

# Gaze-LLE: Gaze Target Estimation via Large-Scale Learned Encoders

Fiona Ryan<sup>1</sup> Ajay Bati<sup>1</sup> Sangmin Lee<sup>2</sup> Daniel Bolya<sup>1</sup> Judy Hoffman<sup>1\*</sup> James M. Rehg<sup>2\*</sup>  
 Georgia Institute of Technology<sup>1</sup> University of Illinois Urbana-Champaign<sup>2</sup>  
 {fkryan, abati7, dbolya, judy}@gatech.edu {sangminl, jrehg}@illinois.edu

## Abstract

We address the problem of gaze target estimation, which aims to predict where a person is looking in a scene. Predicting a person's gaze target requires reasoning both about the person's appearance and the contents of the scene. Prior works have developed increasingly complex, hand-crafted pipelines for gaze target estimation that carefully fuse features from separate scene encoders, head encoders, and auxiliary models for signals like depth and pose. Motivated by the success of general-purpose feature extractors on a variety of visual tasks, we propose **Gaze-LLE**, a novel transformer framework that streamlines gaze target estimation by leveraging features from a frozen DINOv2 encoder. We extract a single feature representation for the scene, and apply a person-specific positional prompt to decode gaze with a lightweight module. We demonstrate state-of-the-art performance across several gaze benchmarks and provide extensive analysis to validate our design choices. Our code is available at: <https://github.com/fkryan/gazelle>.

## 1. Introduction

Gaze is an important component of human behavior, giving insight into how a person interacts with the world around them. A person's visual attention indicates intent during daily activities [18, 28, 37, 68], and plays a key role in social interactions [13]. Humans can perform *gaze-following*, which is the ability to assess where another person is looking. From childhood, we learn to follow the gaze of a social partner to engage in joint attention [19, 43, 60]. In conversations, we use gaze to infer who someone is talking to, or resolve what object they are talking about. Thus, the ability to estimate gaze targets is an essential building block for developing systems that understand human behavior.

A significant number of prior works have proposed specialized architectures and datasets for gaze target estimation. A key property of these architectures is a *multi-branch*

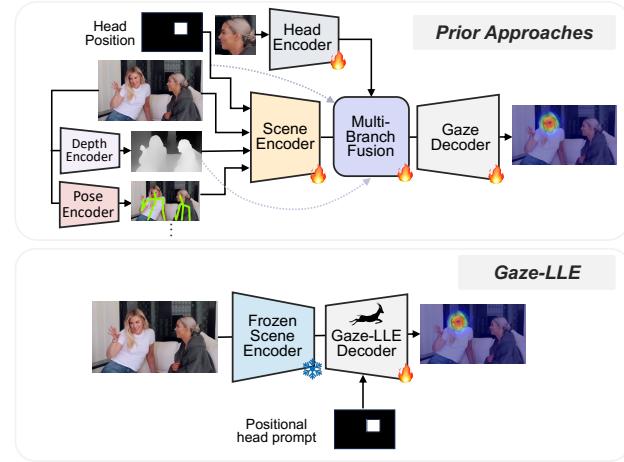


Figure 1. Prior approaches for gaze target estimation carefully fuse features from a separate head encoder, scene encoder, and auxiliary models for multimodal cues like depth and pose. We propose Gaze-LLE, a novel, streamlined approach that uses a single feature representation from a frozen image encoder and injects a person-specific positional prompt to decode gaze targets.

design, consisting of a head branch that extracts visual features from a crop of person's head and a scene branch that extracts features from the full image [4, 8, 9, 36, 50, 52, 59, 71]. More recent works incorporate additional modalities such as depth [2, 16, 21, 30, 42, 58, 62] and human pose [2, 22]. While these models have achieved impressive performance, they are limited to training on *small-scale* datasets obtained by asking human annotators to label gaze targets in images. In contrast, tasks such as segmentation and depth estimation have benefited substantially from self-supervised foundation models trained on *large-scale* data. It is thus natural to ask: Can gaze target estimation similarly benefit from a foundation model-based approach?

In this paper, we demonstrate for the first time that pretrained visual feature representations, produced by *transformer-based foundation models* like DINOv2, can be leveraged via a novel architecture to yield *state-of-the-art* gaze estimation performance. We introduce our model for Gaze estimation via Large-scale Learned Encoders (Gaze-LLE), and show that it has two benefits. First, by establish-

\*Equal advising.

ing the feasibility of the foundation model-based approach, we enable gaze estimation to join the ranks of other dense prediction tasks in terms of leveraging and benefiting from steady advances in the performance of foundation models.

Second, we show that by leveraging a powerful general purpose backbone we can simplify the model architecture significantly, reducing the learnable parameters by 1-2 orders of magnitude, and reducing the complexity of the training process while improving its efficiency. These architectural benefits are illustrated in Fig. 1, which contrasts Gaze-LLE against the standard multi-branch approach that uses multiple specialized encoders to capture gaze cues. This approach requires the careful fusion of different learned representations, taking their spatial relationships into account, with complex loss functions and training procedures. In contrast, Gaze-LLE leverages advances in general-purpose foundational feature representations that can solve dense prediction tasks, such as depth estimation, using only linear projection without representation tuning [45].

Interestingly, merely substituting DINOv2 as a backbone in prior gaze architectures doesn’t perform well. In fact, this results in significantly *worse* performance than the original backbones (Tab. 1). Our solution is the design of our novel Gaze-LLE decoder, which adapts DINOv2 for gaze prediction. In addition, we provide substantial analysis of the challenges arising in leveraging foundation models for our task, along with extensive empirical experiments to quantify architectural decisions and tradeoffs.

Our contributions are as follows: We introduce the novel **Gaze-LLE architecture** (Sec. 3.1) containing a specially-designed decoder that solves the problem of leveraging vision foundation models for gaze target estimation. We identify the key technical challenges and explain why the naive use of large-scale models is ineffective (Sec. 3.2), and we validate the optimality of our design decisions (Sec. 4.2). The Gaze-LLE model is **streamlined**, with just  $\sim 5\%$  of the trainable parameters used in most prior methods (Tab. 3); **powerful**, achieving state-of-the-art performance across the three main gaze estimation benchmarks (Sec. 4.1); **general**, exhibiting strong cross-dataset performance *without finetuning* (Tab. 5); and **easy to train**, achieving state-of-the-art in  $< 1.5$  GPU-hours (Fig. 4). We release our code and models in the hope that even more powerful gaze estimators can be developed from Gaze-LLE.

## 2. Related work

The dominant approach to gaze target estimation is a *multi-branch fusion* approach, in which an initial encoder is followed two or more analysis branches that work in parallel to extract specific cues for gaze estimation. These branches converge in a fusion module which produces an integrated representation which is then decoded into the output heatmap, with end-to-end training of the entire pipeline.

Method	Encoder	AUC $\uparrow$	Avg L2 $\downarrow$	Min L2 $\downarrow$
Chong et al. [9]	<b>Original (Res50)</b>	<b>0.921</b>	<b>0.137</b>	<b>0.077</b>
	Trained DINOv2 ViT-B	0.908	0.167	0.101
	Frozen DINOv2 ViT-B	0.875	0.191	0.125
Miao et al. [42]	<b>Original (Res50)</b>	<b>0.934</b>	<b>0.123</b>	<b>0.065</b>
	Trained DINOv2 ViT-B	0.910	0.152	0.093
	Frozen DINOv2 ViT-B	0.892	0.173	0.109
Gupta et al. [22] IMAGE	<b>Original</b>	<b>0.933</b>	<b>0.134</b>	<b>0.071</b>
	Trained DINOv2 ViT-B	0.912	0.155	0.090
	Frozen DINOv2 ViT-B	0.894	0.184	0.116

Table 1. Existing gaze architectures do not leverage features from large transformer models effectively. We replace the scene encoder in 3 existing open source methods with the DINOv2 ViT-B backbone and evaluate on GazeFollow (see Appendix for details). Using DINOv2 does *not* improve performance—whether or not its parameters are frozen.

In contrast, the goal of this paper is to demonstrate that *multi-branch architectures are unnecessary and inefficient*. We show that SotA performance can be obtained by processing the feature representation produced by a frozen foundational visual encoder, using a novel decoder architecture which incorporates the unique properties of gaze estimation. We are the first to demonstrate the feasibility of using a frozen large-scale encoder for this task, and the design of our decoder architecture is novel relative to prior gaze estimation works.

The origin of the multi-branch approach is Recasens et al. [50], which also introduced the GazeFollow dataset. Their two-branch architecture consisted of a **scene branch** to estimate scene saliency, and a **head branch** to refine the saliency map for a specific person. This approach was adopted by many subsequent works [4, 7, 9, 27, 36, 51, 52, 59, 67, 71]. More recent works extended the paradigm by incorporating additional cues via auxiliary models for depth [2, 16, 21, 30, 42, 58, 62], body pose [2, 22], 3D head direction [16, 25], eye location [16], and object detections [27].

A key property of gaze estimation is the need to integrate features extracted from the head region of the target person with other scene cues. This ensures that head pose, for example, is correctly interpreted in the context of the scene. Multi-branch approaches [2, 9, 16, 22, 25, 30, 36, 42, 58, 59, 62] solve this problem by using the head representation as an input to the scene branch, thereby requiring the scene to be encoded separately for each person, and by carefully crafting fusion mechanisms that combine feature representations across the head branch, scene branch, and other branches. Additionally, many train with a complex multitask objective in order to supervise each encoder differently [2, 16, 22, 30, 58, 59, 62]. These complex architectures can be challenging to train and often converge slowly, as illustrated in Fig. 4. In contrast, we provide a head position prompt as a separate input to our unified decoder archi-

- Why head pose is important?
- How does the estimation of head pose differs in multiple-branch approach vs GazeLLE?

ture (see Fig. 2), and all scene cues are extracted within the decoder, eliminating the need for a separate head branch and for multi-task objectives and fusion modules.

Among the multi-branch architectures, the two-branch approach of Tafasca et. al. [59] is the most closely-related to this paper. They use a large transformer-based backbone for the scene analysis branch, which receives the head branch representation as input. While they initialize their backbone with pretrained weights, it is still trained end-to-end, with the head branch producing a specialized feature representation. In contrast, we demonstrate that 1) head analysis can also be directly integrated into the decoder, eliminating the need for a head branch and further simplifying and streamlining the architecture, and 2) frozen large-scale foundational encoders give superior performance with two orders of magnitude fewer learned parameters. Gaze-LLE produces higher accuracy on all datasets (see Table 3).

Some prior works [63–65] have explored an alternative formulation of gaze target estimation as a set detection problem, where a model based on DETR [3] jointly predicts the location of all heads and their accompanying gaze targets. This avoids the need for a separate head detection step, but creates additional complications and existing works require the input of ground truth gaze data at testing time - preventing comparison with most methods (see Appendix). Other prior works have estimated 3D gaze direction from facial appearance [15, 31, 70] without identifying the gaze target, and some early approaches estimated head orientation [56, 69] to identify gaze targets.

Another area of gaze behavior recognition involves the joint analysis of multi-person social gaze behaviors, such as shared attention (when 2 people are looking at the same gaze target) [14, 24, 44, 47, 54, 57], mutual gaze (when a pair of people is looking at each other) [11, 38–41, 46], and other gaze-related social structures [15, 17]. Such analysis can be used in assessing and understanding social behaviors for conditions like autism [7, 9, 35], and may benefit from a gaze estimation approach like ours where social context is naturally encoded in a shared scene representation.

### 3. Gaze-LLE

**Problem Definition** Given an RGB image  $x_{\text{img}} \in \mathbb{R}^{3 \times H_{\text{in}} \times W_{\text{in}}}$  and the bounding box for a particular person’s head  $x_{\text{bbox}} \in \mathbb{R}^4$ , we predict a heatmap  $\mathcal{H} \in [0, 1]^{H_{\text{out}} \times W_{\text{out}}}$ , where each value represents the probability that the pixel is a gaze target. The VideoAttentionTarget and ChildPlay benchmarks include the additional task of predicting a value  $y \in [0, 1]$  that represents the probability that the given person’s gaze target is inside the frame.

#### 3.1. Model Architecture

Fig. 2 illustrates our novel Gaze-LLE architecture, consisting of a frozen, large-scale general-purpose *scene encoder*

and a learned Gaze Decoder module. Our gaze decoder performs *head prompting* to condition outputs on a particular person, updates the feature representation with a small transformer module, and predicts a gaze heatmap and if the target is in-frame. We describe each component in detail:

**Scene Encoder** A core component of our approach is leveraging strong visual features from a frozen, pretrained feature extractor  $\mathcal{F}$ , instead of learning a feature extractor end-to-end or using auxiliary models for signals like depth and pose.  $\mathcal{F}$  can be any visual feature extractor (see Sec. 4.2), but we primarily use DINOv2. From  $\mathcal{F}(x_{\text{img}})$ , we obtain a lower resolution feature map of size  $d_{\mathcal{F}} \times H \times W$ , which we then use a linear layer to project to a smaller dimension  $d_{\text{model}}$ , yielding a feature map  $x_{\mathcal{F}} \in \mathbb{R}^{d_{\text{model}} \times H \times W}$ .

