datasets/applications.

We train our method on GazeCapture [29] dataset and demonstrate its efficacy on two target domains: MPI-IGaze [62] and Columbia [46] and obtain improved qualitative and quantitative results over state-of-the-art methods [38, 64]. Our experimental results exhibit a higher quality in preserving photo-realism of the generated images while faithfully rendering the desired gaze direction and head pose orientation. Overall, our contributions can be summarized as follows:

1. We propose a domain adaptation framework for jointly learning disentanglement and domain adaptation in latent space, using labels only from the source domain.

2. Our method utilizes auto-encoding behavior to maintain implicit factors and enable fine-grained control over gaze and head pose directions and outperforms the baseline methods on various evaluation metrics.

3. We demonstrate the effectiveness of generated redirected images in improving the downstream task performance on gaze and head pose estimation.

## 2. Related Work

This section provides a brief overview of the works on learning disentangled representations and gaze redirection methods.

### 2.1. Disentangled Representations

The goal of learning disentangled representations is to model the variability of implicit and explicit factors prevalent in the data generating process [35]. Fully supervised methods [42, 59, 6] exploit the semantic knowledge gained from the available annotations to learn these disentangled representations. On the other hand, unsupervised methods [13, 3] aim to learn the same behavior without relying on any labeled information. However, these methods provide limited flexibility to choose a specific factor of variation and are predominantly focused on a single domain representation learning problems [4].

Unsupervised cross-domain disentangled representation learning methods [32, 18] exploit the advantage of domain-shared and domain-specific attributes in order to provide fine-grained control on the appearance and content of the image. For instance, synthetic data is utilized by a few recent works [51, 28] to control various visual attributes while relying on the pre-defined label information associated with the rendered image obtained through a graphics pipeline. On the other hand, Liu *et al*. [33] provide control over different image attributes using the images from both source and target domains and is trained in a semi-supervised setting. However, their approach only considers categorical labels and thus has limited applicability. In contrast, our method allows controllable manipulation of continuous-valued image attributes (i.e., gaze and head pose) in the cross-domain setting.

### 2.2. Gaze Redirection Methods

Numerous methods have been developed for gaze redirection to attain a large amount of labeled synthetic data for the gaze estimation task. Kononenko *et al*. [27] use random forests to predict the flow field for gaze correction. More recently, several works [9, 2, 58] employ a deep neural network to learn the warping flow field between images along with a correction term. However, these warping-based methods cannot generalize well to large gaze and
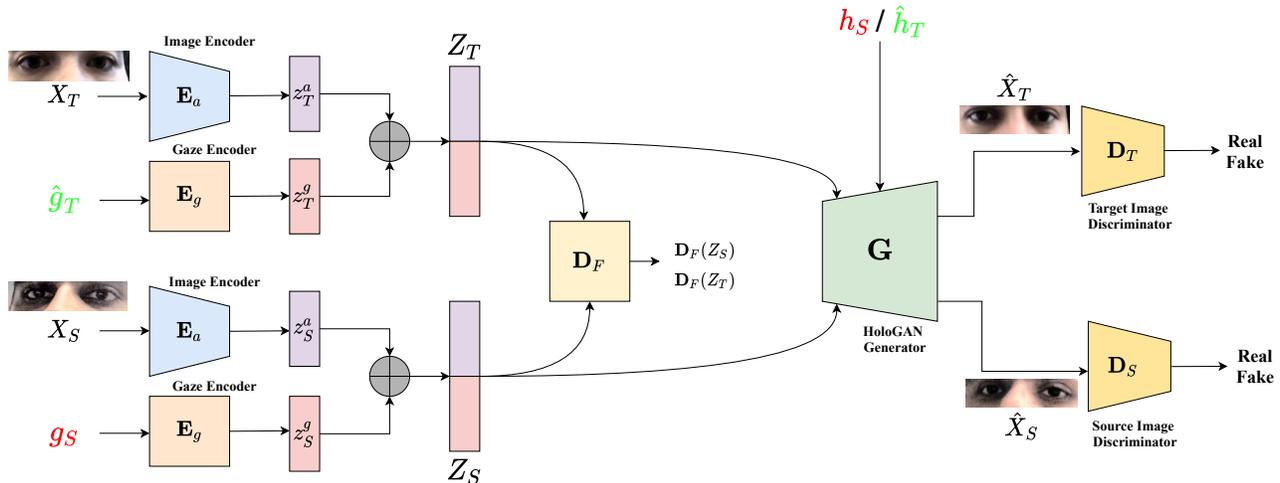
Figure 2: **Overview of CUDA-GHR.** The framework consists of two encoders $\mathbf{E}_a$ and $\mathbf{E}_g$ shared by both source and target domains. $\mathbf{E}_a$ encodes the target domain image $X_T$ to $z_T^a$, and the source domain image $X_S$ to $z_S^a$ while $\mathbf{E}_g$ encodes the target pseudo gaze label $\hat{g}_T$ and ground-truth source gaze label $g_S$ to $z_T^g$ and $z_S^g$, respectively. The overall image representations are formed as $Z_S = z_S^a \oplus z_S^g$ and $Z_T = z_T^a \oplus z_T^g$ (where, $\oplus$ is concatenate operation). These domain-specific encoded embeddings $Z_T$ and $Z_S$ are passed through a shared generator network $\mathbf{G}$ along with the corresponding head poses (pseudo head pose label $\hat{h}_T$ for the target domain, and ground-truth head pose label $h_S$ for source domain). These embeddings are also passed through a feature domain discriminator $\mathbf{D}_F$. $\mathbf{D}_T$ and $\mathbf{D}_S$ represent two domain-specific image discriminators. The whole framework is learned in an end-to-end manner. The labels in red color are the ground-truth labels while in green color are the generated pseudo-labels.

head pose directions. He *et al*. [12] propose a GAN-based framework that utilizes a cycle consistency loss to learn gaze redirection and generate images with a high resolution.

In addition, Wood *et al*. [54] uses a graphics pipeline to redirect eye images by fitting morphable models. However, these modeling-based methods make assumptions that do not hold in practice. Mask-based generator networks [39] have also been explored for the task of gaze redirection, though their performance is highly dependent on the accuracy of the segmentation module [25]. Park *et al*. [38] utilize a transforming encoder-decoder based network [14, 56] to learn disentanglement in the latent space. Recently, Xia *et al*. [57] and Zheng *et al*. [64] proposed controllable gaze redirection method using conditional image-to-image translation. However, these methods use a pair of labeled samples during training. As mentioned earlier, our method does not require any paired input-output samples and can be adapted to the target domain without any label data.

## 3. Proposed Method

Our goal is to learn a controller network $\mathbf{C}$ such that given an input image $X_T$ and subset of explicit factors $\{e_i\}$ (*e.g.*, gaze and head pose directions), it generates an image $X_O$ satisfying the attributes described by $\{e_i\}$, i.e., $C : (X_T, e_i) \rightarrow X_O$. To achieve this, we design a framework that learns to disentangle the latent space and

manipulate each explicit factor independently. We start with the assumption that there are three factors of variations: 1) appearance-related, including illumination, shadows, person-specific, etc., which might or might not be explicitly labeled with the dataset, 2) eye gaze direction and 3) head pose orientation. We train our network in an unsupervised domain adaptation setting by utilizing a fully labeled source domain and an unlabeled target domain considering distribution shift across datasets into account. Recall that we have the gaze and head pose labels only for the source domain. Therefore, we aim to disentangle and control these three factors of variations in the latent space and simultaneously transfer the learned behavior to the unsupervised target domain. We named our framework as CUDA-GHR.

### 3.1. Model

The overall architecture of the CUDA-GHR is shown in Figure 2. We denote $S$ as the source domain and $T$ as the target domain. Further, following the notations used in [38], we represent the appearance-related latent factor as $z^a$ and gaze latent factor as $z^g$.

