

Learning Ego-Exo Visual Representations for Conversational Gaze Estimation

Anshul Gupta^{1,2,3*}

Yijun Qian¹

Ruohan Gao⁴

Ishwarya Ananthabhotla¹

Jean-Marc Odobez^{2,3}

Vamsi Krishna Ithapu¹

Calvin Murdock¹

¹Meta Reality Labs Research, USA

²Idiap Research Institute, Martigny, Switzerland

³Ecole Polytechnique Fédérale de Lausanne, Switzerland

⁴University of Maryland, USA

Abstract

Egocentric gaze estimation is essential for understanding human attention in first-person scenarios, with applications in augmented reality, social interaction analysis, and assistive technology. While most existing methods rely on a sequence of video frames, real-world hardware constraints often necessitate single-frame inference. In this work, we introduce novel approaches that leverage exocentric gaze information of other individuals in the scene to improve single-frame egocentric gaze estimation. During training, our method leverages simultaneous views from a pair of people to jointly learn ego-exo gaze representations, with the exo representations learned via self-supervision. During inference, the model can leverage the learned exo representations to improve egocentric gaze estimation from a single view. Our experiments demonstrate that single frame models can achieve strong egocentric gaze performance, our approaches enable effective learning of exocentric gaze representations, and that learning these representations leads to improved egocentric gaze predictions.

1. Introduction

Augmented Reality and Virtual Reality wearables have seen rapid advancements and widespread adoption in recent years, driven by innovative products such as Ray-Ban Meta glasses [38], Snapchat Spectacles [22], and Apple Vision Pro [21]. These devices promise immersive, gaze-driven interactions for enabling intuitive experiences. In particular, detecting gaze towards faces is crucial for conversation-aware applications. For example, the looked at person is often the auditory focus of attention [1, 40]. In noisy environments, recognizing this can help steer microphones toward the intended speaker and suppress background noise.

However, these devices are often constrained by current hardware limitations. Specifically, integrated eye-tracking systems, though accurate, typically incur complex

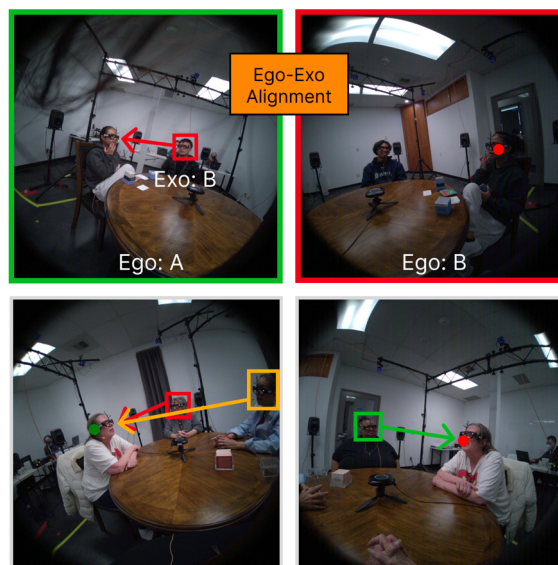


Figure 1. Estimating the egocentric gaze target of a person from a single frame is challenging, but can be improved using exocentric gaze cues of other individuals in the scene. During training, we employ a siamese-style architecture: one branch captures a person’s ego gaze features (top right), the other captures the *same person’s* exo features (top left), which are then aligned. Through symmetric ego-exo alignment and shared weights, the encoder learns to exploit exocentric gaze information from *other individuals* to improve egocentric gaze estimation from a single view (bottom).

calibration procedures, substantial power consumption and higher hardware costs, making them infeasible for many lightweight, cost-effective wearable platforms.

As a result, alternative approaches have emerged that estimate egocentric gaze using scene saliency and contextual cues from the wearer’s visual environment. Recent work in this domain [19, 28, 47] has demonstrated that such methods can achieve competitive performance compared to traditional eye-tracking systems. However, these approaches predominantly use a sequence of video frames, leveraging temporal cues to improve accuracy. Yet temporal models can take an order of magnitude more compute and memory

*Work completed during an internship at Meta.

than static models, which may limit their practicality (see discussion in supplementary).

Motivation. In this work, we focus on the single-frame egocentric gaze estimation task. In particular, we aim to leverage exocentric gaze (also known as gaze following) which involves predicting the gaze targets of other individuals in the scene from a third-person perspective. This additional context can help disambiguate between multiple potential gaze targets. For instance, in Figure 1, predicting egocentric gaze is ambiguous due to the presence of multiple individuals in the field of view. However, by incorporating exocentric gaze cues, we can better understand social interactions (e.g. eye-contact, shared attention¹ towards a person) allowing us to significantly reduce ambiguity.