**Head Position Embedding** A key consideration in our architecture is how to incorporate head position via head prompting. We find that incorporating head position *after* the scene encoder (rather than before as in prior work) gives the best performance (see Sec. 3.2 for a detailed discussion). We construct a downsampled, binarized mask  $M$  of size  $H \times W$  from the given head bounding box  $x_{\text{bbox}}$  within the extracted scene feature map. Using  $M$ , we add a learned position embedding  $p_{\text{head}} \in \mathbb{R}^{d_{\text{model}}}$  to the scene tokens containing the head (see Sec. 4.2 for alternatives). The scene feature map  $S$  is then:

$$S = x_{\mathcal{F}} + (M * p_{\text{head}}) \quad (1)$$

**Transformer Layers** To update the feature representation for our task, we train a small learnable transformer module,  $\mathcal{T}$ , which uses self-attention to process the head-conditioned scene features. As input to  $\mathcal{T}$ , we flatten the feature map with the added head position  $S$  into a scene token list  $[s_1, s_2, \dots, s_{H \times W}]$ . For the VideoAttentionTarget and ChildPlay benchmark settings, where the model also must classify whether the queried person’s gaze is in or out of the frame, we prepend a learnable *task token*,  $t_{\text{in/out}}$ , to the token list. Our token list is then:

$$[ \underbrace{t_{\text{in/out}}}_{\text{task token}}, \underbrace{s_1, s_2, \dots, s_{H \times W}}_{\text{scene tokens}} ] \quad (2)$$

Due to the spatial nature of our task, we add absolute 2d sinusoidal position embeddings [12]  $P$  to the scene features before they are input to  $\mathcal{T}$ , i.e.,  $\mathcal{T}(S + P)$ . By default,  $\mathcal{T}$  consists of 3 standard transformer encoder layers [66].

**Prediction Heads** From  $\mathcal{T}(S + P)$ , we obtain the updated scene features  $S'$ , and the updated task token,  $t'_{\text{in/out}}$ . We reconstruct  $S'$  into a feature map of size  $d_{\text{model}} \times H \times W$ , and pass it to the gaze heatmap decoder  $\mathcal{D}_{\text{hm}}$ .  $\mathcal{D}_{\text{hm}}$  consists of 2 convolutional layers to upsample the feature map to the output size  $H_{\text{out}} \times W_{\text{out}}$  and produce a classification score

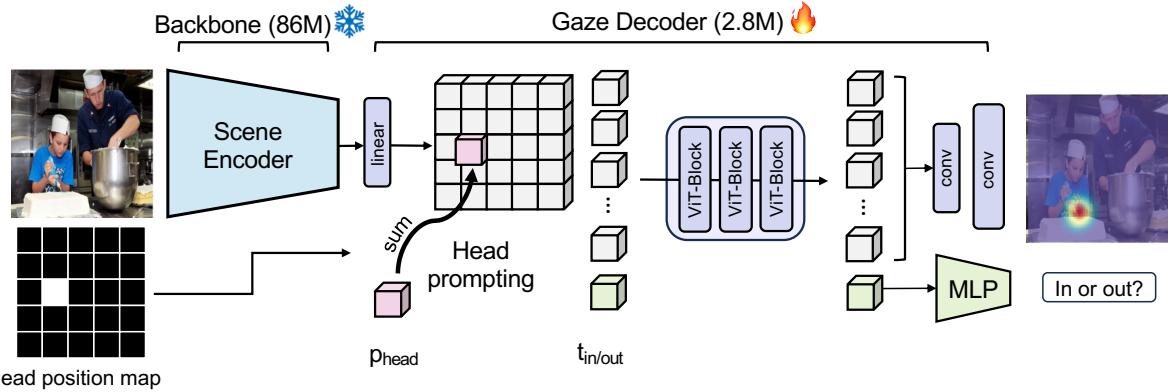


Figure 2. We introduce Gaze-LLE, a new framework for gaze estimation that learns a small gaze decoder on top of a frozen DINOv2 backbone. Using this backbone, we first extract scene tokens from an RGB image and project them to  $d_{\text{model}}$  with a linear layer. We then perform *head prompting* by adding a learned head position embedding  $p_{\text{head}}$  to the scene tokens at a given person’s head location. Next, we update the scene tokens and an optional learnable auxiliary in/out prediction task token  $t_{\text{in/out}}$  with 3 transformer layers. Finally, we upsample and decode the scene tokens into a heatmap and use the in/out task token to predict if the gaze target is in or out of frame.

for each pixel as being a gaze target or not. A 2-layer MLP  $\mathcal{D}_{\text{in/out}}$  takes  $t_{\text{in/out}}$  and outputs a classification score for if the queried person’s gaze target is in or out of frame.

**Training Objective** We train our model using pixel-wise binary cross-entropy loss for the heatmap. Following prior work [9, 50], the supervisory signal is an  $H_{\text{out}} \times W_{\text{out}}$  heatmap constructed by placing a 2D Gaussian distribution with  $\sigma = 3$  around each ground truth  $(x, y)$  gaze annotation. For benchmark settings where the model must jointly predict if the gaze is in or out of frame, we use a multitask loss

$$\mathcal{L} = \mathcal{L}_{\text{hm}} + \lambda \mathcal{L}_{\text{in/out}} \quad (3)$$

where  $\mathcal{L}_{\text{hm}}$  is pixel-wise binary cross entropy loss and  $\mathcal{L}_{\text{in/out}}$  is binary cross entropy loss for the in/out prediction task weighted by  $\lambda \in \mathbb{R}$ . This loss is much simpler and easier to optimize than the complex multi-task losses employed in prior works. The backbone  $\mathcal{F}$  is frozen during training. Our model with a ViT-B backbone has  $\sim 2.8\text{M}$  learnable parameters—significantly fewer than all prior works.

### 3.2. Key Design Decisions for Foundation Models

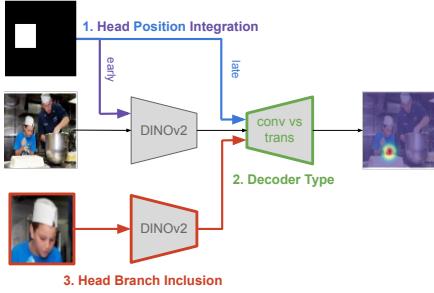
A key component of our approach is the use of a pretrained visual encoder (e.g. DINOv2 [45]) as a *single* backbone for gaze target estimation, without any other auxiliary models. There are many possible ways to incorporate such an encoder, and in this section we systematically identify the relevant issues and explore the design space, providing empirical support for the architectural choices in Sec. 3.1.

Our first finding is that a straightforward substitution of DINOv2 into prior gaze architectures leads to consistently poor performance. In Tab. 1, we show the result of swapping the scene encoder in three open source gaze estimation methods with DINOv2 (see Appendix for more results).

Whether or not we finetune the DINOv2 backbone, it is outperformed by the supposedly “weaker” backbones in these prior works. This is not necessarily surprising, as prior gaze works [30, 64, 65] have found ResNet-50 to sometimes outperform more powerful architectures. But this finding reinforces the need for our Gaze-LLE solution. To gain further insight, we conducted a set of experiments on GazeFollow with a simple baseline: extract scene and head features with a frozen DINOv2, concatenate the results, and decode into a gaze heatmap (see Appendix for full details). With this baseline, we quantify the impact of our three key architectural choices: integration of the head position, design of feature decoding, and use of a head branch (Tab. 2).

**Where should we inject the head position?** The position of a person’s head is an important cue in determining their gaze. Almost all prior works give head position as an extra channel to the scene branch (*i.e.*, RGB + head position), which requires the scene encoder to learn how to use it when finetuning on gaze. This is problematic if we want to exploit pre-trained frozen encoders without finetuning them. We find that simply concatenating the head position channel *after* extracting DINOv2 features boosts performance significantly (Tab. 2: a v.s. c) compared to retraining the input projection to accept it as an additional channel.

**How should we decode the DINOv2 features?** Most prior works decode combined scene and head features using a stack of conv layers. This may work well when using a gaze-specialized scene encoder. However, when using a frozen DINOv2, the receptive field of a few conv layers may be too small to extract long range gaze targets in the scene. We compare using a traditional 6 conv stack to decode heatmaps vs. one transformer layer into a 2-layer conv decoder. Both arrangements have the same number of pa-



	(1) Head Integration	(2) Decoder	(3) Branches	GazeFollow		
				AUC ↑	Avg L2 ↓	Min L2 ↓
a.	early	conv	H+S	0.854	0.254	0.168
b.	early	tran	H+S	0.904	0.178	0.113
c.	late	conv	H+S	0.932	0.155	0.089
d.	late	tran	H+S	<b>0.954</b>	<b>0.113</b>	<b>0.053</b>
e.	late	conv	S	0.916	0.184	0.115
f.	late	tran	S	<b>0.953</b>	<b>0.114</b>	<b>0.054</b>

Table 2. We investigate design choices across 3 axes: (1) early vs. late head integration, (2) convolutional vs. transformer decoder, and (3) using a head & scene branch (H+S) vs. a scene branch alone (S). Row a is the setting most similar to prior work. Conversely, we develop our final Gaze-LLE design from row f.

rameters, but the transformer layer can make use of global information, thus performing better (Tab. 2: c v.s. d).

**Do we need a head branch?** Prior works use a separate encoder that inputs a crop of the head, which is useful for understanding gaze direction. We hypothesize that a large-scale encoder like DINOv2 already captures gaze direction in its representation. We compare performance with and without a head branch and find it to be nearly the same when using a transformer-based decoder (d v.s. f). Notably, this doesn't occur with the conv decoder (c v.s. e), indicating that the relevant features are already there, but we need a transformer's global information propagation to extract them. This experiment motivated our novel head prompting design in Gaze-LLE.

## 4. Experiments

**Datasets** We conduct experiments on GazeFollow [50] and VideoAttentionTarget [9], which are the primarily used benchmarks for gaze target estimation. To assess our model's generalizability to other domains, we also include experiments on ChildPlay [58], a recent benchmark focusing on the gaze behaviors of children, and GOO-Real [61], which captures gaze in a retail environment.

**Evaluation Metrics** We evaluate our model's performance by calculating heatmap AUC, which uses each heatmap pixel as a confidence score for an ROC curve, and pixel L2, which is the Euclidean distance between the

argmax of the predicted heatmap and the ground truth gaze target. For GazeFollow, which contains  $\sim 10$  unique annotations per image, we report the distance to the average of the annotations (Avg L2), and the distance to the closest annotation (Min L2). For VideoAttentionTarget, ChildPlay, and GOO-Real, AUC is calculated by defining a tolerance region around the ground truth gaze point; for all, we follow the benchmark's specific calculation. For VideoAttentionTarget and ChildPlay, we also report the average precision (AP) for the in/out of frame prediction task.

**Technical Details** As in prior work, our model produces a gaze heatmap of size  $64 \times 64$ . Because our model does not include a separate head branch that operates on a high resolution crop of each head, we use an input image size of  $448 \times 448$  to capture dense details like eyes, while maintaining a small token list for computational efficiency. We conduct experiments with frozen DINOv2 ViT-B and ViT-L backbones. With DINOv2's patch size of 14, the internal feature map is size  $32 \times 32$ . We use an internal dimension of  $d_{model} = 256$  and 3 transformer layers with 8 attention heads and MLP dimension 1024. We train our model on GazeFollow for 15 epochs using the Adam optimizer, cosine scheduling with initial learning rate 1e-3, and batch size 60. We use random crop, flip, and bounding box jitter as data augmentation during training, and drop path regularization [34] with  $p = 0.1$ . We finetune our GazeFollow model on VideoAttentionTarget and ChildPlay with the multitask loss in Eq. 3. For VideoAttentionTarget, we train for 8 epochs with lr=1e-2 (in/out params) & lr=1e-5 (other gaze decoder params), and  $\lambda = 1$ . For ChildPlay, we train for 3 epochs with lr=2e-4 & 1e-4 and  $\lambda = 0.1$ .

## 4.1. Main Results

**Comparison to State-of-the-Art** Tab. 3 compares the performance of Gaze-LLE with existing methods on GazeFollow and VideoAttentionTarget. Our model achieves SotA on the AUC and L2 metrics for both datasets, while using only using a single image encoding branch and with a small fraction of the learnable parameters of prior approaches. We include results for both DINOv2 ViT-B and DINOv2 ViT-L, observing that the ViT-L backbone produces stronger results while ViT-B is still sufficient for obtaining SotA. Importantly, our ViT-B model outperforms Tafasca et al. [59], which also uses a ViT-B backbone, but trains the backbone end-to-end along with a separate head branch and large DPT decoder. Our approach's design effectively leverages the power of the pretrained backbone, achieving stronger results with  $\approx 2\%$  of the learned parameters. An additional benefit of our approach is dramatically reduced training time. As illustrated in Fig. 4, our model converges much faster than prior methods, achieving SotA in less than 1.5 hours on a single RTX4090. 😊😊😊

Method	Learnable Params	Input	GazeFollow			VideoAttentionTarget		
			AUC $\uparrow$	Avg L2 $\downarrow$	Min L2 $\downarrow$	AUC $\uparrow$	L2 $\downarrow$	AP <sub>in/out</sub> $\uparrow$
<i>One Human</i>								
Recasens et al. [50]	50M*	I	0.878	0.19	0.113	-	-	-
Chong et al. [8]	51M*	I	0.896	0.187	0.112	0.833	0.171	0.712
Lian et al. [36]	55M	I	0.906	0.145	0.081	-	-	-
Chong et al. [9]	61M	I	0.921	0.137	0.077	0.860	0.134	0.853
Chen et al. [4]	50M*	I	0.908	0.136	0.074	-	-	-
Fang et al. [16]	68M	I+D+E	0.922	0.124	0.067	0.905	0.108	0.896
Bao et al. [2]	29M*	I+D+P	0.928	0.122	-	0.885	0.120	0.869
Jin et al. [30]	>52M*	I+D+P	0.920	0.118	0.063	0.900	<b>0.104</b>	0.895
Tonini et al. [62]	92M	I+D	0.927	0.141	-	0.862 <sup>‡</sup>	0.125	0.742
Hu et al. [27]	>61M*	I+D+O	0.923	0.128	0.069	0.880	0.118	0.881
Gupta et al. [22]	35M	I+D+P	0.943	0.114	0.056	0.914	0.110	0.879
Horanyi et al. [25] <sup>†</sup>	46M <sup>†</sup>	I+D	0.896 <sup>†</sup>	0.196 <sup>†</sup>	0.127 <sup>†</sup>	0.832 <sup>†</sup>	0.199 <sup>†</sup>	0.800 <sup>†</sup>
Miao et al. [42]	61M	I+D	0.934	0.123	0.065	0.917	<b>0.109</b>	<b>0.908</b>
Tafasca et al. [58]	>25M*	I+D	0.939	0.122	0.062	0.914	<b>0.109</b>	0.834
Tafasca et al. [59]	>135M*	I	0.944	0.113	0.057	-	<b>0.107</b>	0.891
<b>Gaze-LLE (ViT-B)</b>	2.8M	I	<b>0.956</b>	<b>0.104</b>	<b>0.045</b>	<b>0.933</b>	<b>0.107</b>	<b>0.897</b>
<b>Gaze-LLE (ViT-L)</b>	2.9M	I	<b>0.958</b>	<b>0.099</b>	<b>0.041</b>	<b>0.937</b>	<b>0.103</b>	<b>0.903</b>

Table 3. Gaze target estimation results on GazeFollow and VideoAttentionTarget. We report the number of learnable parameters for each model, and if auxiliary models are used for inputs: I is image, D is depth, and P is pose, O is objects, and E is eyes. (\*Parameter estimate. <sup>†</sup>Our reimplementations, see Appendix. <sup>‡</sup>Metric re-evaluated to match benchmark’s calculation protocol [9].)