The initial stage of our network consists of two encoders: (a) an image encoder $\mathbf{E}_a$ encodes the implicit (appearance-related) factors of an image $X_i$ and outputs $z_i^a$ such that $i \in \{S, T\}$, and (b) a separate MLP-based gaze encoder $\mathbf{E}_g$ encodes the input gaze $g_i$ corresponding to the image $X_i$ to

a latent factor $z_i^g$. For the source domain, we use ground-truth gaze label $g_S$ as input to $\mathbf{E}_g$ while for the unlabeled target domain, we input pseudo gaze labels $\hat{g}_T$ obtained from a pre-trained task network $\mathcal{T}$ that predicts gaze and head pose of an image. Note that $\mathcal{T}$ is trained only on source domain data. Thus, the overall embedding $Z_i$ related to an image $X_i$ can be formed by concatenating these two latent factors, i.e., $Z_i = z_i^a \oplus z_i^g$ (here $\oplus$ denotes concatenation). Further, $Z_i$ and head pose label $h_i$ are given as input to a decoder $\mathbf{G}$ based on the generator used in HoloGAN [36] as it allows the head pose to be separately controlled without any encoder. This generator $\mathbf{G}$ decodes the latent $Z_i$ and head pose $h_i$ to an output image given by $\hat{X}_i$ and is trained in an adversarial manner with the discriminator network $\mathbf{D}_i$. Note again that for labeled source images, we use ground-truth head pose label $h_S$ while we take pseudo head pose label $\hat{h}_T$ produced by task network $\mathcal{T}$ for unlabeled target domain inputs. In addition, we use a feature domain discriminator $\mathbf{D}_F$ to ensure that the latent distributions of $Z_S$ and $Z_T$ are similar.

At inference time, the gaze and head pose directions are controlled by passing an image from the target domain $X_T$ through the encoder $\mathbf{E}_a$ and desired gaze direction $g$ through $\mathbf{E}_g$, giving us $\mathbf{E}_a(X_T)$ and $\mathbf{E}_g(g)$ respectively. These two latent factors are concatenated and passed through the generator $\mathbf{G}$ along with the desired head pose $h$ to generate an output image $\hat{X}_T^{g,h}$ with gaze $g$ and head pose $h$, i.e.,

$$\hat{X}_T^{g,h} = \mathbf{G}(\mathbf{E}_a(X_T) \oplus \mathbf{E}_g(g), \ h) \qquad (1)$$

Likewise, we can also control the individual factor of gaze (or head pose) by providing desired gaze (or head pose) direction and pseudo head pose (or gaze) label obtained from $\mathcal{T}$ to generate gaze redirected image given as

$$\hat{X}_T^g = \mathbf{G}(\mathbf{E}_a(X_T) \oplus \mathbf{E}_g(g), \ \hat{h}_T) \qquad (2)$$

and head redirected image given as

$$\hat{X}_T^h = \mathbf{G}(\mathbf{E}_a(X_T) \oplus \mathbf{E}_g(\hat{g}_T), \ h) \qquad (3)$$

### 3.2. Learning Objectives

The overall objective of our method is to learn a common factorized latent space for both source and target domain such that the individual latent factors can be easily controlled to manipulate target images. To ensure this, we train our framework using multiple objective functions, each of which are explained in detail below.

**Reconstruction Loss.** We apply pixel-wise L1 reconstruction loss between the generated image $\hat{X}_i$ and input image $X_i$ to ensure the auto-encoding behavior.

$$\mathcal{L}_{\mathcal{R}}(\hat{X}_i, X_i) = \frac{1}{|X_i|}||\hat{X}_i - X_i||_1 \qquad (4)$$

Thus, the total reconstruction loss is defined as

$$\mathcal{L}_{recon} = \sum_{i \in \{S,T\}} \mathcal{L}_{\mathcal{R}}(\hat{X}_i, X_i) \qquad (5)$$

**Perceptual Loss.** To ensure that our generated images perceptually match the input images, we apply the perceptual loss [24] which is defined as a mean-square loss between the activations of a pre-trained neural network applied between the generated image $\hat{X}_i$ and input image $X_i$. For this, we use VGG-16 [45] network trained on ImageNet [31].

$$\mathcal{L}_{\mathcal{P}}(\hat{X}_i, X_i) = \sum_{l=1}^{4} \frac{1}{|\psi_l(X_i)|} ||\psi_l(\hat{X}_i) - \psi_l(X_i)||_2 \qquad (6)$$

where $\psi$ denotes VGG-16 network. Therefore, our overall perceptual loss becomes

$$\mathcal{L}_{perc} = \sum_{i \in \{S,T\}} \mathcal{L}_{\mathcal{P}}(\hat{X}_i, X_i) \qquad (7)$$

**Consistency Loss.** To ensure disentangled behavior between implicit and explicit factors, we apply a consistency loss between the generated image $\hat{X}_i$ and input image $X_i$. For this, we use a pre-trained task network $\mathcal{T}$ which predicts the pseudo-labels (gaze and head pose) for an image. The consistency loss consists of two terms: (a) *label consistency loss* is applied between pseudo-labels for input and the generated images to preserve the gaze and head pose information, and (b) *redirection consistency loss* guarantees to preserve the pseudo-labels for redirected images. For (b), we generate gaze and head redirected images using Equation 2 and 3 respectively, by applying gaze and head pose labels from source domain. We enforce the gaze prediction consistency between $\hat{X}_T^g$ and $X_S$, and head pose prediction consistency between $\hat{X}_T^g$ and $X_T$, i.e., $\mathcal{T}^g(\hat{X}_T^g) = \mathcal{T}^g(X_S)$ and $\mathcal{T}^h(\hat{X}_T^g) = \mathcal{T}^h(X_T)$. A similar argument holds for the head redirected image $\hat{X}_T^h$, i.e., $\mathcal{T}^g(\hat{X}_T^h) = \mathcal{T}^g(X_T)$ and $\mathcal{T}^h(\hat{X}_T^h) = \mathcal{T}^h(X_S)$. Here, $\mathcal{T}^g$ and $\mathcal{T}^h$ represent the gaze and head pose predicting layers of $\mathcal{T}$. The overall gaze consistency loss will become

$$\mathcal{L}_{gc} = \underbrace{\mathcal{L}_a(\mathcal{T}^g(\hat{X}_S), \mathcal{T}^g(X_S)) + \mathcal{L}_a(\mathcal{T}^g(\hat{X}_T), \mathcal{T}^g(X_T))}_{label\ consistency\ loss}$$
$$+ \underbrace{\mathcal{L}_a(\mathcal{T}^g(\hat{X}_T^g), \mathcal{T}^g(X_S)) + \mathcal{L}_a(\mathcal{T}^g(\hat{X}_T^h), \mathcal{T}^g(X_T))}_{redirection\ consistency\ loss}$$
$$(8)$$

Similarly, we can compute the head pose consistency

loss $\mathcal{L}_{hc}$ as follows:

$$\mathcal{L}_{hc} = \underbrace{\mathcal{L}_a(\mathcal{T}^h(\hat{X}_S), \mathcal{T}^h(X_S)) + \mathcal{L}_a(\mathcal{T}^h(\hat{X}_T), \mathcal{T}^h(X_T))}_{\text{label consistency loss}}$$
$$+ \underbrace{\mathcal{L}_a(\mathcal{T}^h(\hat{X}_T^g), \mathcal{T}^h(X_T)) + \mathcal{L}_a(\mathcal{T}^h(\hat{X}_T^h), \mathcal{T}^h(X_S))}_{\text{redirection consistency loss}} \tag{9}$$

Here, $\mathcal{L}_a$ is defined as:

$$\mathcal{L}_a(\hat{\boldsymbol{u}}, \boldsymbol{u}) = \arccos\left(\frac{\hat{\boldsymbol{u}} \cdot \boldsymbol{u}}{||\hat{\boldsymbol{u}}|| \cdot ||\boldsymbol{u}||}\right) \tag{10}$$