Although existing models for egocentric gaze estimation may implicitly learn representations related to gaze following, the complexity of the task and the limitations of current datasets may restrict the extent to which these representations are effectively captured. Indeed, prior research has demonstrated the benefits of explicitly integrating additional contextual information such as hand actions [20] for egocentric gaze estimation. Guided by these insights, we explore whether explicitly learning exocentric gaze cues can enhance single-frame egocentric gaze estimation.

Contribution. Given these motivations, we propose novel approaches that leverage simultaneous views from a pair of individuals to jointly learn ego-exo gaze representations. Egocentric representations are learned via supervised training, while exocentric representations are learned through a self-supervised alignment task. Specifically, this task aims to match the ego representation of a person with the exo representation of the *same person* as captured from another person’s view (Figure 1, top).

We adopt a Siamese architecture where one branch captures egocentric gaze information and the other exocentric gaze information. Through symmetric ego-exo alignment and weight sharing, the same encoder learns to capture not only the egocentric gaze features of an individual but also the exocentric gaze features of *other individuals* in the scene. As a result, during inference, we can use a single branch to improve egocentric gaze estimation from a single view by leveraging the learned exocentric representations (Figure 1, bottom). This benefits practical use cases where simultaneous views of other people are often unavailable.

Our contributions can be summarized as follows:

- *Exploring single-frame egocentric gaze estimation:* We show that single-frame methods can achieve strong performance by leveraging modern CNN and transformer architectures.
- *Learning ego-exo gaze representations:* We propose three ego-exo alignment approaches: time synchronization, implicit matching and explicit matching.

¹We provide definitions of key terms in the supplementary

These approaches jointly learn ego and exo gaze representations, using self-supervision for the exo features. Our results show that these models improve egocentric gaze estimation, likely by utilizing the learned exocentric representations.

- *Probing for exocentric gaze:* We further probe the models to assess their ability to capture exocentric gaze representations, confirming that they indeed learn meaningful exocentric gaze features.
- *Additional metrics for egocentric gaze estimation:* We propose a suite of metrics inspired from gaze following literature, to enable more comprehensive analysis of model performance.

In addition, we perform an initial exploration of spatial audio for improving egocentric gaze performance, with promising preliminary results.

2. Related Work

Egocentric Gaze Estimation. There have been several works for egocentric gaze estimation using deep learning [19, 20, 28, 31, 47, 48]. In particular, Huang et al. [20] demonstrated the benefit of jointly modelling egocentric gaze with hand actions. Lai et al. [28] proposed the first transformer based model for this task, achieving state of the art results. Meanwhile, Thakur et al. [48] incorporated auxiliary information in the form of IMU measurements for improved performance. In the context of joint attention, Park et al. [44] leveraged data from multiple egocentric views to reconstruct the 3D scene and predict joint attention based on social formation. A related task is that of egocentric gaze anticipation first explored by Zhang et al. [55]. Lai et al. [29] performed autoregressive gaze anticipation, leveraging audio information for improved performance. While these works have pushed the state of the art, they can still fail on cases where gaze following information could have otherwise helped disambiguate the target.

Gaze Following. Recasens et al. [39] first introduced this task, proposing a two-branch architecture for processing the head crop and the scene. This design was then continued by several follow-up works [3, 8, 12, 25, 26, 33, 45], with some leveraging additional inferred modalities such as pose [12], audio [18] and depth [2, 8, 12, 26, 45] for improved performance. More recently, transformer based architectures [11, 41, 43, 46] have achieved state of the art results. In particular, [11, 46] perform multi-person gaze following, with [11] additionally leveraging temporal information and jointly modelling social gaze. On the other hand, [41, 43] leverage higher resolution scene images to allow discarding the head crop branch. A small number of works have also attempted to tackle data limitations by leveraging pseudo labels [35, 37] extracted using pre-trained models. We do not explicitly target the gaze following task, but in-

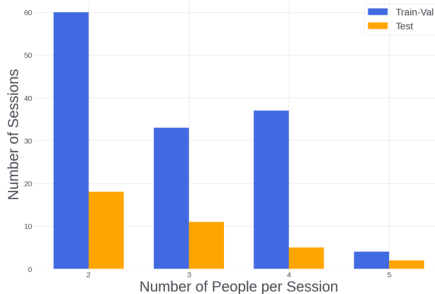


Figure 2. RLR-CHAT session distribution by number of people.