On VideoAttentionTarget’s in/out of frame prediction task, we obtain second-best results; however our method obtains the best results when considering all metrics together. Our approach’s strong, SotA results validate our hypothesis that DINOv2 features do indeed capture appropriate information for gaze target estimation, and that a single-stream design can outperform multi-stream architectures with additional modalities. To assess our model’s performance on a more specialized dataset, we include results on ChildPlay in Tab. 4. Our ViT-B and ViT-L models both achieve SotA results across all AUC, L2, and AP. We also include the benchmark’s P.Head metric [58], which assesses the precision predicting when gaze targets lie within a head bounding box; however due to its reliance on detections which aren’t always accurate, we find this metric tends not to correlate with the others.

Method	AUC $\uparrow$	L2 $\downarrow$	AP $\uparrow$	P.Head $\uparrow$
Gupta et al. [22]	0.919	0.113	0.983	<u>0.694</u>
Tafasca et al. [58]	<u>0.935</u>	<u>0.107</u>	0.986	0.663
Tafasca et al. [59]	-	<u>0.106</u>	0.990	0.600
<b>Gaze-LLE (ViT-B)</b>	<b>0.949</b>	<u>0.106</u>	<b>0.994</b>	<b>0.715</b>
<b>Gaze-LLE (ViT-L)</b>	<b>0.951</b>	<b>0.101</b>	<b>0.994</b>	0.662

Table 4. Gaze target estimation results on ChildPlay.

**Cross-dataset Results** We include results of our GazeFollow-trained model applied to VideoAttentionTarget, GOO-Real, and ChildPlay *without finetuning* in Tab. 5 and Fig. 3. Our approach achieves strong cross-dataset results across diverse domains and exhibits better general-

Method	VAT		GOO-Real		ChildPlay	
	AUC $\uparrow$	L2 $\downarrow$	AUC $\uparrow$	L2 $\downarrow$	AUC $\uparrow$	L2 $\downarrow$
Chong et al.[9]*	0.906	0.119	0.670	0.334	0.912	0.121
Jin et al. [30]	0.900	<u>0.104</u>	-	-	-	-
Tonini et al. [62] w/ UDA	-	-	0.840	0.238	-	-
Miao et al.[42]*	0.923	0.109	0.869	<u>0.202</u>	0.933	0.113
Gupta et al. [22]	0.907	0.137	-	-	0.923	0.142
Tafasca et al. [58]	0.911	0.123	-	-	0.932	0.115
<b>Gaze-LLE (B)</b>	<u>0.932</u>	<u>0.105</u>	<b>0.901</b>	<b>0.174</b>	0.946	0.114
<b>Gaze-LLE (L)</b>	<b>0.937</b>	<b>0.100</b>	<u>0.898</u>	<u>0.175</u>	<b>0.951</b>	<b>0.101</b>

Table 5. Cross-dataset results on VideoAttentionTarget (VAT), GOO-Real, and ChildPlay. (\*Results we evaluated ourselves from the official code releases.)

ization than approaches that achieve high results on GazeFollow but experience larger performance drops in cross-dataset settings (*e.g.*, Gupta et al. [22]). We attribute the strong generalizability of our method to using an encoder that is not specialized to a task or dataset, learning minimal parameters and thus not overfitting to a particular dataset, and not depending on auxiliary models, which may generalize poorly themselves. Like other methods, our model experiences the largest performance drop on GOO-Real due to (1) the large domain gap, as GOO-Real contains a unique retail environment where the user rarely faces the camera, and (2) the difference in annotation scheme - GOO-Real’s ground truth is sourced from instructing participants to look at certain objects, rather than what a human annotator can reasonably infer from an image. We obtain SotA cross-dataset results, surpassing Tonini et al.’s [62] method with unsupervised domain adaptation, which requires access to

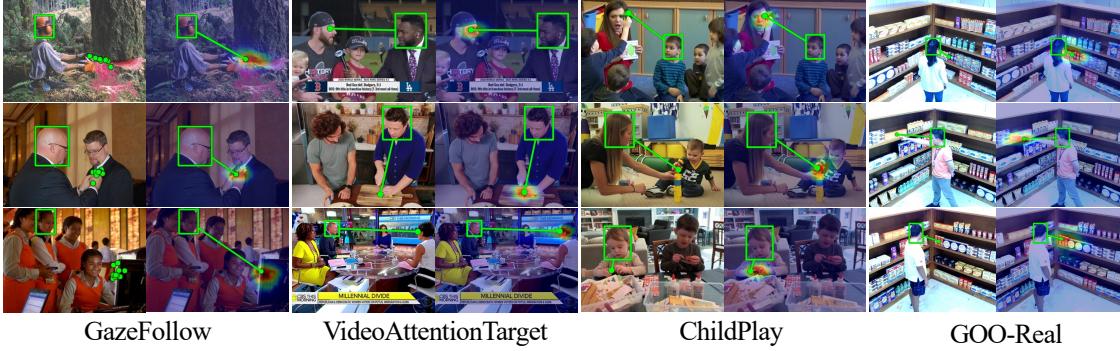


Figure 3. Qualitative results of our GazeFollow-trained ViT-B model on GazeFollow and applied **without finetuning** to VideoAttentionTarget, ChildPlay, and GOO-Real. We show ground truth on the left and the predicted heatmap & maximal point on the right.

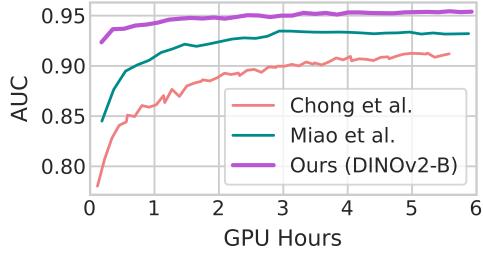


Figure 4. Training convergence: our method achieves strong results in fewer GPU hours than prior approaches.

Backbone	AUC $\uparrow$	Avg L2 $\downarrow$	Min L2 $\downarrow$
Supervised [55]	0.928	0.151	0.086
MAE [23]	0.947	0.126	0.061
CLIP [48]	0.953	0.107	0.049
DINOv2 [45]	<b>0.958</b>	<b>0.099</b>	<b>0.041</b>

Table 6. Ablation of different pretrained Vit-L backbones with Gaze-LLE on GazeFollow.

in-domain data at train time.

## 4.2. Analysis

In this section, we provide further insight into the optimality of our design choices for Gaze-LLE. We investigate different backbone feature extractors and alternative strategies for head prompting and evaluate on GazeFollow.

**Portability Across Backbones** While we use DINOv2 in our main experiments, Gaze-LLE can be used with any backbone. Table 6 reports our model’s performance with different pretrained encoders. The supervised [55] and MAE [23] models are pretrained on ImageNet-1k [10], while CLIP [48] and DINOv2 [45] are trained on much larger data sources. Unsurprisingly, DINOv2, which is the state-of-the-art for general-purpose feature extraction on dense downstream tasks, performs best, but CLIP also achieves strong results. As new backbones are developed,

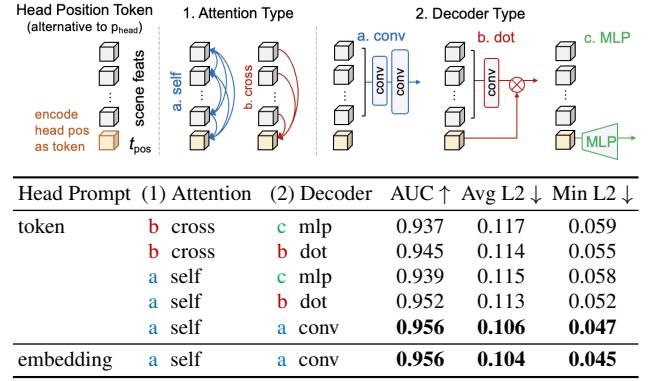


Table 7. As an alternative to adding the head position embedding  $p_{\text{head}}$  to the scene tokens, we explore representing the head’s center position as an additional token,  $t_{\text{pos}}$ . We consider self attention vs. cross attention across the token list, and different ways to decode the heatmap from the scene tokens and  $t_{\text{pos}}$ .

Gaze-LLE provides a framework for adapting them to gaze estimation.

**An Alternative Head Prompting Method** We also consider integrating the head position as its own token during attention as an alternative to our added position embedding  $p_{\text{head}}$ , inspired by works in point tracking and segmentation that represent positional queries as tokens [29, 32]. We construct a *head position token*,  $t_{\text{pos}}$ , by sampling the position embedding  $P$  at the head bounding box’s center point and summing this with a learned embedding. We concatenate this token to the scene token list,  $S$ . To fuse positional information with the scene features, we consider two types of attention in the transformer layers: self attention across the full token list, which updates both  $S$  and  $t_{\text{pos}}$ , and cross attention from the scene tokens to the position token, which updates only  $t_{\text{pos}}$ . Finally, we consider 3 methods for producing the heatmap from  $S$  and  $t_{\text{pos}}$ : a 2-layer convolutional decoder on  $S$  (as used in our default method), replacing

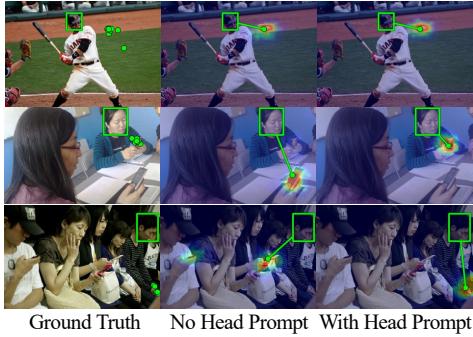


Figure 5. Without head prompting, our model succeeds on single-person cases, but cannot effectively condition gaze target estimation on the correct person in multi-person scenarios.

Prompt type	AUC ↑	Avg L2 ↓	Min L2 ↓
No prompting	0.926	0.169	0.105
With prompting	<b>0.956</b>	<b>0.104</b>	<b>0.045</b>

(a) Head prompt ablation

Prompt location	AUC ↑	Avg L2 ↓	Min L2 ↓
Layer 3	0.955	0.108	0.048
Layer 2	0.955	0.106	0.047
Layer 1 (default)	<b>0.956</b>	<b>0.104</b>	<b>0.045</b>

(b) Head prompt location

Table 8. We demonstrate the effectiveness of our head prompting mechanism (17), and find that injecting the head prompt before the first transformer layer in our gaze decoder module slightly outperforms later layers (8b)

the second convolutional layer with the dot product between  $t_{\text{pos}}$  and the scene feature map (like transformer segmentation methods [5, 6, 32]), and directly regressing the  $64 \times 64$  heatmap with a 2-layer MLP. We show results in Tab. 7 and find that given the right settings, this position token can be made almost as effective as our default embedding approach—however, with the added benefit of potentially being able to decode multiple head positions at the same time (with each new position being another token). For our lightweight decoder, additional head locations already add negligible compute (see Appendix), but future work with heavier gaze decoders may benefit from this token design.

**Ablating the Head Prompt** To assess the effectiveness of our head prompting mechanisms in providing the necessary information to decode a person’s gaze in the absence of a head branch, we perform ablation in Tab. 17 and Fig. 5 by inferencing our model without providing a head prompt. The resulting performance drop shows that our head prompting strategy is effective and necessary. However, we also observe an interesting phenomenon visible in Fig. 5: Without a head prompt, the model still predicts a

valid gaze target for at least one person in the scene – and in scenes with only one person, a head bounding box is not actually needed as input to effectively predict gaze! This result indicates that the scene representation implicitly detects heads and uses them to reason about potential gaze targets, providing further evidence for our hypothesis that a standalone head branch is not necessary. Our head prompting mechanism serves the purpose of identifying which person’s gaze should be decoded in multi-person scenes.