Hence, total consistency loss becomes

$$\mathcal{L}_{consistency} = \mathcal{L}_{gc} + \mathcal{L}_{hc} \tag{11}$$

**GAN Loss.** To enforce photo-realistic output from the generator $\mathbf{G}$, we apply the standard GAN loss [10] to image discriminator $\mathbf{D}_i$.

$$\mathcal{L}_{GAN_D}(\mathbf{D}_i, X_i, \hat{X}_i) = \log \mathbf{D}_i(X_i) + \log(1 - \mathbf{D}_i(\hat{X}_i))$$
$$\mathcal{L}_{GAN_G}(\mathbf{D}_i, \hat{X}_i) = \log \mathbf{D}_i(\hat{X}_i) \tag{12}$$

The final GAN loss is defined as

$$\mathcal{L}_{disc} = \sum_{i \in \{S,T\}} \mathcal{L}_{GAN_D}(\mathbf{D}_i, X_i, \hat{X}_i)$$
$$\mathcal{L}_{gen} = \sum_{i \in \{S,T\}} \mathcal{L}_{GAN_G}(\mathbf{D}_i, \hat{X}_i) \tag{13}$$

**Feature Domain Adversarial Loss.** We employ a latent domain discriminator network $\mathbf{D}_F$ and train it using the following domain adversarial loss [49] to push the distribution of $Z_T$ closer to $Z_S$.

$$\mathcal{L}_{feat}(\mathbf{D}_F, Z_T, Z_S) = \log \mathbf{D}_F(Z_S) + \log(1 - \mathbf{D}_F(Z_T)) \tag{14}$$

**Overall Loss.** Altogether, our final loss function for training encoders and generator network is

$$\mathcal{L}_{overall} = \lambda_R \mathcal{L}_{recon} + \lambda_P \mathcal{L}_{perc} + \lambda_C \mathcal{L}_{consistency} + \lambda_G \mathcal{L}_{gen} + \lambda_F \mathcal{L}_{feat} \tag{15}$$

# 4. Experiments

## 4.1. Datasets

**GazeCapture** [29] is the largest publicly available gaze dataset consisting of around 2M frames taken from unique 1474 subjects. Following the split defined in [29], we use data from 1274 subjects for training, 50 for validation, and 150 for the test.

**MPIIGaze** [62] is the most challenging dataset for the in-the-wild gaze estimation and includes higher within-subject variations in appearance, for example, illumination, make-up, and facial hair. We use the images from the standard evaluation subset MPIIFaceGaze [63] provided by MPI-IGaze containing 37667 images captured from 15 subjects. **Columbia** [46] contains 5880 high-resolution images from 56 subjects and displays larger diversity within participants. The images are collected in a constrained laboratory setting, with limited variations of head pose and gaze directions.

## 4.2. Implementation Details

The architecture of the encoder $\mathbf{E}_a$ is DenseNet-based blocks as used in Park *et al*. [38] and the decoder network $\mathbf{G}$ is HoloGAN based generator [36]. The gaze encoder $\mathbf{E}_g$ consists of four MLP layers with hidden dimensions equal to the input dimension and output dimension is set to 8. The task network $\mathcal{T}$ is a ResNet-50 [11] based model trained on GazeCapture [29] training subset and gives 4-D output, two angles for each gaze and head direction. The two image discriminators $\mathbf{D}_S$ and $\mathbf{D}_T$ share a similar PatchGAN [22] based architecture. The domain discriminator $\mathbf{D}_F$ consists of four MLP layers. Note that $\mathcal{T}$ remains fixed during training of our whole pipeline. More implementation details can be found in the supplementary materials.

All the datasets are pre-processed by a data normalization algorithm as described in Zhang *et al*. [61]. Our input is a single image containing both eyes and is of size $256 \times 64$. We use a data processing pipeline as employed in Park *et al*. [38] to extract the eye image strip. The inputs gaze $g$ and head pose $h$ are 2-D pitch and yaw angles. We train our framework in two settings: *GazeCapture→MPIIGaze*, trained with GazeCapture as source domain and MPIIGaze as target domain, and *GazeCapture→Columbia* is trained with Columbia as the target domain. For GazeCapture, we use the training subset from the data split as labeled source domain data. From MPIIGaze and Columbia, we respectively choose the first 11 and 50 subjects as unlabeled target domain data for training. We call them as **'Seen'** subjects as our network sees them during training while remaining users fall into **'Unseen'** category. We evaluate our method on three test subsets: 'Unseen', 'Seen' and 'All'. 'All' includes both 'Seen' and 'Unseen' participants data.

**Hyper-parameters.** We use a batch size of 10 for both *GazeCapture→MPIIGaze* and *GazeCapture→Columbia* and are trained for 200K and 81K iterations, respectively. All network modules are optimized through Adam [26] optimizer with a weight decay coefficient of $10^{-4}$. The initial learning rate is set to $0.0005$ which is decayed by a factor of $0.8$ after approximately 34K iterations. For *GazeCapture→MPIIGaze*, we restart the learning rate scheduler after around 160K iterations for better conver-

Table 1: **Quantitative Evaluation.** Comparison of CUDA-GHR with the state-of-the-art methods [38, 64]. *GazeCapture→MPIIGaze* is evaluated on MPIIGaze subsets and *GazeCapture→Columbia* is evaluate on Columbia subsets. All errors are in degrees (°) except LPIPS, and lower is better.

| Test Set | Method | GazeCapture→MPIIGaze | | | | | GazeCapture→Columbia | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | LPIPS ↓ | Gaze Redir. ↓ | Head Redir. ↓ | $g \rightarrow h \downarrow$ | $h \rightarrow g \downarrow$ | LPIPS ↓ | Gaze Redir. ↓ | Head Redir. ↓ | $g \rightarrow h \downarrow$ | $h \rightarrow g \downarrow$ |
| Unseen | FAZE | 0.311 | 6.131 | 6.408 | 6.925 | 4.909 | 0.435 | 9.008 | 6.996 | 6.454 | 4.295 |
| | ST-ED | 0.274 | 2.355 | 1.605 | 1.349 | 2.455 | 0.265 | 2.283 | 1.651 | 1.364 | 2.190 |
| | ST-ED+PS | 0.266 | 2.864 | 1.576 | 1.472 | 2.346 | 0.266 | 2.117 | 1.437 | **1.124** | 2.356 |
| | CUDA-GHR | **0.261** | **2.023** | **1.154** | **1.161** | **1.829** | **0.255** | **1.449** | **0.873** | 1.209 | **1.514** |
| Seen | FAZE | 0.382 | 5.778 | 6.899 | 5.311 | 5.172 | 0.486 | 10.368 | 7.231 | 7.302 | 4.788 |
| | ST-ED | 0.315 | 2.405 | 1.669 | 1.209 | 2.341 | 0.319 | 2.484 | 1.616 | 1.343 | 2.456 |
| | ST-ED+PS | 0.288 | 2.269 | 1.888 | 1.179 | 2.229 | 0.299 | 2.071 | 1.536 | 1.088 | 2.330 |
| | CUDA-GHR | **0.278** | **1.905** | **0.979** | **0.761** | **1.236** | **0.282** | **1.328** | **0.831** | **0.646** | **0.996** |
| All | FAZE | 0.370 | 5.840 | 6.828 | 5.613 | 5.123 | 0.481 | 10.214 | 7.226 | 7.214 | 4.737 |
| | ST-ED | 0.307 | 2.392 | 1.660 | 1.232 | 2.359 | 0.314 | 2.473 | 1.618 | 1.350 | 2.435 |
| | CUDA-GHR | **0.275** | **1.922** | **1.012** | **0.844** | **1.341** | **0.279** | **1.337** | **0.832** | **0.707** | **1.045** |

gence. The weights of the objective function are set as $\lambda_R$ = 200, $\lambda_P$ = 10, $\lambda_C$ = 10, $\lambda_G$ = 5 and $\lambda_F$ = 5.