Split	Number of Frames	LAH Pairs	
		Positive	Negative
Train	1848555	273464	1107699
Val	448173	77914	309426
Test	385039	36361	351216

Table 1. RLR-CHAT number of frames and exo LAH statistics.

stead aim to learn gaze following representations via self-supervision for improving egocentric gaze estimation.

Ego-Exo Learning. This is a relatively nascent domain with works focusing on either person-matching between ego and exo views [6, 50, 51], learning of view-invariant features [32, 42, 52, 53], and conversational dynamics [23]. Our work shares similarities to [6, 51] as they also leverage contrastive losses to learn to associate the camera wearer in the egocentric view to the person in the exocentric view. However, our method learns to associate not just person identity, but also their gaze between ego and exo views. Another interesting work [53] also learns view invariant features, with the learned representations showing improvements for downstream applications including gaze angle prediction. However, they focus on scenes with a single person, and do not predict egocentric gaze or the gaze following target. Also, unlike all the above works, the exocentric view in our case does not come from a static camera but from another egocentric video which increases complexity.

3. The RLR-CHAT Dataset

The Reality Labs Research Conversations for Hearing Augmentation Technology (RLR-CHAT) dataset [17, 36, 54] is a large-scale collection of egocentric multisensory recordings captured from individuals engaging in natural conversations. Each conversation session is approximately one hour in duration and is recorded using Aria glasses [5], which capture RGB frames at 5Hz, 7-channel spatial audio at 48kHz, and eye-tracking data at 30Hz, among other modalities. Participants were recruited to ensure that no individual appeared in more than one session. Given the conversational setting, faces tend to be a common gaze target

(~40% cases). To maximize visual diversity, we subsample RGB frames by selecting every third frame and align them with the nearest eye-tracking annotations in time. The distribution of session sizes by participant count is illustrated in Figure 2. The dataset contains a total of 170 sessions, the majority of which involve two participants.

A distinctive feature of RLR-CHAT is the synchronized availability of modalities from all participants within the conversation (examples in Figure 1). This synchronization uniquely enables the exploration of ego-exo alignment techniques to learn richer gaze representations. To our knowledge, the only comparable accessible dataset is the Aria Everyday Activities dataset [34], which is significantly smaller and primarily focuses on single-person activities.

We augment RLR-CHAT by incorporating head box detections and automatically assigning identities to these boxes using spatial audio cues (details in supplementary). The test set includes manually corrected head bounding boxes and high-quality, OptiTrack-based head identity matching. By leveraging these identity-aligned head bounding boxes alongside eye-tracking annotations, we first determine if person A is looking at person B from A’s egocentric perspective. This information is then mapped to another person’s viewpoint (e.g., person C) to obtain exocentric annotations indicating whether person A is looking at the head of person B (denoted as $LAH_{A \rightarrow B}$). This annotation process is applied to all pairs of individuals present in the scene. The resulting LAH statistics are summarized in Table 1.

4. Method

Our training architecture, illustrated in Figure 3, follows a siamese design. The bottom branch predicts person A’s egocentric gaze, while the top branch predicts person B’s. Each branch processes an egocentric image frame, denoted as I^A and I^B , extracting features F^A and F^B using an encoder V . F^A captures A’s egocentric features, as well as others’ exocentric features (including B) and vice-versa for F^B . We then align the ego-exo gaze features of the same individual across views. We explore two approaches: (1) *Time Synchronization*, which encodes others’ exocentric features within a single global representation, and (2) *Head Matching*, which encodes others’ exocentric features in local representations extracted via head bounding boxes. Finally, the aligned features are passed through a decoder D_{ego} to generate each person’s egocentric gaze heatmap, H^A and H^B .

Due to weight sharing and symmetric ego-exo alignment, we can use a single branch of the network—termed EgoGazeViT—at inference.