**Where to Perform Head Prompting** We also investigate *where* in the decoding process to inject the head prompt. To explore the tradeoff between performance and the reduction of person-specific computation, we move the head prompting to later transformer layers in our gaze decoder, and show that this yields only small performance drops (see Table 8b). This demonstrates SotA results on GazeFollow with the only per-person computation taking place in the final transformer layer and 2 convolutional layers. The ability to minimize person-specific computation while maintaining strong performance could provide efficient scaling for multi-person gaze analysis. We also investigate the performance when head detections from YOLOv5 are used instead of ground truth head bounding boxes, and observe *almost no degradation*. See the Appendix for the details.

## 5. Discussion

**Limitations** By leveraging a frozen encoder without end-to-end training, our performance is inherently tied to the encoder quality. We find it is important to select an encoder trained on a large, diverse dataset with a dense objective (see Tab. 6). Additionally, while our method is reasonably efficient (>50 fps on an RTX4090, see Appendix), the overall efficiency depends on the use of a large encoder, which may pose a challenge for embedded systems. We note that recent approaches that depend on auxiliary transformer-based depth/pose models also experience this limitation. However, as stronger, faster general-purpose feature extractors become available, Gaze-LLE provides a way to harness them for gaze estimation.

**Conclusion** In this work, we are the first to demonstrate that frozen foundational feature encoders can be leveraged for gaze target estimation. We propose Gaze-LLE, a new architecture that learns a gaze decoder with a novel head prompting design on top of a single, frozen DINOv2 encoder. We validate our design by achieving state-of-the-art results across four benchmarks and conducting experiments to validate the necessity and optimality of our design decisions. We hope our work opens a new chapter on gaze estimation by eliminating the need for complex multi-branch approaches via a streamlined and adaptable method that can be easily applied to new tasks and integrated into larger systems for understanding human behavior.

## Summary

### Problem Statement

Given an RGB image and the bounding box for a person's head in the image, predict a heatmap  $H \in [0, 1]$ , where each value represents the probability that the pixel is a gaze target. Optionally, predict  $y \in [0, 1]$  if the gaze target is inside the frame.

### Model Architecture

- Consists of a scene encoder and a gaze decoder module.
- The scene encoder is frozen, and the decoder is trainable. The gaze decoder performs head prompting to condition outputs on a particular person, updates the feature representation with a small transformer module, predicts a gaze heatmap, and (optionally) if the target is in-frame.
- Though the Gaze-LLE framework is flexible enough to incorporate any feature extractor, the authors use DiNOv2 as the scene encoder for extracting features. The feature map obtained ( $d_f \times H \times W$ ) is then projected using a linear layer to a smaller dimension  $d_{model}$  resulting in a feature map of size ( $d_{model} \times H \times W$ ).
- The head position is incorporated after the scene encoder. It provides the best performance compared to inserting it before the scene encoder. The authors construct a downsampled, binarized mask  $M$  of size ( $H \times W$ ) from the given head bounding box ( $x_{bbox}$ ) within the extracted scene feature map. They also add a learned position embedding to the scene tokens containing the head resulting in a final scene feature map as shown below.

$$S = x_{\mathcal{F}} + (M * p_{head})$$

$$\left[ \underbrace{t_{in/out}}_{\text{task token}}, \underbrace{s_1, s_2, \dots, s_{H \times W}}_{\text{scene tokens}} \right]$$

- The scene feature maps are then flattened into a list of scene tokens [ $s_1, s_2, \dots$ ]. If the in/out frame prediction is necessary for a person's gaze prediction, a learnable task token  $t_{(in/out)}$ , is also prepended to the token list. 2D sinusoidal positional encodings are added to these tokens and are fed into a trainable 3-encoder-layer transformer.
- The output feature map  $S'$  from the transformer is reconstructed into a feature map of size ( $d_{model} \times H \times W$ ), which is then passed to the gaze heatmap decoder  $D_{hm}$ .  $D_{hm}$  consists of 2 convolutional layers to upsample the feature map to the output size ( $H_{out} \times W_{out}$ ) and produce a classification score for each pixel as being a gaze target or not. A 2-layer MLP  $D_{(in/out)}$  takes ( $t_{in/out}$ ) and outputs a classification score for if the person's gaze target is in or out of frame.

## Training Objective

- pixel-wise binary cross-entropy for the heatmap
- If gaze in/out of the frame is also to be detected, then a binary cross-entropy loss for the in/out prediction task weighted by  $\lambda$  is also used.

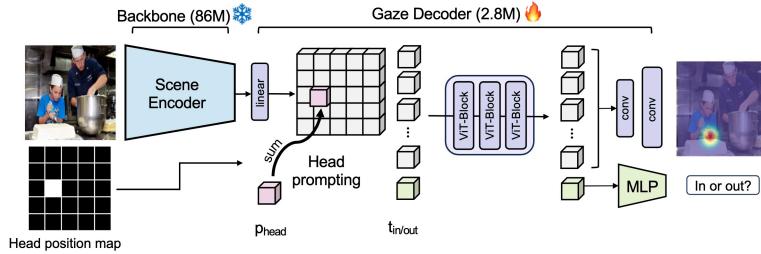


Figure 2. We introduce Gaze-LLE, a new framework for gaze estimation that learns a small gaze decoder on top of a frozen DINOv2 backbone. Using this backbone, we first extract scene tokens from an RGB image and project them to  $d_{\text{model}}$  with a linear layer. We then perform **head prompting** by adding a learned head position embedding  $p_{\text{head}}$  to the scene tokens at a given person's head location. Next, we update the scene tokens and an optional learnable auxiliary in/out prediction task token  $t_{\text{in/out}}$  with 3 transformer layers. Finally, we upsample and decode the scene tokens into a heatmap and use the in/out task token to predict if the gaze target is in or out of frame.

## Experimental Details

- Datasets: GazeFollow, VideoAttentionTarget, ChildPlay, and GOO-Rea
- AUCROC metric, pixel L2 loss
- DINOv2 with ViT-B and ViTL backbones, patch size=14,  $d_{\text{model}}=256$ , 3 transformer layers with 8 attention heads, and MLP dimension 1024.
- Adam optimizer with lr=1e-3 and batch\_size=60, and 15 epochs
- Random crop, flip, and bounding box jitter as data augmentation and drop path regularization with p = 0.1
- For VideoAttentionTarget, they fine-tuned it for 8 epochs with lr=1e-2 (in/out params), lr=1e-5 (for other gaze decoder params), and  $\lambda=1$ .
- For ChildPlay, they fine-tuned it for 3 epochs with lr=2e-4 and 1e-4 and  $\lambda=0.1$ .
- The model achieves SOTA in less than 1.5 hours on a single RTX4090 GPU

## Interesting findings

- Straightforward substitution of DINOv2 into prior gaze architectures leads to poor performance. The authors find that concatenating the head position channel after extracting DINOv2 features boosts performance significantly compared to retraining the input projection to accept it as an additional channel.
- Without head prompting, the model succeeds in single-person cases, but cannot effectively condition gaze target estimation on the correct person in multi-person scenarios.
- A separate head branch with DINOv2 is unnecessary, though adding it boosts performance.
- A gaze decoder without transformer layers does not provide optimal performance.

**Acknowledgments** The authors thank Stefan Stojanov for helpful discussions, and the members of the Hoffman Lab and Rehg Lab for their feedback. This work is in part supported by NIH R01 MH114999, NSF #2144194, and the NSF Graduate Research Fellowship under Grant No. DGE-2039655. Any opinion, findings, and conclusions or recommendations expressed in this material are those of the authors(s) and do not necessarily reflect the views of the National Science Foundation.

## References

- [1] <http://github.com/deepakcrk/yolov5-crowdhuman>. 6
- [2] Jun Bao, Buyu Liu, and Jun Yu. Escnet: Gaze target detection with the understanding of 3d scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14126–14135, 2022. 1, 2, 6
- [3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020. 3, 4
- [4] Wenhe Chen, Hui Xu, Chao Zhu, Xiaoli Liu, Yinghua Lu, Caixia Zheng, and Jun Kong. Gaze estimation via the joint modeling of multiple cues. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(3):1390–1402, 2021. 1, 2, 6
- [5] Bowen Cheng, Anwesa Choudhuri, Ishan Misra, Alexander Kirillov, Rohit Girdhar, and Alexander G Schwing. Mask2former for video instance segmentation. *arXiv preprint arXiv:2112.10764*, 2021. 8
- [6] Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. *Advances in Neural Information Processing Systems*, 34:17864–17875, 2021. 8
- [7] Eunji Chong, Katha Chanda, Zhefan Ye, Audrey Southerland, Nataniel Ruiz, Rebecca M Jones, Agata Rozga, and James M Rehg. Detecting gaze towards eyes in natural social interactions and its use in child assessment. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 1(3):1–20, 2017. 2, 3
- [8] Eunji Chong, Nataniel Ruiz, Yongxin Wang, Yun Zhang, Agata Rozga, and James M Rehg. Connecting gaze, scene, and attention: Generalized attention estimation via joint modeling of gaze and scene saliency. In *Proceedings of the European conference on computer vision (ECCV)*, pages 383–398, 2018. 1, 6
- [9] Eunji Chong, Yongxin Wang, Nataniel Ruiz, and James M Rehg. Detecting attended visual targets in video. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5396–5406, 2020. 1, 2, 3, 4, 5, 6, 7
- [10] 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. 7
- [11] Bardia Doosti, Ching-Hui Chen, Raviteja Vemulapalli, Xuhui Jia, Yukun Zhu, and Bradley Green. Boosting image-based mutual gaze detection using pseudo 3d gaze. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1273–1281, 2021. 3
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 3
- [13] Nathan J Emery. The eyes have it: the neuroethology, function and evolution of social gaze. *Neuroscience & biobehavioral reviews*, 24(6):581–604, 2000. 1
- [14] Lifeng Fan, Yixin Chen, Ping Wei, Wenguan Wang, and Song-Chun Zhu. Inferring shared attention in social scene videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6460–6468, 2018. 3
- [15] Lifeng Fan, Wenguan Wang, Siyuan Huang, Xinyu Tang, and Song-Chun Zhu. Understanding human gaze communication by spatio-temporal graph reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5724–5733, 2019. 3
- [16] Yi Fang, Jiapeng Tang, Wang Shen, Wei Shen, Xiao Gu, Li Song, and Guangtao Zhai. Dual attention guided gaze target detection in the wild. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11390–11399, 2021. 1, 2, 6
- [17] Alircza Fathi, Jessica K Hodgins, and James M Rehg. Social interactions: A first-person perspective. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1226–1233. IEEE, 2012. 3
- [18] Alireza Fathi, Yin Li, and James M Rehg. Learning to recognize daily actions using gaze. In *Computer Vision–ECCV 2012: 12th European Conference on Computer Vision, Florence, Italy, October 7–13, 2012, Proceedings, Part I* 12, pages 314–327. Springer, 2012. 1
- [19] Alexandra Frischen, Andrew P Bayliss, and Steven P Tipper. Gaze cueing of attention: visual attention, social cognition, and individual differences. *Psychological bulletin*, 133(4):694, 2007. 1
- [20] Clément Godard, Oisin Mac Aodha, Michael Firman, and Gabriel J Brostow. Digging into self-supervised monocular depth estimation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3828–3838, 2019. 7
- [21] Jian Guan, Liming Yin, Jianguo Sun, Shuhan Qi, Xuan Wang, and Qing Liao. Enhanced gaze following via object detection and human pose estimation. In *MultiMedia Modeling: 26th International Conference, MMM 2020, Daejeon, South Korea, January 5–8, 2020, Proceedings, Part II* 26, pages 502–513. Springer, 2020. 1, 2
- [22] Anshul Gupta, Samy Tafasca, and Jean-Marc Odobezi. A modular multimodal architecture for gaze target prediction: Application to privacy-sensitive settings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 5041–5050, 2022. 1, 2, 6