### 4.3. Evaluation Metrics

We evaluate our framework using three evaluation metrics as previously adopted by [64]: perceptual similarity, redirection errors, and disentanglement errors.

**Learned Perceptual Image Patch Similarity (LPIPS)** [60] is used to measure the pairwise image similarity by calculating the distance in AlexNet [30] feature space.
**Redirection Errors** are computed as angular errors between the estimated direction obtained from our task network $\mathcal{T}$ and the desired direction. It measures the accomplishment of the explicit factors, i.e., gaze and head directions in the image output.
**Disentanglement Error** measures the disentanglement of explicit factors like gaze and head pose. We evaluate $g \rightarrow h$, the effect of change in gaze direction on the head pose, and vice versa ($h \rightarrow g$). To compute $g \rightarrow h$, we first calculate the joint probability distribution function of the gaze direction values from the source domain and sample random gaze labels. We apply this gaze direction to the input image while keeping the head pose unchanged and measure the angular error between head pose predictions from task network $\mathcal{T}$ of the redirected image and the original reconstructed image. Similarly, we compute $h \rightarrow g$ by sampling random head pose orientations from the source labeled data.

### 4.4. Comparison to the state-of-the-art

We adopt FAZE [38] and ST-ED [64] as our baseline methods. Both FAZE and ST-ED are based on transforming encoder-decoder architecture [14, 56] and apply known differences in gaze and head rotations to the embedding space

for translating the input image to a redirected output image. FAZE inputs an image containing both eyes, which is the same as our method, thus necessary to compare. We use original implementation[1] and trained models provided by the FAZE authors for comparison. In addition, we retrain the ST-ED network on images containing both eyes for a fair comparison. FAZE learns to control only explicit factors (gaze and head pose orientations) while ST-ED controls implicit factors too. Note that for the ST-ED baseline, we compare only by altering explicit factors. Furthermore, we also compare CUDA-GHR to baseline ST-ED+PS which is trained with source data GazeCapture and using pseudo-labels for target dataset (MPIIGaze or Columbia). The pseudo-labels are obtained in same manner as of CUDA-GHR. For more details, please refer to the supplementary materials.

**Quantitative Evaluation.** Table 1 summarizes the quantitative evaluation of both our experiments *GazeCapture→MPIIGaze* and *GazeCapture→Columbia*. The left half of Table 1 shows evaluation on MPIIGaze test subsets {'Seen', 'Unseen', 'All'}, and we observe that our method outperforms the baselines (even ST-ED+PS) on all the evaluation metrics for each test subset. We get lower LPIPS (even on 'Unseen' users), indicating the generation of better quality images while achieving the desired gaze and head directions attested by lower gaze and head redirection errors. We also obtain better disentanglement errors exhibiting that our method successfully controls each explicit factor individually. The improved performance on 'Unseen' users shows the superiority and generalizability of our method over baselines. We also notice improvements over ST-ED+PS baseline, exhibiting that domain adaptation is essential to achieve better performance.

---

[1]https://github.com/NVlabs/few_shot_gaze

We show evaluation of *GazeCapture→Columbia* experiment on right half of Table 1. Note that due to the small size of the Columbia dataset, we initialize the model for this experiment with the previously trained weights on *GazeCapture→MPIIGaze* for better convergence. Recall that we do not use any labels from the target domain dataset in any experiment. As shown in Table 1, our method is consistently better than other baselines on all evaluation metrics, showing the generalizability of our framework on different domains and thus, can be adapted to new datasets without the requirement of any labels.
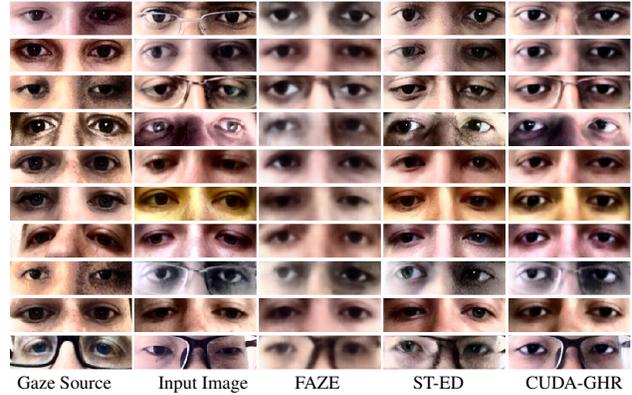
**Qualitative Evaluation.** We also report the qualitative comparison of generated images in Figure 3 using a model trained with *GazeCapture→MPIIGaze*. The results are shown on MPIIGaze dataset images which is the target domain dataset in this setting. As can be seen, our method produces better quality images while preserving the appearance information (*e.g.*, skin color, eye shape) and faithfully manipulating the gaze and head pose directions when compared with FAZE [38] and ST-ED [64]. It is also worth noting that our method generates higher-quality images for people with glasses, *e.g.*, row 3 in Figure 3a and row 6 in Figure 3b. These results are consistent with our findings in quantitative evaluation, thus showing that our method is more faithful in reproducing the desired gaze and head pose directions. Additional results are provided in the supplementary materials.
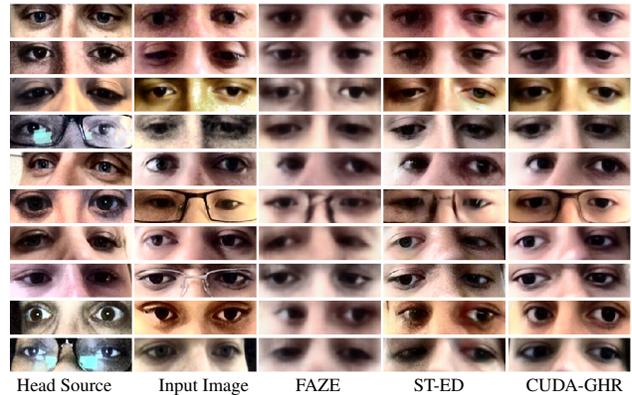
### 4.5. Ablation Study

To understand the role of individual components of the objective function, we provide following ablation study. In Table 2, we compare against the ablations of individual loss terms. The ablation on the perceptual loss is shown in the first row ($\lambda_P = 0$). The second row ($\lambda_C = 0$) represents when consistency loss is set to zero, while the third row ($\lambda_F = 0$) shows results when feature domain adversarial loss is not enforced during training. The fourth and fifth row shows an ablation on reconstruction ($\lambda_R = 0$) and GAN ($\lambda_G = 0$) loss, respectively. As can be seen, all of these loss terms are critical for the improvements in performance. We see a substantial improvement with the addition of $\mathcal{L}_{consistency}$. The ablation study is performed for *GazeCapture → MPIIGaze* on the 'Seen' subset of MPIIGaze.

### 4.6. Controllability

Figure 4 shows the effectiveness of our method in controlling the gaze and head pose directions. We vary pitch and yaw angles from $-30°$ to $+30°$ for gaze and head redirections. We can see that our method faithfully renders the desired gaze direction (or head pose orientation) while retaining the head pose (or gaze direction), therefore, exhibiting the efficacy of disentanglement. Furthermore, note that



(a) Gaze Redirected images



(b) Head Redirected images

Figure 3: **Qualitative Evaluation:** Comparison of the generated images from CUDA-GHR (*GazeCapture→MPIIGaze*) with the baseline methods FAZE [38] and ST-ED [64]. The quality of gaze redirected images is depicted in 3a, while head redirected images are shown in 3b. The first column represents the gaze/head pose source image from which gaze/head pose information is used to redirect. The second column shows the input image from the target domain. Our method (column 5) produces better quality images and preserves the implicit factors than the baseline methods (columns 3 and 4). Best viewed in color.

the range of yaw and pitch angles $[-30°, 30°]$ is out-of-label distribution of source dataset (GazeCapture), showing the extrapolation capability of CUDA-GHR in the generation process.