4.1. Feature Extraction

The feature extraction module is responsible for obtaining gaze-relevant features from an input egocentric image frame, denoted as I . We employ a Vision Transformer (ViT)

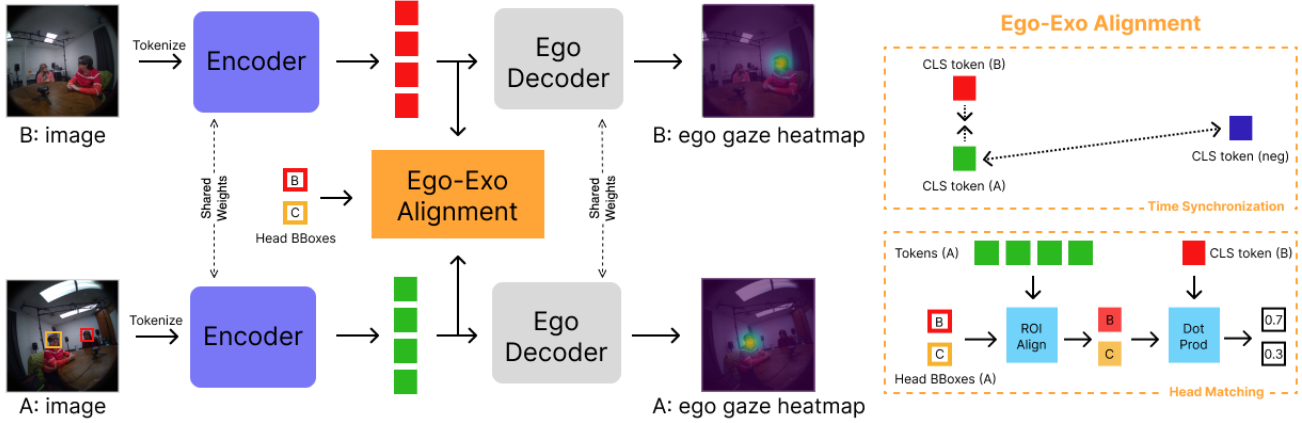


Figure 3. Our proposed architecture for ego-exo gaze representation learning. The Encoder first extracts features for each person’s view, which are then aligned using one of two Ego-Exo Alignment techniques: (1) time synchronization or (2) head matching. Finally, the Ego Decoder processes the aligned features to predict each person’s egocentric gaze heatmap. During inference, we leverage a single branch of the architecture, termed **EgoGazeViT**, for egocentric gaze estimation.

encoder, denoted as V , to extract these features. Specifically, we utilize the output from the last layer of the ViT.

$$\mathbf{F} = V(\mathbf{I}) \quad (1)$$

This module is applied independently to the egocentric images of both individuals, \mathbf{I}^A and \mathbf{I}^B , yielding corresponding feature representations, \mathbf{F}^A and \mathbf{F}^B .

4.2. Ego-Exo Alignment

The ego-exo alignment module enables self-supervised learning of exocentric gaze representations by aligning egocentric and exocentric features. The key idea is that a person’s egocentric gaze features already encode valuable information about where they are looking, which can be leveraged to supervise learning of their corresponding exocentric representations captured from another person’s viewpoint.

Through this alignment, we expect the ego and exo features corresponding to an individual to capture complementary information. For instance, in cases of eye contact, the exocentric features from A’s FoV can capture B’s head orientation and gaze direction, while B’s egocentric features can encode complementary information regarding B’s own head pose through cues like body orientation. Or in shared attention scenarios, both egocentric and exocentric features can capture similar visual cues about the attended item.

In practice, persons A and B are randomly sampled from a given timestamp. We explore two alignment approaches:

Time Synchronization. Inspired by prior work [53], we align egocentric features of two participants from the same timestamp in a session. Egocentric features for each person (e.g., person A) are obtained from the CLS token of the ViT output:

$$\mathbf{G}_{ego}^A = \text{CLS}(\mathbf{F}^A) \quad (2)$$

Here, the exocentric features for person A are directly captured in the egocentric features from person B :

$$\mathbf{G}_{exo}^A = \mathbf{G}_{ego}^B \quad (3)$$

We compute similarity between these ego-exo features using the L2 distance:

$$S = \|\mathbf{G}_{ego}^A - \mathbf{G}_{exo}^A\|_2 \quad (4)$$

Our triplet loss (Section 4.4) encourages high similarity between egocentric features from the same timestamp. Simultaneously, it minimizes similarity against negative samples drawn from the batch. These negatives may include features from different timestamps of the same session or features from entirely different sessions. Due to symmetric alignment, after training, we expect the same CLS token to capture *both* ego and others’ exo gaze information.