- [23] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022. 7
- [24] Matthew W Hoffman, David B Grimes, Aaron P Shon, and Rajesh PN Rao. A probabilistic model of gaze imitation and shared attention. *Neural Networks*, 19(3):299–310, 2006. 3
- [25] Nora Horanyi, Linfang Zheng, Eunji Chong, Aleš Leonardis, and Hyung Jin Chang. Where are they looking in the 3d space? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 2677–2686, 2023. 2, 6, 7, 8
- [26] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 7
- [27] Zhengxi Hu, Kunxu Zhao, Bohan Zhou, Hang Guo, Shichao Wu, Yuxue Yang, and Jingtao Liu. Gaze target estimation inspired by interactive attention. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(12):8524–8536, 2022. 2, 6
- [28] Chien-Ming Huang, Sean Andrist, Allison Sauppé, and Bilge Mutlu. Using gaze patterns to predict task intent in collaboration. *Frontiers in psychology*, 6:1049, 2015. 1
- [29] Wei Jiang, Eduard Trulls, Jan Hosang, Andrea Tagliasacchi, and Kwang Moo Yi. Cotr: Correspondence transformer for matching across images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6207–6217, 2021. 7
- [30] Tianlei Jin, Qizhi Yu, Shiqiang Zhu, Zheyuan Lin, Jie Ren, Yuanhai Zhou, and Wei Song. Depth-aware gaze-following via auxiliary networks for robotics. *Engineering Applications of Artificial Intelligence*, 113:104924, 2022. 1, 2, 4, 6
- [31] Petr Kellnhofer, Adria Recasens, Simon Stent, Wojciech Matusik, and Antonio Torralba. Gaze360: Physically unconstrained gaze estimation in the wild. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6912–6921, 2019. 3, 7
- [32] 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. *arXiv preprint arXiv:2304.02643*, 2023. 7, 8
- [33] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955. 3
- [34] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. Fractalnet: Ultra-deep neural networks without residuals. *ICLR*, 2017. 5
- [35] Peitong Li, Hui Lu, Ronald W Poppe, and Albert Ali Salah. Automated detection of joint attention and mutual gaze in free play parent-child interactions. In *International Conference on Multimodal Interaction*, pages 374–382. 2023. 3
- [36] Dongze Lian, Zehao Yu, and Shenghua Gao. Believe it or not, we know what you are looking at! In *Asian Conference on Computer Vision*, pages 35–50. Springer, 2018. 1, 2, 6
- [37] Kristian Lukander, Miika Toivanen, and Kai Puolamäki. Inferring intent and action from gaze in naturalistic behavior: a review. *International Journal of Mobile Human Computer Interaction (IJMHCI)*, 9(4):41–57, 2017. 1
- [38] Manuel J Marin-Jimenez, Andrew Zisserman, and Vittorio Ferrari. Here’s looking at you, kid. *Detecting people looking at each other in videos*. In *BMVC*, 5, 2011. 3
- [39] Manuel Jesús Marin-Jimenez, Andrew Zisserman, Marcin Eichner, and Vittorio Ferrari. Detecting people looking at each other in videos. *International Journal of Computer Vision*, 106:282–296, 2014.
- [40] Manuel J Marin-Jimenez, Vicky Kalogeiton, Pablo Medina-Suarez, and Andrew Zisserman. Laeo-net: revisiting people looking at each other in videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3477–3485, 2019.
- [41] Manuel J. Marín-Jiménez, Vicky Kalogeiton, Pablo Medina-Suárez, , and Andrew Zisserman. LAEO-Net++: revisiting people Looking At Each Other in videos. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 3
- [42] Qiaomu Miao, Minh Hoai, and Dimitris Samaras. Patch-level gaze distribution prediction for gaze following. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 880–889, 2023. 1, 2, 6, 5
- [43] Chris Moore, Philip J Dunham, and Phil Dunham. *Joint attention: Its origins and role in development*. Psychology Press, 2014. 1
- [44] Chihiro Nakatani, Hiroaki Kawashima, and Norimichi Ukita. Interaction-aware joint attention estimation using people attributes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10224–10233, 2023. 3
- [45] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 2, 4, 7
- [46] Cristina Palmero, Elsbeth A van Dam, Sergio Escalera, Mike Kelia, Guido F Lichtert, Lucas PJJ Noldus, Andrew J Spink, and Astrid van Wieringen. Automatic mutual gaze detection in face-to-face dyadic interaction videos. In *Proceedings of Measuring Behavior*, page 2, 2018. 3
- [47] Hyun Soo Park, Eakta Jain, and Yaser Sheikh. Predicting primary gaze behavior using social saliency fields. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3503–3510, 2013. 3
- [48] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 7
- [49] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE transactions on pattern analysis and machine intelligence*, 44(3):1623–1637, 2020. 1, 5

- [50] Adria Recasens, Aditya Khosla, Carl Vondrick, and Antonio Torralba. Where are they looking? *Advances in neural information processing systems*, 28, 2015. 1, 2, 4, 5, 6
- [51] Adria Recasens, Carl Vondrick, Aditya Khosla, and Antonio Torralba. Following gaze in video. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1435–1443, 2017. 2
- [52] Akanksha Saran, Srinjoy Majumdar, Elaine Schaertl Short, Andrea Thomaz, and Scott Niekum. Human gaze following for human-robot interaction. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 8615–8621. IEEE, 2018. 1, 2
- [53] Shuai Shao, Zijian Zhao, Boxun Li, Tete Xiao, Gang Yu, Xiangyu Zhang, and Jian Sun. Crowdhuman: A benchmark for detecting human in a crowd. *arXiv preprint arXiv:1805.00123*, 2018. 6
- [54] Hyun Soo Park and Jianbo Shi. Social saliency prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4777–4785, 2015. 3
- [55] Andreas Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit, and Lucas Beyer. How to train your vit? data, augmentation, and regularization in vision transformers. *arXiv preprint arXiv:2106.10270*, 2021. 7
- [56] Rainer Stiefelhagen, Jie Yang, and Alex Waibel. Modeling focus of attention for meeting indexing. In *Proceedings of the seventh ACM international conference on Multimedia (Part 1)*, pages 3–10, 1999. 3
- [57] Omer Sumer, Peter Gerjets, Ulrich Trautwein, and Enkelejda Kasneci. Attention flow: End-to-end joint attention estimation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3327–3336, 2020. 3
- [58] Samy Tafasca, Anshul Gupta, and Jean-Marc Odobez. Childplay: A new benchmark for understanding children’s gaze behaviour. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20935–20946, 2023. 1, 2, 5, 6
- [59] Samy Tafasca, Anshul Gupta, and Jean-Marc Odobez. Sharingan: A transformer architecture for multi-person gaze following. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2008–2017, 2024. 1, 2, 3, 5, 6
- [60] Claudia Thoermer and Beate Sodian. Preverbal infants’ understanding of referential gestures. *First Language*, 21(63):245–264, 2001. 1
- [61] Henri Tomas, Marcus Reyes, Raimarc Dionido, Mark Ty, Jonric Mirando, Joel Casimiro, Rowel Atienza, and Richard Guinto. Goo: A dataset for gaze object prediction in retail environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 3125–3133, 2021. 5
- [62] Francesco Tonini, Cigdem Beyan, and Elisa Ricci. Multi-modal across domains gaze target detection. In *Proceedings of the 2022 International Conference on Multimodal Interaction*, pages 420–431, 2022. 1, 2, 6
- [63] Francesco Tonini, Nicola Dall’Asen, Cigdem Beyan, and Elisa Ricci. Object-aware gaze target detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 21860–21869, 2023. 3, 4, 5, 6
- [64] Danyang Tu, Xiongkuo Min, Huiyu Duan, Guodong Guo, Guangtao Zhai, and Wei Shen. End-to-end human-gaze-target detection with transformers. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2192–2200. IEEE, 2022. 4, 3, 5
- [65] Danyang Tu, Wei Shen, Wei Sun, Xiongkuo Min, and Guangtao Zhai. Joint gaze-location and gaze-object detection. *arXiv preprint arXiv:2308.13857*, 2023. 3, 4, 5
- [66] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3
- [67] Binglu Wang, Tao Hu, Baoshan Li, Xiaojuan Chen, and Zhi-jie Zhang. Getecto: A unified framework for gaze object prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19588–19597, 2022. 2
- [68] Ping Wei, Yang Liu, Tianmin Shu, Nanning Zheng, and Song-Chun Zhu. Where and why are they looking? jointly inferring human attention and intentions in complex tasks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6801–6809, 2018. 1
- [69] Zeynep Yücel and Albert Ali Salah. Resolution of focus of attention using gaze direction estimation and saliency computation. In *2009 3rd International Conference on Affective Computing and Intelligent Interaction and Workshops*, pages 1–6. IEEE, 2009. 3
- [70] 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. 3
- [71] Hao Zhao, Ming Lu, Anbang Yao, Yurong Chen, and Li Zhang. Learning to draw sight lines. *International Journal of Computer Vision*, 128:1076–1100, 2020. 1, 2

# Gaze-LLE: Gaze Target Estimation via Large-Scale Learned Encoders

## Supplementary Material

### Table of Contents

- 6 Integration of DINOv2 into Existing Methods
- 7 Experiment Details for Section 3.2
- 8 Comparison to Detection Methods
- 9 Runtime Analysis
- 10 Performance with Estimated Head Bounding Boxes
- 11 Reimplementation of Horanyi et al.
- 12 Additional Ablation Studies
- 13 LoRA Backbones
- 14 Additional Visualizations & Failure Modes

### 6. Integration of DINOv2 into Existing Methods

In this section, we provide further details on our experiments in Table 1, which integrate DINOv2 into three existing methods: Chong et al. [9], Miao et al. [42], and Gupta et al. [22].

**Chong et al.** Chong et al. [9]’s method employs separate head and scene encoders, each of which is composed of a ResNet50 + 1 additional ResNet layer. The input to the head branch is a  $224 \times 224$  crop of the head and the input to the scene branch is the  $224 \times 224$  scene image concatenated channel-wise with a binary map of the person’s head bounding box position. The output of each encoder is a  $1024 \times 7 \times 7$  feature map (channels  $\times$  height  $\times$  width). For our experiments, we replace the scene encoder with a ViT-Base DINOv2 encoder. Because the DINOv2 encoder produces a  $768 \times 16 \times 16$  feature map, we apply average pooling with kernel size=3 and stride=2 followed by a convolutional layer with kernel size=1 and stride=1 to transform the feature map to the model’s expected size of  $1024 \times 7 \times 7$ . We follow the rest of the existing method, which consists of an attention mechanism to re-weight the scene features based on the head features and head position, concatenation of the head and scene features, 2 convolutional encoding layers, and a 4-layer convolutional decoder. We consider 3 learning settings for the DINOv2 encoder:

1. **Frozen:** We simply replace the scene encoder with the DINOv2 encoder and freeze it during training. Because the DINOv2 takes in a 3-channel RGB image, we do not concatenate the head position map to the input as in the original method.
2. **Frozen + proj:** We alter the DINOv2 encoder’s patch projection layer to take in 4 channels so that the input to the scene encoder is the concatenated RGB image and head position map like in the original method. We freeze all weights of the DINOv2 during training *except* the

DINOv2 Training	Learning rate	AUC $\uparrow$	Avg L2 $\downarrow$	Min L2 $\downarrow$
Original Method	2.5e-4	<b>0.921</b>	<b>0.137</b>	<b>0.077</b>
Frozen	2.5e-4	0.858	0.196	0.133
	1.0e-4	0.857	0.201	0.145
	1.0e-5	0.808	0.230	0.166
	1.0e-6	0.726	0.287	0.218
Frozen + proj	2.5e-4	0.875	0.191	0.125
	1.0e-4	0.872	0.198	0.129
	1.0e-5	0.850	0.212	0.143
	1.0e-6	0.766	0.282	0.208
Trained + proj	2.5e-4	0.876	0.185	0.120
	1.0e-4	0.908	0.167	0.101
	1.0e-5	0.870	0.199	0.132
	1.0e-6	0.805	0.260	0.187

Reason why  
out of the box  
pretrained Dino  
alone is not  
enough

Table 9. Comparison of integrating DINOv2 into Chong et al. [9] with different training configurations (DINOv2 encoder learning strategy & learning rate) on GazeFollow.

patch projection layer.

3. **Trained + proj:** We include the altered 4-channel patch projection layer and train the full DINOv2 encoder during training.

Tab. 9 shows our results from trying different training strategies for the DINOv2 encoder and different learning rates. We see that learning the projection layer to integrate head position as an input to the scene encoder has a significant performance gain over using the DINOv2 with RGB-only inputs, and that training the DINOv2 fully performs best. Importantly, we do not observe overfitting - the trained results are better than using the frozen DINOv2. For this method, regular training outperforms LoRA. However, all results using a DINOv2 encoder in place of the ResNet50-based scene encoder perform worse than the original method.

Regular  
training if  
possible,  
will always  
outperform  
LoRA

**Miao et al.** Miao et al. [42] is a more recent work that expands upon Chong et al.’s architecture by integrating estimated depth into the scene encoding and feature fusion, a global attention mechanism over the scene prior to decoding, and an additional patch-level training objective. Similar to Chong et al., Miao et al. employ head and scene encoders composed of a ResNet50 + 1 additional ResNet layer. The input to the scene branch is a 5-channel concatenation of the RGB scene image, the binary head position map, and an estimated depth map from MiDaS [49]. Like with Chong et al., we replace the scene encoder with a DINOv2 encoder, and use average pooling and a convolutional layer to transform the scene feature map to size  $1024 \times 7 \times 7$ . We consider the same training configurations as we did with Chong et al.,

DINOv2 Training	Learning rate	AUC ↑	Avg L2 ↓	Min L2 ↓
Original Method	2.5e-4	<b>0.934</b>	<b>0.123</b>	<b>0.065</b>
Frozen	2.5e-4	0.858	0.207	0.141
	1.0e-4	0.859	0.203	0.138
	1.0e-5	0.807	0.236	0.169
	1.0e-6	0.702	0.297	0.228
Frozen + proj	2.5e-4	0.892	0.173	0.109
	1.0e-4	0.887	0.176	0.113
	1.0e-5	0.859	0.203	0.137
	1.0e-6	0.761	0.286	0.213
Trained + proj	2.5e-4	0.899	0.165	0.103
	1.0e-4	0.910	0.152	0.093
	1.0e-5	0.900	0.161	0.098
	1.0e-6	0.847	0.220	0.149

Table 10. Comparison of integrating DINOv2 into Miao et al. [42] with different training configurations (DINOv2 encoder learning strategy & learning rate) on GazeFollow.