### 4.7. Evaluation of Downstream Tasks

We also demonstrate the utility of generated images from our framework in improving the performance of the downstream gaze and head pose estimation task. For this, we conduct experiments for cross-subject estimation on both MPIIGaze and Columbia datasets. The main goal of this

Table 2: **Ablation Study:** An ablation study on different loss terms for *GazeCapture → MPIIGaze* on MPIIGaze 'Seen' subset. All errors are in degrees (°) except LPIPS, and lower is better.

| Ablation term | LPIPS ↓ | Gaze Redir. ↓ | Head Redir. ↓ | $g \to h$ ↓ | $h \to g$ ↓ |
|---|---|---|---|---|---|
| $\lambda_P = 0$ | 0.307 | 6.450 | 0.922 | 0.655 | 1.334 |
| $\lambda_C = 0$ | 0.326 | 15.183 | 3.412 | **0.106** | 11.616 |
| $\lambda_F = 0$ | 0.281 | 4.791 | **0.787** | 0.636 | **0.826** |
| $\lambda_R = 0$ | 0.304 | 4.958 | 0.911 | 0.463 | 0.876 |
| $\lambda_G = 0$ | 0.309 | 11.130 | 0.942 | 0.355 | 0.868 |
| Ours | **0.278** | **1.905** | 0.979 | 0.761 | 1.236 |



(a) Gaze redirected images with (pitch, yaw) $\in [-30°, 30°]$



(b) Head redirected images with (pitch, yaw) $\in [-30°, 30°]$

Figure 4: **Controllable Generation:** Illustration of controllable gaze and head redirection showing the effectiveness of disentanglement between various explicit factors.

experiment is to show that the generated "free" labeled data from our framework can be used to obtain a good pretrained model to further fine-tune on cross-subject estimation task. We compare it against three initializations: random, ImageNet [7], and pretrained model obtained using ST-ED [64] generated images.

We generate around 10K samples per user from MPI-IGaze dataset using *GazeCapture→MPIIGaze* trained generator and train a ResNet-50 [11] network (initialized with ImageNet pre-trained weights) with batch normalization [20] replaced by instance normalization [50] layers. Afterward, we fine-tune this network on MPIIGaze dataset using leave-one-subject-out cross-validation for both gaze and head pose estimation and report the mean angular error. A similar method is followed for ST-ED generated images. We compare the errors obtained from four initialization methods: random, ImageNet, ST-ED, and CUDA-GHR. Analogously, we train gaze and head pose estimation mod-

Table 3: **Downstream Task Evaluation:** Comparison of mean angular errors (*mean ± std* in degrees) for various initialization methods on gaze and head pose estimation task. Lower is better.

| Initialization Method | Head Pose Estimation Errors↓ | | Gaze Estimation Errors↓ | |
|---|---|---|---|---|
| | Columbia | MPIIGaze | Columbia | MPIIGaze |
| Random | 6.8 ± 1.2 | 6.7 ± 0.7 | 6.7 ± 0.7 | 6.7 ± 1.3 |
| ImageNet | 5.9 ± 1.3 | 5.7 ± 2.8 | 5.5 ± 0.1 | 5.7 ± 1.4 |
| ST-ED | 5.7 ± 1.1 | 5.1 ± 2.4 | 5.4 ± 0.4 | **5.5 ± 1.3** |
| CUDA-GHR | **5.3 ± 1.1** | **4.9 ± 2.5** | **5.1 ± 0.4** | 5.5 ± 1.4 |

els on generated images for Columbia data subjects (∼1.6K samples each) using *GazeCapture→Columbia* model and fine-tune on Columbia dataset using 4-fold cross-validation. The comparison of different initialization methods on two datasets is shown in Table 3.

It can be seen that the model trained with CUDA-GHR gives around 7% and 4% relative improvements over ST-ED initialization on Columbia and MPIIGaze, respectively, for the head pose estimation task. We also show results for the gaze estimation task in Table 3 giving a relative improvement of around 5.5% on the Columbia dataset while performing similar to the ST-ED baseline on MPIIGaze. We hypothesize that this is because the gaze and head pose label distribution of GazeCapture is closer to MPIIGaze distribution than Columbia [5] and thus, performs closely for both ST-ED and CUDA-GHR. This indicates that domain adaptation is more advantageous for the Columbia dataset. Hence, it shows the effectiveness of our method over baselines when performing domain adaptation across datasets with significant distribution shifts.

## 5. Conclusion

We present an unsupervised domain adaptation framework trained using cross-domain datasets for gaze and head redirection tasks. The proposed method takes advantage of both supervised source domain and unsupervised target domain to learn the disentangled factors of variations. Experimental results demonstrate the effectiveness of our model in generating photo-realistic images in multiple domains while truly adapting the desired gaze direction and head pose orientation. Because of removing the requirement of annotations in the target domain, the applicability of our work increases for new datasets where manual annotations are hard to collect. Our framework is relevant to various applications such as video conferencing, photo correction, and movie editing for redirecting gaze to establish eye contact with the viewer. It can also be extended to improve performances on the downstream task of gaze and head pose estimation.

# References

[1] Guy Thomas Buswell. How people look at pictures: a study of the psychology and perception in art. 1935. 1

[2] Jingjing Chen, Jichao Zhang, Enver Sangineto, Tao Chen, Jiayuan Fan, and Nicu Sebe. Coarse-to-fine gaze redirection with numerical and pictorial guidance. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3665–3674, 2021. 2

[3] Ricky TQ Chen, Xuechen Li, Roger Grosse, and David Duvenaud. Isolating sources of disentanglement in variational autoencoders. *arXiv preprint arXiv:1802.04942*, 2018. 2

[4] Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. *arXiv preprint arXiv:1606.03657*, 2016. 2

[5] Yihua Cheng, Haofei Wang, Yiwei Bao, and Feng Lu. Appearance-based gaze estimation with deep learning: A review and benchmark. *arXiv preprint arXiv:2104.12668*, 2021. 8

[6] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8789–8797, 2018. 2

[7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 8, 12

[8] Jiankang Deng, Yuxiang Zhou, Shiyang Cheng, and Stefanos Zaferiou. Cascade multi-view hourglass model for robust 3d face alignment. In *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*, pages 399–403. IEEE, 2018. 12

[9] Yaroslav Ganin, Daniil Kononenko, Diana Sungatullina, and Victor Lempitsky. Deepwarp: Photorealistic image resynthesis for gaze manipulation. In *European conference on computer vision*, pages 311–326. Springer, 2016. 2

[10] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *arXiv preprint arXiv:1406.2661*, 2014. 1, 5

[11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5, 8, 12

[12] Zhe He, Adrian Spurr, Xucong Zhang, and Otmar Hilliges. Photo-realistic monocular gaze redirection using generative adversarial networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6932–6941, 2019. 3

[13] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. 2016. 2

[14] Geoffrey E Hinton, Alex Krizhevsky, and Sida D Wang. Transforming auto-encoders. In *International conference on artificial neural networks*, pages 44–51. Springer, 2011. 3, 6

[15] Peiyun Hu and Deva Ramanan. Finding tiny faces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 951–959, 2017. 12

[16] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017. 12

[17] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1501–1510, 2017. 12

[18] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 172–189, 2018. 2

[19] Patrik Huber, Guosheng Hu, Rafael Tena, Pouria Mortazavian, P Koppen, William J Christmas, Matthias Ratsch, and Josef Kittler. A multiresolution 3d morphable face model and fitting framework. In *Proceedings of the 11th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications*, 2016. 12

[20] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015. 8, 12

[21] Ryo Ishii, Kazuhiro Otsuka, Shiro Kumano, and Junji Yamato. Prediction of who will be the next speaker and when using gaze behavior in multiparty meetings. *ACM Transactions on Interactive Intelligent Systems (TIIS)*, 6(1):1–31, 2016. 1