Head Matching. We explore a novel approach that aligns an individual’s egocentric features with their exocentric features encoded locally in the region corresponding to their head box in another person’s view. Specifically, given a participant B , we first extract exocentric features for all people visible in B ’s FoV (including A) using ROI-Align [15]:

$$\mathbf{G}_{exo}^B = \text{ROI Align}(\mathbf{F}^B, \mathbf{B}^B) \quad (5)$$

where \mathbf{B}^B represents the head bounding boxes in B ’s FoV. We again obtain egocentric features from the CLS token:

$$\mathbf{G}_{ego}^A = \text{CLS}(\mathbf{F}^A) \quad (6)$$

Similarity between the normalized egocentric and exocentric features is computed using a dot product:

$$\mathbf{S}^A = \mathbf{G}_{exo}^B \cdot \mathbf{G}_{ego}^A \quad (7)$$

Our loss function (Section 4.4) maximizes the similarity of matched ego-exo pairs ($\mathbf{S}^A(A)$) and minimizes it for unmatched pairs. This alignment is again performed symmetrically, so the same feature map learns to encode egocentric (in the CLS token) and others’ exocentric (in tokens corresponding to head box regions) gaze information.

4.3. Prediction

The Prediction Module processes the features from the Feature Extraction module, integrating both egocentric and others’ exocentric information to generate an egocentric gaze heatmap for each person. It consists of four transformer layers, followed by a linear projection layer that maps the processed token representations to the spatial dimensions of the gaze heatmap. Specifically, given the extracted feature representation \mathbf{F} from the Feature Extraction module, the egocentric gaze heatmap \mathbf{H} is predicted as follows:

$$\mathbf{H} = D_{\text{ego}}(\mathbf{F}) \quad (8)$$

where D_{ego} represents the transformer-based decoder.

This operation is applied independently to the feature representations of both person A and person B, producing their respective egocentric gaze heatmaps $\mathbf{H}^A, \mathbf{H}^B$.

4.4. Losses

Our training objective combines an egocentric gaze estimation loss ($\mathcal{L}_{\text{gaze}}$) and an ego-exo alignment loss ($\mathcal{L}_{\text{ego-exo}}$). $\mathcal{L}_{\text{gaze}}$ is a pixel-wise cross-entropy applied independently to the predicted heatmaps \mathbf{H}^A and \mathbf{H}^B , comparing them to the corresponding ground truth heatmaps.

The total loss is given by:

$$\mathcal{L} = \mathcal{L}_{\text{gaze}}^A + \mathcal{L}_{\text{gaze}}^B + \mathcal{L}_{\text{ego-exo}} \quad (9)$$

We explore two formulations for the ego-exo alignment loss ($\mathcal{L}_{\text{ego-exo}}$):

Time Synchronization Loss. This loss uses a triplet formulation based on the similarity between egocentric features at the same timestamp:

$$\mathcal{L}_{\text{ego-exo}} = \frac{e^{S^+}}{e^{S^+} + e^{S^-}} \quad (10)$$

where S^+ is the distance between matched ego features from simultaneous views, and S^- is the distance from $\mathbf{G}_{\text{ego}}^A$ to a randomly sampled negative ego feature from the batch.

Head Matching Loss. We explore two variants for the head matching loss: an *implicit* approach that automatically learns the alignment between egocentric and exocentric features, and an *explicit* approach leveraging ground truth head-box identities when available.

- **Explicit matching:** This is also a cross-entropy loss, applied to the similarity scores \mathbf{S} . The correct "class" corresponds to the similarity score of the same person

viewed from the other perspective. Specifically, for \mathbf{S}^A , the correct class is $\mathbf{S}^A(A)$, and vice versa for \mathbf{S}^B .

- **Implicit matching:** We apply an entropy loss on the similarity scores \mathbf{S} . This encourages the model to select exactly one exocentric feature with maximum similarity. Additionally, for two-person sessions, where head-box identities are trivially identifiable, we apply the explicit cross-entropy loss described above.

In both variants, the total ego-exo loss $\mathcal{L}_{\text{ego-exo}}$ is the sum of the losses computed independently for persons A and B .

5. Experiments

5.1. Datasets

We perform experiments on two datasets:

RLR-CHAT. The RLR-CHAT dataset, as described in Section 3, is divided into training, validation, and test splits. In particular, we refer to the test split as the "golden subset" as it has higher quality annotations. Initially, we train and evaluate various baselines on this golden subset, as training on the entire dataset is computationally expensive. Subsequently, we leverage the best performing approach for training on the full dataset.

Ego4D [10]. It is a large-scale, publicly available egocentric dataset that captures individuals performing daily life activities. We use the subset with gaze annotations introduced by Lai et al. [28], which consists of 27 approximately hour-long videos featuring 80 participants engaged in social interactions such as playing board games.