DINOv2 Training	Learning rate	AUC ↑	Avg L2 ↓	Min L2 ↓
Original Method	2.5e-4	<b>0.933</b>	<b>0.134</b>	<b>0.071</b>
Frozen + proj	2.5e-4	0.893	0.180	0.113
	1.0e-3	0.894	0.184	0.116
	1.0e-4	0.897	0.175	0.108
	1.0e-5	0.874	0.199	0.129
	1.0e-6	0.818	0.228	0.161
Trained + proj	2.5e-4	0.908	0.165	0.099
	1.0e-3	0.912	0.155	0.091
	1.0e-4	0.911	0.159	0.095
	1.0e-5	0.899	0.167	0.101
	1.0e-6	0.842	0.219	0.149

Table 11. Comparison of integrating DINOv2 into Gupta et al. [22] (Image-only variant) with different training configurations (DINOv2 encoder learning strategy & learning rate) on GazeFollow.

however we change the learned patch projection to have 5 input channels to account for Miao et al.’s inclusion of depth as input. As shown in Tab. 10, we achieve the best results by fully training the DINOv2. However, all configurations still perform worse than the original method with the ResNet50 backbone.

**Gupta et al.** Gupta et al. [22]’s approach consists of of a head-centric module, scene-centric module, and heatmap decoder. The head-centric module is a ResNet18 encoder which is supervised to predict 3D gaze from the head crop. This 3D gaze prediction is processed along with the head location into spatial gaze cone, which is passed to the scene-centric module along with the image. The scene-centric module consists of a separately trained EfficientNet encoder from different scene modalities: image, predicted depth, or predicted pose. Optionally, the encoders for the different modalities may be used together with a learned weighted attention module for fusion. As the training process calls for

Method	Input size	AUC	Avg L2	Min L2
Chong et al. - Original	224	0.921	0.137	0.077
Chong et al. - Original	448	0.923	0.138	0.076
Chong et al. - Trained DINOv2	224	0.908	0.170	0.101
Chong et al. - Trained DINOv2	448	0.897	0.169	0.105
Miao et al. - Original	224	0.934	0.123	0.065
Miao et al. - Original	448	0.923	0.151	0.086
Miao et al. - Trained DINOv2	224	0.910	0.152	0.093
Miao et al. - Trained DINOv2	448	0.908	0.154	0.094
Gupta et al. - Original	224	0.943	0.114	0.056
Gupta et al. - Original	448	0.939	0.108	0.052
Gupta et al. - Trained DINOv2	224	0.912	0.155	0.091
Gupta et al. - Trained DINOv2	448	0.908	0.170	0.103

Table 12. Effect of increasing the input scene image size for Chong et al., Miao et al., and Gupta et al.’s original methods and best variants with DINOv2. We do not observe clear gains from using a larger input size.

separately training each modality, we consider the image-only variant for our DINOv2 integration experiments. We replace the EfficientNet-B1 image encoder with DINOv2, and add an additional learned projection layer to reduce the dimension from DINOv2’s output dimension of 768 to the model’s internal dimension of 64. We consider both training the full encoder and freezing the encoder (with the exception of the input projection, which must accept the extra gaze cone channel). We report performance in Tab. 11. Like the other methods, training the DINOv2 performs better than freezing it, but still underperforms compared to the original method.

**Input Size** Because we do not include a separate head branch that operates on a higher-resolution crop of the head in our main method, we use an input size of  $448 \times 448$  instead of  $224 \times 224$  like these prior works. To validate that our method’s gains are not only a result of the larger input to the scene encoder, we retrain Chong et al., Miao et al.’s, and Gupta et al.’s original methods as well as the best variant with a DINOv2 scene encoder with scene input size  $448 \times 448$  in Tab. 12. For Gupta et al.’s original method, we use their full multimodal model. We perform average pooling on the resultant scene feature maps when necessary to reduce the spatial dimensions to the expected shape for compatibility with the rest of the model. For Chong et al.’s method, the results are largely the same between using 224 vs. 448, while for Miao et al., using 448 actually decreases performance. For Gupta et al.’s architecture, increasing the resolution to 448 results in worse AUC, which is the primary metric on GazeFollow, but achieves slight gains on the L2 metrics. We thus do not see clear improvements from using an increased input size, illustrating that a larger scene input size is not necessary when a high-resolution head crop is already provided to the model.

Yes, increasing resolution only helps when there are so minute details in the picture that gets squashed by a low resolution



Transformer Decoder
Linear ( $d \rightarrow 256$ )
Trans. Layer (dim=256, heads=8, mlp_dim=1024)
ConvT(256 → 256, k=2, s=2)
Conv(256 → 1, k=1, s=1)
Sigmoid

Conv Decoder
Conv ( $d \rightarrow 768$ , k=1, s=1)
Conv(768 → 384, k=1, s=1)
Conv(384 → 192, k=2, s=2)
ConvT(192 → 96, k=2, s=2)
ConvT(96 → 1, k=2, s=2)
Conv(1 → 1, k=1, s=1)
Sigmoid

Table 13. Architecture details for Transformer Decoder and Convolutional Decoder for experiments in Section 3.1

## 7. Experiment Details for Section 3.2

In this section, we provide further details about our experiments in Section 3.2 that investigate early vs. late head position integration, transformer vs. convolutional decoding, and head & scene branch vs. scene-branch only design.

**Scene & Head Backbones** We use a frozen DINOv2 ViT-Base backbone for both the scene branch and the head branch. For the scene branch, we use input size  $448 \times 448$ , yielding a feature map  $x_{\text{scene}} \in \mathbb{R}^{768 \times 32 \times 32}$ . Because the head occupies a smaller portion of the full-resolution image, we use input size  $224 \times 224$  for the head branch and upsample the resulting feature map to  $x_{\text{head}} \in \mathbb{R}^{768 \times 32 \times 32}$  so it can be concatenated with  $x_{\text{scene}}$ . We concatenate  $x_{\text{scene}}$  and  $x_{\text{head}}$  channel-wise to form the combined features  $x \in \mathbb{R}^{1536 \times 32 \times 32}$ . For the scene-only variant, we set  $x = x_{\text{scene}} \in \mathbb{R}^{768 \times 32 \times 32}$ .

**Head Position Integration** For “early” integration of the head position, we change the patch projection layer of the DINOv2 scene branch to have 4 input channels (RGB + binary head position map) instead of 3. During training, we learn this patch projection layer while keeping the rest of the DINOv2 frozen. For “late” integration, we do not alter or train the projection layer. Instead, we downsample the binary head position map to size  $1 \times 32 \times 32$  and concatenate it with  $x$  to form  $x' \in \mathbb{R}^{d \times 32 \times 32}$  where  $d = 1537$  or  $d = 769$  depending on the inclusion of the head branch. For “early” integration, we do not concatenate the head position map, so  $x' = x \in \mathbb{R}^{d \times 32 \times 32}$  where  $d = 1536$  (head & scene branch) or  $d = 768$  (scene branch only).

**Decoder** We provide architecture details for the transformer and convolutional heatmap decoders in Tab. 13. Each produce a  $64 \times 64$  gaze heatmap from  $x'$ . The convolutional decoder is based on the network design used by Chong et al. [9] and several subsequent methods, consisting

of 6 convolutional layers (each followed by batch normalization and a ReLU activation) to progressively project the feature map to a smaller dimension while upscaling it to the output heatmap size. The transformer decoder consists of a single transformer layer of dimension 256 followed by 2 shallow convolutional layers. Both decoders have approximately the same number of learned parameters (1.85M for the scene-branch only model with late head position integration).

**Training Details** We train the models on GazeFollow for 15 epochs using the Adam optimizer, cosine scheduling with initial learning rate 1e-3, and batch size 60. We use the same data augmentations during training that we use in our main experiments (random crop, flip, and bounding box jitter).

## 8. Comparison to Detection Methods

A set of recent works formulate gaze target estimation as a set detection problem, jointly predicting a set of head bounding boxes and their corresponding gaze locations [63–65]. We exclude these works from our main comparisons in Section 4.2 due to differences in the evaluation setting, as these methods perform bipartite matching *using the ground truth gaze targets at test time*. In this section, we provide further details about the difference in evaluation setting, and provide quantitative comparison with Tonini et al. [63] in our setting using their open source codebase.

**Formulation** Tu et al. [64] provided the first set detection formulation for joint head and gaze target detection by proposing HGTTR, a DETR [3]-based transformer detection framework. Given an image  $x_{\text{img}} \in \mathbb{R}^{3 \times H_{\text{in}} \times W_{\text{in}}}$ , HGTTR predicts a fixed number of  $N$  human-gaze instances, where each instance  $y$  is composed of a head bounding box prediction  $y_{\text{bbox}} \in [0, 1]^4$ , a binary classification score  $y_{\text{class}} \in [0, 1]$  indicating the probability that the instance is indeed a head, a prediction of if the gaze is in or out of frame  $y_{\text{in/out}} \in [0, 1]$ , and a gaze heatmap  $y_{\text{heatmap}} \in [0, 1]^{H_{\text{out}} \times W_{\text{out}}}$ . Notably, the difference between this setting and the traditional problem formulation (which we follow) is that a head bounding box is not given as input. Instead, the model predicts all head bounding boxes along with their associated gaze target as output.

**Matching Algorithm** Like DETR, HGTTR uses the Hungarian algorithm [33] to determine a one-to-one mapping between the predicted instances and the ground truth instances in order to calculate loss at train time. The optimal matching is found by considering all possible mappings  $w$  between the predicted instances and ground truth instances and selecting the one that minimizes

$$\mathcal{L}_{\text{cost}} = \sum_{i=1}^N \mathcal{L}_{\text{match}}(y_i, \hat{y}_{w(i)}) \quad (4)$$

where  $\mathcal{L}_{\text{match}}(y_i, \hat{y}_{w(i)})$  is a pairwise matching cost function between the  $i$ -th ground truth instance  $y_i$ , and the predicted instance with index  $w(i)$ ,  $\hat{y}_{w(i)}$ . In HGTTR,  $\mathcal{L}_{\text{match}}$  is defined as a weighted sum of loss functions:

$$\lambda_1 \mathcal{L}_{\text{bbox}} + \lambda_2 \mathcal{L}_{\text{class}} + \lambda_3 \mathcal{L}_{\text{in/out}} + \lambda_4 \mathcal{L}_{\text{heatmap}} \quad (5)$$

where  $\mathcal{L}_{\text{bbox}}$  is an IoU loss on the predicted bounding box,  $\mathcal{L}_{\text{class}}$  is a binary classification loss on predicting if the instance is a head or not,  $\mathcal{L}_{\text{in/out}}$  is a binary classification loss on predicting if the gaze is in or out of frame, and  $\mathcal{L}_{\text{heatmap}}$  is heatmap loss between the predicted and ground truth gaze target. The model always predicts  $N$  instances, with  $N$  being chosen to exceed the number of ground truth instances present in each image in the dataset (all existing methods use  $N = 20$ , which is significantly larger than the typical number of people present in a single image in GazeFollow). Because there are always less ground truth instances than predicted instances, the ground truth instance list is padded with  $\emptyset$  so that it is length  $N$ . Predicted instances mapped to  $\emptyset$  are excluded from cost and loss calculation. See Tu et al. [64] and DETR (Carion et al.) [3] for further details on matching.

Tonini et al. [63] expand upon this formulation, training their model to also predict all objects in the scene as an auxiliary training objective, and including depth as an input. They also add a term  $y_{\text{vector}}$  to each instance, which is a predicted gaze vector, and use this as auxiliary supervision. Their matching cost is defined as:

$$\lambda_1 \mathcal{L}_{\text{bbox}} + \lambda_2 \mathcal{L}_{\text{class}} + \lambda_3 \mathcal{L}_{\text{in/out}} + \lambda_4 \mathcal{L}_{\text{heatmap}} + \lambda_5 \mathcal{L}_{\text{vector}} \quad (6)$$

For all approaches, the mapping between the ground truth and predicted instances is determined by finding the closest subset of predicted instances to the ground truth based on bounding box, class, and gaze.

**Evaluation Setting** At inference time, these methods use the same matching cost  $\mathcal{L}_{\text{cost}}$  to determine which predicted instances are evaluated against which ground truth instances, and use this to calculate the gaze performance metrics (e.g. heatmap AUC, L2 distance). This is inherently a different evaluation setting than ours because *the ground truth gaze labels are used at inference time* to retrieve the predicted instances that are compared against the ground truth instances. Because the fixed number of predicted instances,  $N = 20$ , is much higher than the typical number of ground truth instances per image, a model can predict multiple instances with the same head bounding box, but different gaze targets (see Fig. 6 for visual examples of this). In this case, the matching algorithm will match each ground truth instance to the predicted instance with the closest gaze and calculate the gaze metrics between these pairs,

Method	AUC $\uparrow$	Avg L2 $\downarrow$	Min L2 $\downarrow$
<i>with ground truth gaze matching</i>			
Tu et al. [64]	0.917	0.133	0.069
Tu et al. [65]	0.928	0.114	0.057
Tonini et al. [63]	0.922	0.069	0.029
Tonini et al.* [63]	0.924	0.068	0.030
<i>no ground truth gaze matching</i>			
Tonini et al.* [63]	0.767	0.211	0.148
Ours	0.956	0.104	0.045

Table 14. Quantitative comparison with detection-based methods on GazeFollow. The results *with ground truth gaze matching* use the ground truth gaze labels to perform bipartite matching at test time, and thus are not a direct comparison to our method and prior work. The *no ground truth gaze matching* results report our method compared to Tonini et al.’s model evaluated with the altered matching cost function in Equation 7, which excludes ground truth gaze information. (\*Results we obtained ourselves by running Tonini et al.’s published code.)

discarding the extra incorrect predictions from evaluation. In this way, the gaze performance metrics alone do not penalize overdetection. They characterize recall by assessing the accuracy of the predicted instances that are closest to the ground truth, but do not assess precision by penalizing the model for predicting additional instances with incorrect gaze targets. Additionally, the matching algorithm does not enforce that the ground truth instance is matched to a detection with a similar predicted head bounding box; if the heatmap loss dominates the matching cost, an instance may be selected based only on similarity between the predicted gaze heatmap and ground truth (see Fig 6 for examples). Thus, the model does not need to correctly associate people with their respective gaze targets to achieve high performance. For these reasons, the gaze metrics in this evaluation setting are not a direct comparison against our work and prior methods that follow the traditional problem formulation.