[22] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017. 5, 12

[23] Robert JK Jacob and Keith S Karn. Eye tracking in human-computer interaction and usability research: Ready to deliver the promises. In *The mind's eye*, pages 573–605. Elsevier, 2003. 1

[24] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016. 4

[25] Harsimran Kaur and Roberto Manduchi. Subject guided eye image synthesis with application to gaze redirection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 11–20, 2021. 1, 3

[26] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5, 12

[27] Daniil Kononenko and Victor Lempitsky. Learning to look up: Realtime monocular gaze correction using machine learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4667–4675, 2015. 2

[28] Marek Kowalski, Stephan J. Garbin, Virginia Estellers, Tadas Baltrušaitis, Matthew Johnson, and Jamie Shotton. Config: Controllable neural face image generation. In *European Conference on Computer Vision (ECCV)*, 2020. 2

[29] Kyle Krafka, Aditya Khosla, Petr Kellnhofer, Harini Kannan, Suchendra Bhandarkar, Wojciech Matusik, and Antonio Torralba. Eye tracking for everyone. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2176–2184, 2016. 2, 5

[30] Alex Krizhevsky. One weird trick for parallelizing convolutional neural networks. *arXiv preprint arXiv:1404.5997*, 2014. 6

[31] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25:1097–1105, 2012. 4

[32] Hsin-Ying Lee, Hung-Yu Tseng, Jia-Bin Huang, Maneesh Singh, and Ming-Hsuan Yang. Diverse image-to-image translation via disentangled representations. In *Proceedings of the European conference on computer vision (ECCV)*, pages 35–51, 2018. 2

[33] Yen-Cheng Liu, Yu-Ying Yeh, Tzu-Chien Fu, Sheng-De Wang, Wei-Chen Chiu, and Yu-Chiang Frank Wang. Detach and adapt: Learning cross-domain disentangled deep representation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8867–8876, 2018. 2

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

[35] Michael Mathieu, Junbo Zhao, Pablo Sprechmann, Aditya Ramesh, and Yann LeCun. Disentangling factors of variation in deep representations using adversarial training. *arXiv preprint arXiv:1611.03383*, 2016. 2

[36] Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yong-Liang Yang. Hologan: Unsupervised learning of 3d representations from natural images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7588–7597, 2019. 4, 5, 12

[37] Catharine Oertel, Kenneth A Funes Mora, Joakim Gustafson, and Jean-Marc Odobez. Deciphering the silent participant: On the use of audio-visual cues for the classification of listener categories in group discussions. In *Proceedings of the 2015 ACM on International Conference on Multimodal Interaction*, pages 107–114, 2015. 1

[38] Seonwook Park, Shalini De Mello, Pavlo Molchanov, Umar Iqbal, Otmar Hilliges, and Jan Kautz. Few-shot adaptive gaze estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9368–9377, 2019. 1, 2, 3, 5, 6, 7, 12

[39] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2337–2346, 2019. 3

[40] Anjul Patney, Joohwan Kim, Marco Salvi, Anton Kaplanyan, Chris Wyman, Nir Benty, Aaron Lefohn, and David Luebke. Perceptually-based foveated virtual reality. In *ACM SIGGRAPH 2016 Emerging Technologies*, pages 1–2. 2016. 1

[41] Thies Pfeiffer. Towards gaze interaction in immersive virtual reality: Evaluation of a monocular eye tracking set-up. In *Virtuelle und Erweiterte Realität-Fünfter Workshop der GI-Fachgruppe VR/AR*, 2008. 1

[42] Shengju Qian, Kwan-Yee Lin, Wayne Wu, Yangxiaokang Liu, Quan Wang, Fumin Shen, Chen Qian, and Ran He. Make a face: Towards arbitrary high fidelity face manipulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10033–10042, 2019. 2

[43] Constantin A Rothkopf, Dana H Ballard, and Mary M Hayhoe. Task and context determine where you look. *Journal of vision*, 7(14):16–16, 2007. 1

[44] Ashish Shrivastava, Tomas Pfister, Oncel Tuzel, Joshua Susskind, Wenda Wang, and Russell Webb. Learning from simulated and unsupervised images through adversarial training. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2107–2116, 2017. 1

[45] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 4

[46] Brian A Smith, Qi Yin, Steven K Feiner, and Shree K Nayar. Gaze locking: passive eye contact detection for human-object interaction. In *Proceedings of the 26th annual ACM symposium on User interface software and technology*, pages 271–280, 2013. 2, 5

[47] Yusuke Sugano, Yasuyuki Matsushita, and Yoichi Sato. Learning-by-synthesis for appearance-based 3d gaze estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1821–1828, 2014. 12

[48] Marco Toldo, Andrea Maracani, Umberto Michieli, and Pietro Zanuttigh. Unsupervised domain adaptation in semantic segmentation: a review. *Technologies*, 8(2):35, 2020. 2

[49] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7167–7176, 2017. 5

[50] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016. 8, 12

[51] Ben Usman, Nick Dufour, Kate Saenko, and Chris Bregler. Puppetgan: Cross-domain image manipulation by demonstration. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9449–9457, 2019. 2

[52] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19376–19385, 2022. 2

[53] Erroll Wood, Tadas Baltrušaitis, Louis-Philippe Morency, Peter Robinson, and Andreas Bulling. Learning an appearance-based gaze estimator from one million synthesised images. In *Proceedings of the Ninth Biennial ACM Symposium on Eye Tracking Research & Applications*, pages 131–138, 2016. 1

[54] Erroll Wood, Tadas Baltrušaitis, Louis-Philippe Morency, Peter Robinson, and Andreas Bulling. Gazedirector: Fully articulated eye gaze redirection in video. In *Computer Graphics Forum*, volume 37, pages 217–225. Wiley Online Library, 2018. 3

[55] Erroll Wood, Tadas Baltrusaitis, Xucong Zhang, Yusuke Sugano, Peter Robinson, and Andreas Bulling. Rendering of eyes for eye-shape registration and gaze estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3756–3764, 2015. 1

[56] Daniel E Worrall, Stephan J Garbin, Daniyar Turmukhambetov, and Gabriel J Brostow. Interpretable transformations with encoder-decoder networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5726–5735, 2017. 3, 6

[57] Weihao Xia, Yujiu Yang, Jing-Hao Xue, and Wensen Feng. Controllable continuous gaze redirection. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 1782–1790, 2020. 1, 3

[58] Yu Yu, Gang Liu, and Jean-Marc Odobez. Improving few-shot user-specific gaze adaptation via gaze redirection synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11937–11946, 2019. 2

[59] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor Lempitsky. Few-shot adversarial learning of realistic neural talking head models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9459–9468, 2019. 2

[60] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 6

[61] Xucong Zhang, Yusuke Sugano, and Andreas Bulling. Revisiting data normalization for appearance-based gaze estimation. In *Proceedings of the 2018 ACM Symposium on Eye Tracking Research & Applications*, pages 1–9, 2018. 5, 12

[62] Xucong Zhang, Yusuke Sugano, Mario Fritz, and Andreas Bulling. Appearance-based gaze estimation in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4511–4520, 2015. 2, 5

[63] Xucong Zhang, Yusuke Sugano, Mario Fritz, and Andreas 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. 5

[64] Yufeng Zheng, Seonwook Park, Xucong Zhang, Shalini De Mello, and Otmar Hilliges. Self-learning transformations for improving gaze and head redirection. *Advances in Neural Information Processing Systems*, 33:13127–13138, 2020. 1, 2, 3, 6, 7, 8, 12

[65] Yang Zou, Zhiding Yu, BVK Kumar, and Jinsong Wang. Unsupervised domain adaptation for semantic segmentation via class-balanced self-training. In *Proceedings of the European conference on computer vision (ECCV)*, pages 289–305, 2018. 2

## A. Data Pre-processing

We follow the same data pre-processing pipeline as done in Park *et al*. [38]. The pipeline consists of a normalization technique [61] initially introduced by Sugano *et al*. [47]. It is followed by face detection [15] and facial landmarks detection [8] modules for which open-source implementations are publicly available. The Surrey Face Model [19] is used as a reference 3D face model. Further details can be found in Park *et al*. [38]. To summarize, we use the public code[2] provided by Park *et al*. [38] to produce image patches of size $256 \times 64$ containing both eyes.

## B. Architecture Details

**Our framework CUDA-GHR.** We use DenseNet architecture [16] to implement image encoder $\mathbf{E}_a$. The DenseNet is formed with a growth rate of 32, 4 dense blocks (each with four composite layers), and a compression factor of 1. We use instance normalization [50] and leaky ReLU activation function ($\alpha = 0.01$) for all layers in the network. We remove dropout and $1 \times 1$ convolution layers. The dimension of latent factor $z^a$ is set to be equal to 16. Thus, to project CNN features to the latent features, we use global-average pooling and pass through a fully-connected layer to output 16-dimensional feature from $\mathbf{E}_a$. The gaze encoder $\mathbf{E}_g$ is a MLP-based block whose architecture is shown in Table 4. The dimension of $z^g$ is set as 8.

For the generator network $\mathbf{G}$, we use HoloGAN [36] architecture shown in Table 8. The latent vector $z$ for each AdaIN [17] input is processed by a 1-layer MLP, and the rotation layer is the same as the one used in the original paper [36]. The latent domain discriminator $\mathbf{D}_F$ consists of 4 MLP layers as shown in Table 5. It takes the input of dimension 24 and gives 1-dimensional output. Both image discriminators $\mathbf{D}_T$ and $\mathbf{D}_S$ are PatchGAN [22] based networks as used in Zheng *et al*. [64]. The architecture of the discriminator is described in Table 7.

Table 4: Architecture of gaze encoder $\mathbf{E}_g$

| Layer name | Activation | Output shape |
|---|---|---|
| Fully connected | LeakyReLU ($\alpha = 0.01$) | 2 |
| Fully connected | LeakyReLU ($\alpha = 0.01$) | 2 |
| Fully connected | LeakyReLU ($\alpha = 0.01$) | 2 |
| Fully connected | None | 8 |

The task network $\mathcal{T}$ is a ResNet-50 [11] model with batch normalization [20] replaced by instance normalization [50] layers. It takes an input of $256 \times 64$ and gives a 4-dimensional output describing pitch and yaw angles for gaze and head directions. It is initialized with ImageNet [7] pre-trained weights and is fine-tuned on the GazeCapture

Table 5: Architecture of latent domain discriminator $\mathbf{D}_F$

| Layer name | Activation | Output shape |
|---|---|---|
| Fully connected | LeakyReLU ($\alpha = 0.01$) | 24 |
| Fully connected | LeakyReLU ($\alpha = 0.01$) | 24 |
| Fully connected | LeakyReLU ($\alpha = 0.01$) | 24 |
| Fully connected | None | 1 |

Table 6: Architecture of the task network $\mathcal{T}$

| Module/Layer name | Output shape |
|---|---|
| ResNet-50 layers with MaxPool stride=1 | $2048 \times 1 \times 1$ |
| Fully connected | 4 |

training subset for around 190K iterations. The GazeCapture validation subset is used to select the best-performing model. The initial learning rate is 0.0016, decayed by a factor of 0.8 after about 34K iterations. Adam [26] optimizer is used for optimization with a weight decay coefficient of $10^{-4}$. The architecture of $\mathcal{T}$ is summarized in Table 6.

**Downstream Tasks.** For gaze and head pose estimation, we use similar architecture as employed for $\mathcal{T}$ shown in Table 6. For all the experiments, the initial learning rate is 0.0001 decayed by a factor of 0.5 after every 1500 iterations. The pre-trained models are trained for 10 epochs with a batch size of 64 while fine-tuning is done for 5 epochs with a batch size of 32.

**State-of-the-art Baselines.** We re-implement the STED [64] on images containing both eyes for a fair comparison with our method using the public code[3] available. We use the same hyperparameters as provided by the original implementation. For the accurate comparison, we replaced *tanh* non-linearity with an identity function and removed a constant factor of $0.5\pi$ in all the modules.

## C. Additional Results

In Figures 5 and 6, we show additional qualitative results for both target datasets, namely, MPIIGaze and Columbia. Figure 5a and 6a represent gaze redirected images and Figure 5b and 6b show head redirected images.

---

Table 7: Architecture of the image discriminator networks $\mathbf{D}_T$ and $\mathbf{D}_S$. Note that, both the discriminators has the same architecture.

| Layer name | Kernel, Stride, Padding | Activation | Normalization | Output shape |
|---|---|---|---|---|
| Conv2d | 4×4, 2, 1 | LeakyReLU ($\alpha = 0.2$) | - | 64× 32×128 |
| Conv2d | 4×4, 2, 1 | LeakyReLU ($\alpha = 0.2$) | InstanceNorm | 128×16×64 |
| Conv2d | 4×4, 2, 1 | LeakyReLU ($\alpha = 0.2$) | InstanceNorm | 256×8×32 |
| Conv2d | 4×4, 1, 1 | LeakyReLU ($\alpha = 0.2$) | InstanceNorm | 512×7×31 |
| Conv2d | 4×4, 1, 1 | - | - | 1×6×30 |