**Quantitative Results** We show the reported results of the 3 detection-based methods and our results on GazeFollow in Table 14. To quantitatively characterize the difference in evaluation setting, we also re-evaluate Tonini et al.’s [63] method on GazeFollow with the ground truth gaze label removed from the matching cost, using their published codebase. We alter the matching cost from Equation 6 to exclude the ground truth gaze label. This altered matching cost is defined as:

$$\mathcal{L}'_{\text{match}} = \lambda_1 \mathcal{L}_{\text{bbox}} + \lambda_2 \mathcal{L}_{\text{class}} \quad (7)$$

With this cost, the model retrieves a prediction for each ground truth instance based only on bounding box overlap and class similarity. This reflects our use case, where the model is used to predict a gaze target for a certain per-

son based on their head location, and does not have knowledge of the ground truth gaze. Without the use of ground truth gaze in the matching cost at inference time, we observe a significant performance drop. This quantitatively indicates the overdetection of gaze instances, as the altered matching cost results in the model selecting detections that have more bounding box overlap and class similarity<sup>2</sup> to the ground truth, but a less accurate gaze target. However, it is important to note that we do not exhaustively attempt to adapt their method to this setting (e.g. by developing a new matching algorithm for training or a non-maximal suppression method). We include this result to quantitatively demonstrate the difference in evaluation setting and use case between our method and their method as-is. We note that Tu et al. [64, 65] do not publish code or models so we do not re-evaluate their methods.

**Qualitative Results** We visualize the output instances of Tonini et al.’s default matching algorithm that uses ground truth gaze as part of the cost function, and our altered matching algorithm that does not use ground truth gaze in Figure 6. The first two rows demonstrate cases where the matching algorithm chooses an instance with a predicted bounding box that is not associated with the correct person; the heatmap loss dominates the matching cost. With our altered matching function, an instance with a predicted bounding box for the correct person but incorrect gaze heatmap is retrieved. The third row shows an example of overdetection, where multiple instances are predicted with a head bounding box for the correct person, but different gaze targets. With the use of ground truth gaze during matching, the instance with the most correct heatmap is selected. However, without this ground truth information, the model does not select the best instance and produces an incorrect gaze prediction. These examples visually illustrate the difference in evaluation setting: when ground truth gaze information is used at test time, a model can achieve high performance by producing instances that capture different potential gaze targets and relying on the matching algorithm to retrieve the best instances to evaluate with. However, the gaze metrics do not characterize the model’s ability to determine which of these instances are indeed gaze targets and associate them with the correct person.

## 9. Runtime Analysis

**Inference Speed** Our ViT-Base model runs in 15ms ( $\approx$  66fps) on a single NVIDIA RTX 4090 GPU. We compare the inference time of our model with existing methods in Fig 7a. For Miao et al. [42], we include the auxiliary

<sup>2</sup>We observe that matching is mainly based on bounding box overlap. Changing the weight of class similarity in the matching cost has little effect on performance both in the original setting and our altered setting where gaze is not used in matching.

depth estimation model (DPT-Hybrid[49]) in runtime calculation. Compared to Miao et al., our approach is both faster, and achieves better performance. In fact, the inference time of the DPT-Hybrid depth model (17ms) exceeds the entire inference time of our approach. This result highlights the benefit of using a single encoder, both in inference speed and performance. Chong et al.’s approach [9], which does not use any models for auxiliary modalities like depth, runs faster than our model. However, this comes with a significant drop in performance compared to our method. As shown in Table 3, recent convolutional methods all use at least one auxiliary model to augment performance. While these approaches may use faster backbones than a ViT, requiring auxiliary models ultimately increases runtime.

**Multi-Person Scaling** We also investigate how our model’s runtime scales with estimating the gaze for multiple people per image (Fig. 7b). We measure the inference time for 1-10 people per image for both our default method, and our variant that uses a head position token ( $t_{pos}$ ) and decodes gaze via cross attention and a dot product with the scene features (Table 7 configuration 1b 2b). Because the majority of our model’s computation can be attributed to the DINOv2 scene encoder (>95% of computation), which is run once regardless of the number of people, our model’s runtime does not increase much with addition of more people (15ms for 1 person vs. 19ms for 10 people). The token variant of our model with cross attention scales even better, as it decodes gaze for all people from the same final feature map. However, as shown in Fig. 7a, this is accompanied by a slight performance decrease.

We include Tonini et al.’s [63] detection method for comparison, which is designed to simultaneously predict the gaze and bounding boxes for all people in an image and thus has a constant runtime across different numbers of people. We include both the 2D variant (which does not use depth), and the 3D variant (which uses depth), accounting for the inference time for a DPT-Hybrid depth estimator for the 3D variant. Because our model requires head bounding boxes, we include a YOLOv5 head detector in the displayed runtimes for our model. We observe that our default method is faster than Tonini et al.’s 2D method for up to 7 people, and our token variant with cross attention is faster for all numbers of people. Due to the inclusion of running the depth model and modeling differences to include depth, the 3D version of Tonini et al.’s method is slower than the 2D version and our method.

## 10. Performance with Estimated Head Bounding Boxes

Tu et al. [64] report that 2-stream methods suffer major performance drops when using head bounding boxes from a detector rather than the dataset ground truth. In contrast,

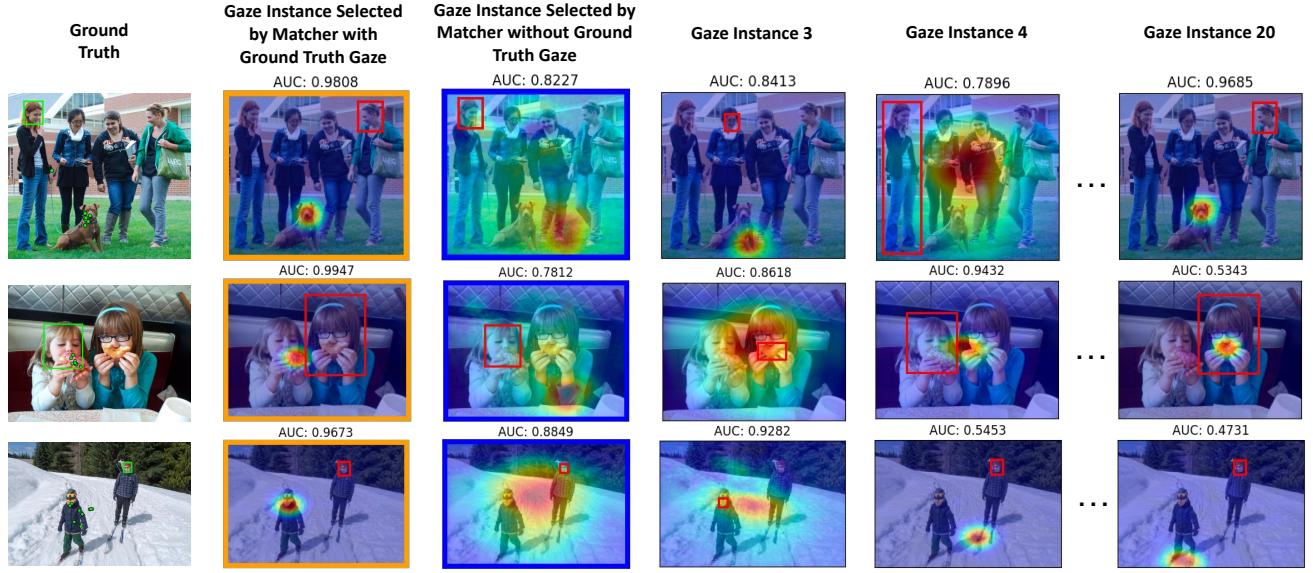


Figure 6. We show the output gaze instances (predicted head bounding box & gaze heatmap) from Tonini et al.’s model [63] for 3 examples. We identify the instances selected by Tonini et al.’s matching cost (which uses the ground truth gaze) and our altered matching cost (which excludes ground truth gaze and instead performs matching based on bounding box overlap). Tonini et al.’s matching algorithm selects the instance with the closest gaze prediction to the ground truth, but the bounding box prediction does not always correspond to the correct person (Rows 1-2). Additionally, we observe *overdetection*, where the algorithm predicts multiple instances for the same person with different gaze heatmaps (Row 3). Without the use of ground truth gaze information, the model cannot determine which of these instances is best.

Method	AUC $\uparrow$	Avg L2 $\downarrow$	Min L2 $\downarrow$
ViT-B + GT	0.956	0.104	0.045
ViT-B + YOLO	0.955	0.106	0.047
ViT-L + GT	0.958	0.099	0.041
ViT-L + YOLO	0.958	0.101	0.043

Table 15. Gaze-LLE achieves consistent results when using head detections from an out-of-the-box YOLOv5 detector instead of head ground truth bounding boxes.

we observe almost no performance degradation when pairing our method with a YOLOv5 head detector trained on CrowdHuman [1, 53] (Tab. 15). This result demonstrates that our single-stream design, which uses a coarse, down-sampled head position map, is less dependent on an exact head crop, and works well with out of the box head detections.

## 11. Reimplementation of Horanyi et al.

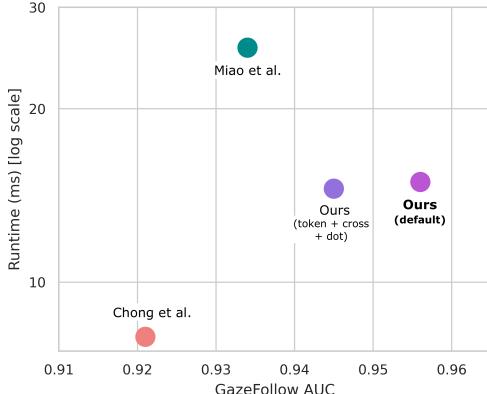
We use our own implementation of Horanyi et al. [25] for our main comparison. We choose to reimplement this method because the reported results are outliers among other methods, and there is imbalance between the reported metrics (*e.g.*, 0.932 AUC on GazeFollow, but very low L2 error). Since the method is largely constructed from ele-

Experiment	GazeFollow			VideoAttentionTarget		
	AUC $\uparrow$	Avg L2 $\downarrow$	Min L2 $\downarrow$	AUC $\uparrow$	L2 $\downarrow$	AP <sub>in/out</sub> $\uparrow$
Frozen Aux. Angle	0.869	0.217	0.146	0.802	0.234	0.720
Trained Aux. Angle	<b>0.896</b>	<b>0.196</b>	<b>0.127</b>	<b>0.832</b>	<b>0.199</b>	<b>0.800</b>

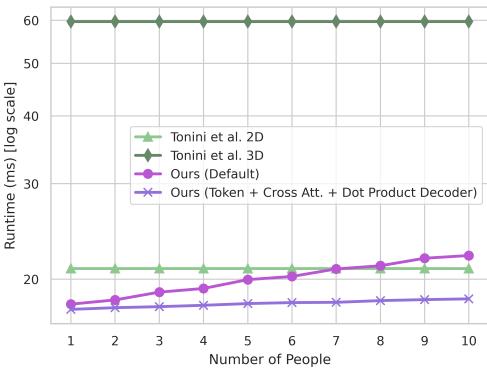
Table 16. Experimental results for our implementation of Horanyi et al.[25] on GazeFollow and VideoAttentionTarget. We consider the setting where we freeze the auxiliary 3D gaze angle model vs. where we train it along with the rest of the network.

ments that are present in other works (*e.g.*, constructing a “gaze cone” from estimated 3D gaze and depth [16, 22, 63], providing estimated depth as input to the scene encoder [16, 22], and using a ResNet50-based scene encoder + 4-layer convolutional decoder [9, 42]), it is difficult to identify the source of large reported performance gains. There is not published code for this work. The original paper provides limited implementation details, so we follow some choices from Chong et al.’s codebase [9].

The model consists of a 3D Field-of-View (FoV) map construction module, a scene encoder, and a gaze decoder. The FoV module uses an auxiliary depth estimator and 3D gaze angle estimator to produce an FoV heatmap for a person over the scene. The estimated depth, FoV map, and  $224 \times 224$  map are passed to a ResNet50-based scene en-



(a) Runtime vs. Performance



(b) Runtime scaling for multi-person inference

Figure 7. Runtime analysis of our approach: we show the tradeoff of inference time vs. performance (7a), and analyze how different variants of our approach paired with a head detector scale for multi-person prediction, compared to detection methods (7b). All experiments are performed on a single NVIDIA RTX 4090 GPU.

coder and decoded into gaze predictions. Figure 8 illustrates the architecture details for our reimplementation. We use the same auxiliary models used in the original approach [25]: Gaze360 [31] and Monodepth2 [20]. We follow the version of their scene encoder without non-local (NL) blocks. The FoV module uses the construction equation from Horanyi et al. [25]:

$$M_{ind} = \text{min\_max\_scaler} \left( \frac{(i - h_x, j - h_y, k - h_z) \cdot (g_x, g_y, g_z)}{\|i - h_x, j - h_y, k - h_z\|^2 \cdot \|g_x, g_y, g_z\|} \right) \quad (8)$$

We make the assumption that  $k$ -coordinate comes from the normalization of the estimated depth map. We follow the high-level architecture described in the text: a ResNet50 trainable scene encoder, two convolutions for encoding, and a 4-layer convolutional decoder. Because details such as hidden dimensions and kernel sizes are not specified, we generally follow Chong et al.’s open-source code [9] since Horanyi et al.’s described architecture mostly matches

$d_{\text{model}}$	Params	AUC $\uparrow$	Avg L2 $\downarrow$	Min L2 $\downarrow$
128	1.2M	<b>0.956</b>	0.106	0.046
256 (default)	2.8M	<b>0.956</b>	<b>0.104</b>	<b>0.045</b>
384	5.0M	<b>0.956</b>	0.105	0.046
512	7.7M	0.953	0.108	0.049
768	14.8M	0.953	0.108	0.049

(a) Dimension of gaze estimation module.

Layers	Params	AUC $\uparrow$	Avg L2 $\downarrow$	Min L2 $\downarrow$
1 layer	1.2M	0.953	0.115	0.054
2 layers	2.0M	0.955	0.108	0.049
3 layers	2.8M	<b>0.956</b>	0.104	<b>0.045</b>
4 layers	3.6M	<b>0.956</b>	<b>0.103</b>	<b>0.045</b>
5 layers	4.4M	<b>0.956</b>	0.104	<b>0.045</b>

(b) Number of transformer layers.

Table 17. We investigate the effect of different internal model dimensions and number of transformer layers for our gaze estimation module with a ViT-Base DINOv2 backbone. We observe diminishing returns as we increase the dimension and number of layers. We select  $d_{\text{model}} = 256$  with 3 transformer layers as our default configuration.

Chong et al.’s. We conduct experiments in two settings on the GazeFollow and VideoAttentionTarget datasets. The first setting keeps both the auxiliary gaze angle and depth estimation models frozen, as suggested in the text [25]. In the second setting, we train the gaze angle model. For the GazeFollow experiments, we use batch size 128, learning rate 4e-4, and the Adam optimizer. For VideoAttentionTarget, we finetune the GazeFollow-trained model with batch size 32 and learning rate 1e-4. The results of these experiments are shown in Tab. 16. We observe training the auxiliary gaze angle model performs better, so we report these results in the main paper.

## 12. Additional Ablation Studies

We provide additional ablations for our ViT-Base model on GazeFollow in Tab. 17. We find there is benefit to using a smaller internal dimension for our gaze estimation module, both in performance and reduction of learnable parameters (Tab. 17a). Our model produces competitive results with prior work using only 1 transformer layer (Tab. 17b); however, we achieve sizeable performance gains by increasing the number of layers to 3. Beyond 3 layers, the performance is largely stable. To balance performance with reducing learnable parameters, we select dimension 256 with 3 layers as our default configuration.

## 13. LoRA Backbones

To investigate if training the backbone improves performance, we explore using Low Rank Adaptation (LoRA) [26] on GazeFollow in Table 19. LoRA updates the backbone while introducing limited additional learnable param-

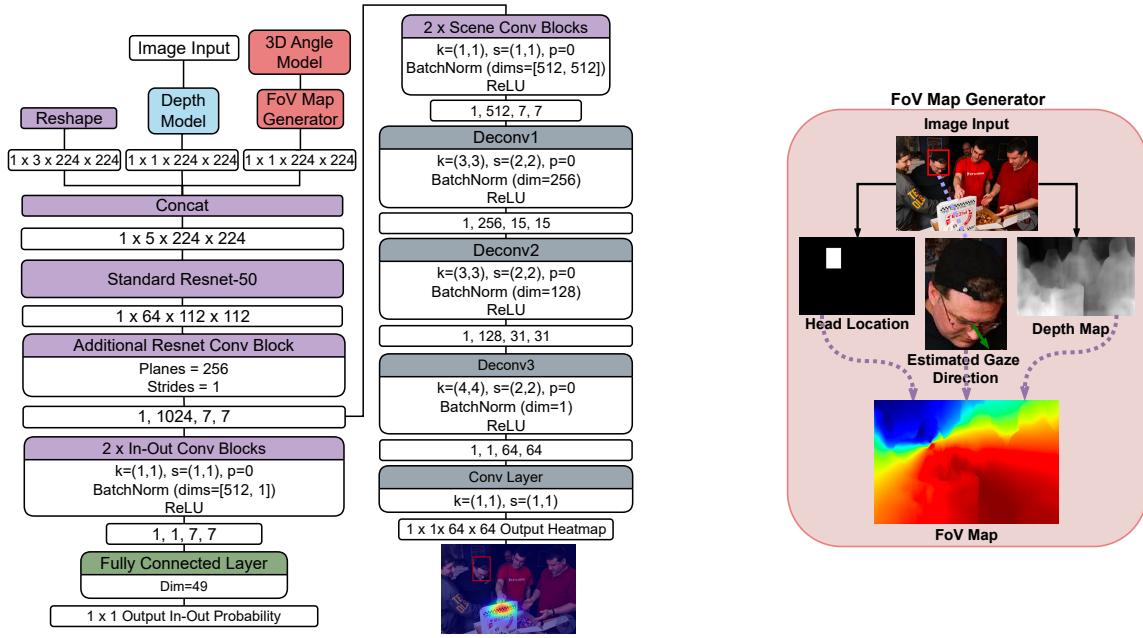


Figure 8. Architecture details for our reimplementation of Horanyi et al.’s model [25]. The model consists of a FoV Map Generator (shown on right), which uses an auxiliary 3D gaze angle estimator and an auxiliary depth model to produce an FoV map for a given person. The FoV map, estimated depth, and image are passed to a ResNet50-based encoder and convolutional decoder to produce a gaze prediction. In our experiments, we consider both freezing vs. training the 3D gaze angle estimator as part of the model.

Backbone	Params	AUC $\uparrow$	Avg L2 $\downarrow$	Min L2 $\downarrow$
<i>One Human</i>		0.924	0.096	0.040
ViT-B	2.8M	0.956	0.104	0.045
ViT-B + LoRA	3.1M	0.957	0.103	0.045
ViT-L	2.9M	0.958	0.099	0.041
ViT-L + LoRA	3.7M	0.960	0.097	0.040

Table 18. LoRA-tuned DINOv2 Backbones

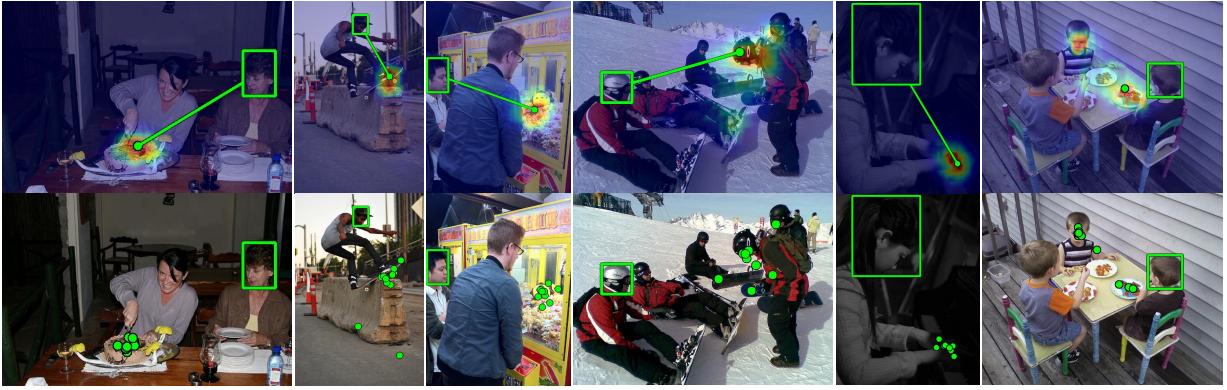
Table 19. Frozen vs. LoRA-tuned DINOv2 backbones with GazeLLE on GazeFollow.

ters by learning weight update matrices as low rank decompositions. We update the query and value projections of the DINOv2 backbone using rank 16. We observe limited improvements, which we attribute to (1) the effectiveness of the frozen encoder’s feature representation for our task and (2) that our models with frozen encoders already achieve extremely close performance to the inter-rater performance of the human annotators, which serves as a soft upper bound on the L2 metrics.

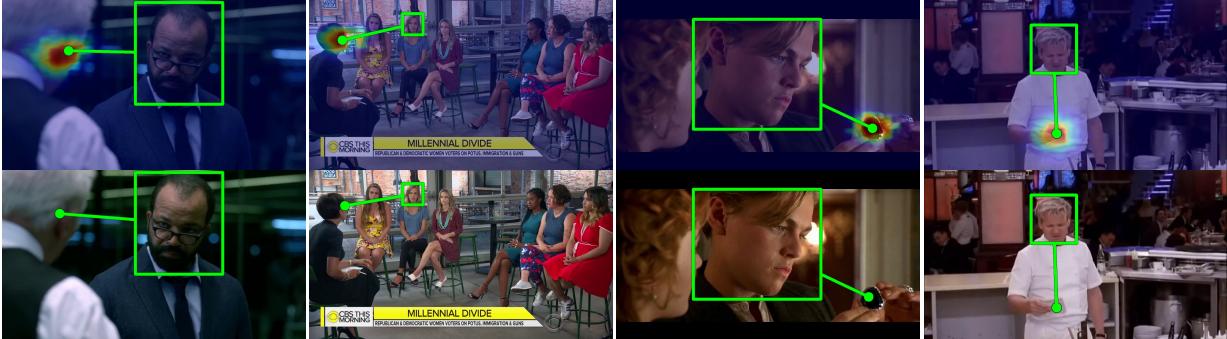
## 14. Additional Visualizations & Failure Modes

We provide additional visualizations of our ViT-B model’s predicted heatmaps on the GazeFollow, VideoAttentionTarget, ChildPlay, and GOO-Real datasets in Figure 9. We

show examples where our model does not perform as well in Figure 10. These cases are representative of error modes we observe across the evaluation datasets. Our model is more likely to exhibit errors when the person is positioned with the back of their head towards the camera (examples 1-2) or their face is heavily occluded (example 3). In these cases, we observe our model selects potential targets (such as faces) that are broadly in person’s field of view, but does not always refine this prediction to the ground truth gaze target. It is not surprising that the model does not perform as well on these cases, as the ground truth is often inherently more ambiguous in such examples. We observe similar errors in cases where the person’s face and eyes are blurred (examples 4-5), which is more common in video datasets like VideoAttentionTarget and ChildPlay. Future work may explore using temporal information from surrounding frames to resolve ambiguities in these cases.



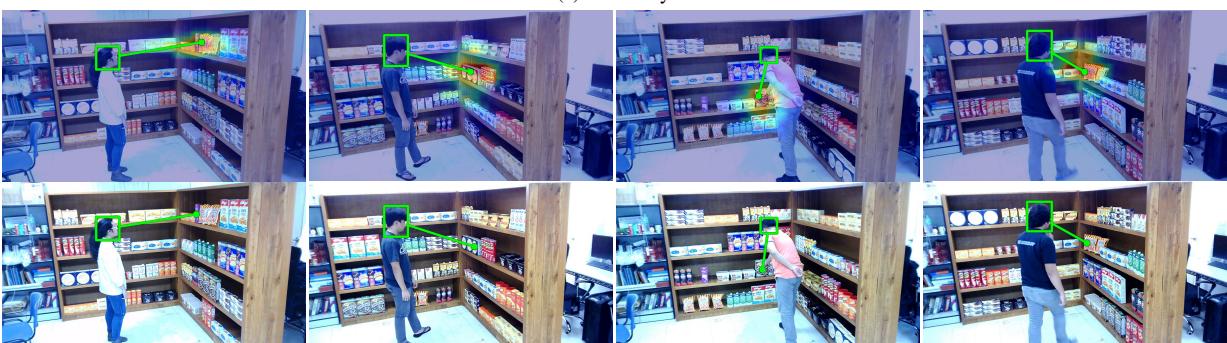
(a) GazeFollow



(b) VideoAttentionTarget



(c) ChildPlay



(d) GOO-Real

Figure 9. Additional qualitative results on the 4 evaluation datasets: For each example, we show our model's predicted heatmap with its maximum point on the top, and the ground truth gaze annotations on the bottom.

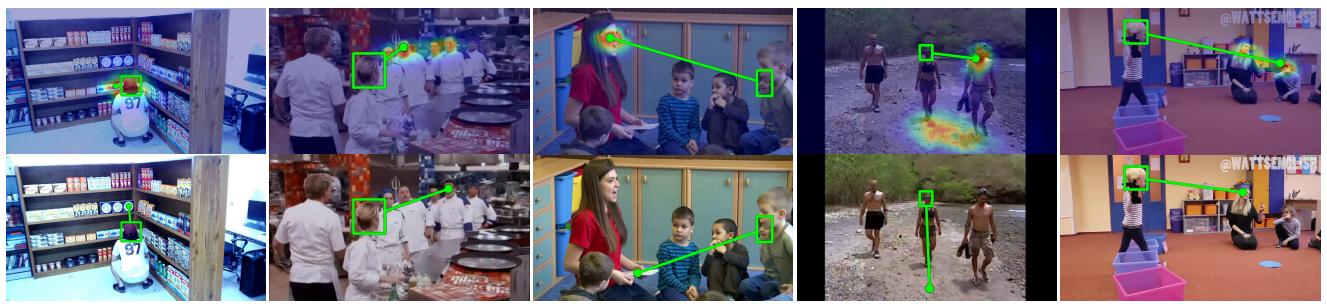


Figure 10. Lower performing cases: we observe errors in some cases where the head is facing away from the camera (examples 1-2), the head is occluded (examples 3), or the face is blurred (examples 4-5).