Table 8: Architecture of the generator network $\mathbf{G}$

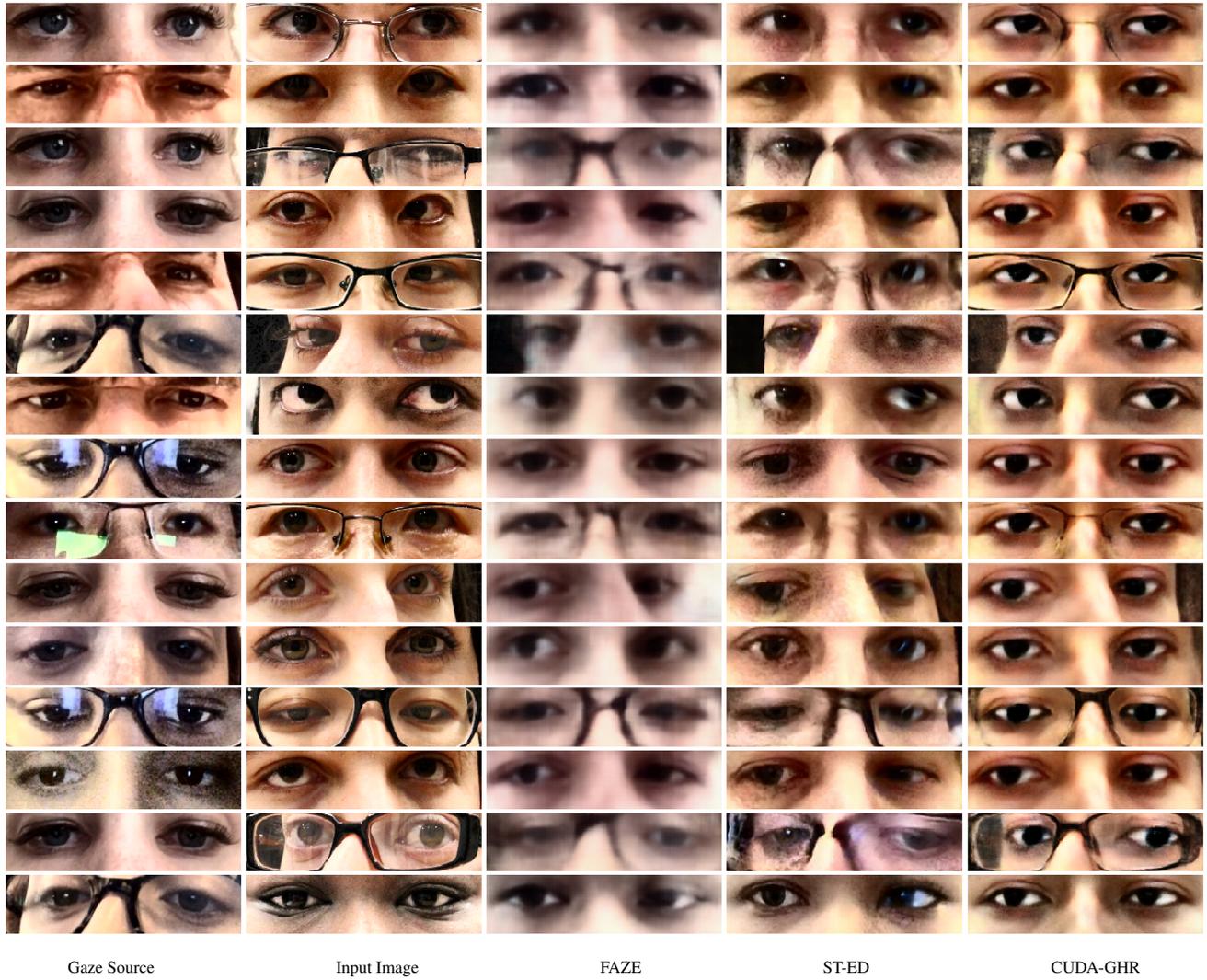| Layer name | Kernel | Activation | Normalization | Output shape |
|---|---|---|---|---|
| Learned Constant Input | - | - | - | 512×4×2×8 |
| Upsampling | - | - | - | 512×8×4×16 |
| Conv3d | 3×3×3 | LeakyReLU | AdaIN | 256×8× 4×16 |
| Upsampling | - | - | - | 256×16×8×32 |
| Conv3d | 3×3×3 | LeakyReLU | AdaIN | 128×16×8×32 |
| Volume Rotation | - | - | - | 128×16×8×32 |
| Conv3d | 3×3×3 | LeakyReLU | - | 64×16×8×32 |
| Conv3d | 3×3×3 | LeakyReLU | - | 64×16×8×32 |
| Reshape | - | - | - | $(64 \cdot 16)$×8×32 |
| Conv2d | 1×1 | LeakyReLU | - | 512×8×32 |
| Conv2d | 4×4 | LeakyReLU | AdaIN | 256×8×32 |
| Upsampling | - | - | - | 256×16×32 |
| Conv2d | 4×4 | LeakyReLU | AdaIN | 64×16×64 |
| Upsampling | - | - | - | 64×32×128 |
| Conv2d | 4×4 | LeakyReLU | AdaIN | 32×32×128 |
| Upsampling | - | - | - | 32×64×256 |
| Conv2d | 4×4 | Tanh | - | 3×64×256 |

| Gaze Source | Input Image | FAZE | ST-ED | CUDA-GHR |

(a) Gaze Redirected images for MPIIGaze dataset (*GazeCapture→MPIIGaze*)
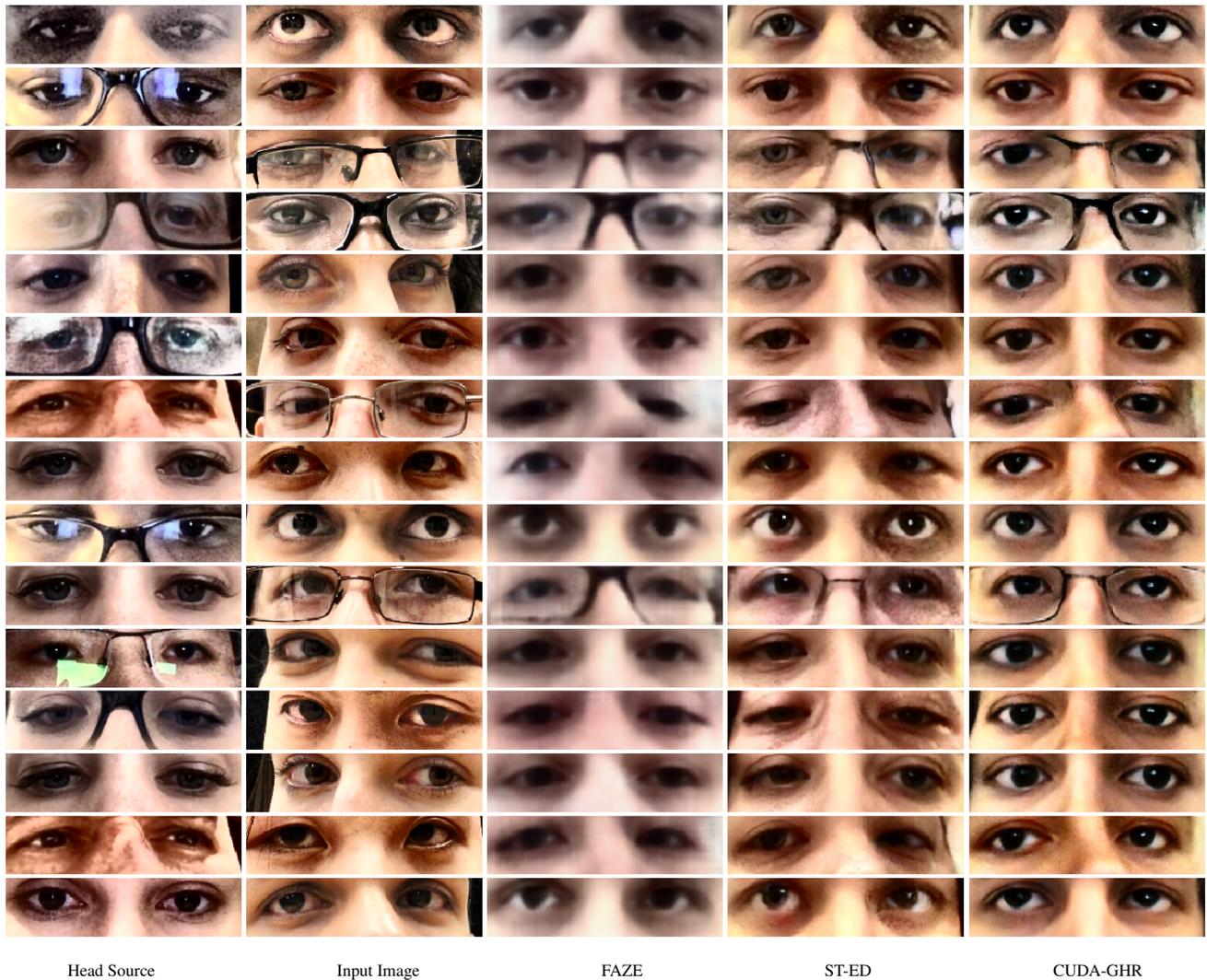
| Head Source | Input Image | FAZE | ST-ED | CUDA-GHR |

(b) Head Redirected images for MPIIGaze dataset (*GazeCapture→MPIIGaze)*

Figure 5: **Additional Qualitative Results (*GazeCapture→MPIIGaze*):** More qualitative results on the MPIIGaze dataset. 5a shows the gaze redirected images and 5b shows the head redirected images. The first column shows the gaze/head pose source image from which gaze/head pose information is used to redirect. The second column shows the input image from the MPIIGaze dataset. Best viewed in color.

| Gaze Source | Input Image | FAZE | ST-ED | CUDA-GHR |

(a) Gaze Redirected images for Columbia dataset (*GazeCapture→Columbia*)

| Head Source | Input Image | FAZE | ST-ED | CUDA-GHR |

(b) Head Redirected images for Columbia dataset (*GazeCapture→Columbia*)

Figure 6: **Additional Qualitative Results (*GazeCapture→Columbia*):** Qualitative results on the Columbia dataset. 6a shows the gaze redirected images and 6b shows the head redirected images. The first column shows the gaze/head pose source image from which gaze/head pose information is used to redirect. The second column shows the input image from the Columbia dataset. Best viewed in color.