We selected Ego4D over alternative datasets such as EGTEA Gaze [30] and Aria [34], as those primarily focus on single-person activities. Since Ego4D emphasizes social settings, it is more suitable for evaluating improvements derived from learning exocentric gaze representations.

However, it is important to note that there is still a significant domain gap between Ego4D and RLR-CHAT. The images have a smaller FoV and include much more diverse settings. Further, as people are playing games instead of mainly conversing, the gaze points tend to fall less on faces.

5.2. Trained Models

We train several egocentric gaze estimation baselines along with our proposed self-supervised methods using the RLR-CHAT dataset.

Egocentric baselines. We compare heuristic baselines, as well as CNN and transformer-based approaches:

- **Heuristic Baselines:** Predicting the image center, using the average gaze point from training, and selecting the head closest to the image center as the gaze target.
- **U-Net:** A CNN-based model with a ResNet-18 encoder and an FPN-style decoder. It corresponds to the image branch of the MAV-Gaze baseline described below and operates on single-frame inputs.

- **MAV-Gaze:** An adaptation of MAV-ASL [24], originally designed for active speaker localization. It processes both visual and auditory cues, taking in a single image frame and a 7-channel, 200ms audio clip.
- **EgoGazeViT:** A transformer-based model comprising a ViT encoder and a transformer decoder. This architecture corresponds to one of the branches in our proposed SSL methods and operates on single-frames.

Self-supervised Approaches. We initialize EgoGazeViT with one of our three self-supervised alignment methods described in Section 4, or with standard training for egocentric gaze prediction (**Standard Training**):

- **Synchronization:** Aligns egocentric features across simultaneous views using temporal correspondence.
- **Implicit Matching:** Aligns ego-exo features without explicit identity annotations.
- **Explicit Matching:** Explicitly aligns ego-exo features using ground truth head-box identities.

5.3. Metrics

Previous works on egocentric gaze estimation evaluate the quality of the predicted gaze heatmap by performing a pixel-level comparison against the generated ground truth heatmap after binarizing both of them [28] (referred to as the **Heatmap** metric). However, this approach is highly sensitive to the choice of the heatmap’s standard deviation and the threshold used for binarization. Similar concerns have led to the avoidance of such metrics in recent gaze following research [11, 46]. Moreover, these metrics assess only localization accuracy, and may not always reflect semantic performance, which is often more valuable for downstream applications.

To address these limitations, we propose a new set of evaluation metrics inspired from gaze following [39, 45]:

- **Distance:** The predicted gaze point is obtained by taking the argmax of the predicted gaze heatmap. We then compute the L2 distance between the predicted and ground truth gaze points, normalized to a unit square (1×1). Both mean and median distances are reported.
- **Looking at Heads (LAH):** This semantic metric evaluates how well the model detects gaze directed at other people’s heads. A prediction is classified as follows:
 - **True Positive:** Both the predicted and ground truth gaze points fall in the same head box.
 - **False Positive:** The predicted gaze point falls in a head box, but the ground truth does not.
 - **False Negative:** The ground truth gaze point falls in a head box, but the predicted gaze point falls on a different head or object.
 - **True Negative:** Neither the predicted nor the ground truth gaze points fall in a head box.

We compute precision, recall, and F1-scores.

Model	Distance↓			LAH↑	
	Mean	Median	Prec	Recall	F1
<i>Heuristic Baselines</i>					
Predict center	0.107	0.093	0.633	0.146	0.237
Predict avg of train data	0.105	0.092	0.638	0.130	0.216
Predict closest head to center	0.131	0.073	0.396	0.863	0.543
<i>CNN Baselines</i>					
U-Net	0.105	0.072	0.520	0.610	0.561
MAV-Gaze	0.098	0.065	0.617	0.724	0.667
<i>Transformer Baseline</i>					
EgoGazeViT (Standard Training)	0.096	0.057	0.507	0.798	0.620

Table 2. Comparison of egocentric gaze estimation baselines on the RLR-CHAT golden subset test split. Best results are in bold.

Initialization	Distance↓			LAH↑	
	Mean	Median	Prec	Recall	F1
Standard training	0.102	0.057	0.538	0.819	0.650
<i>SSL Approaches</i>					
Synchronization	0.100	0.055	0.536	0.843	0.656
Implicit matching	0.101	0.056	0.533	0.833	0.650
Explicit matching	0.101	0.055	0.545	0.836	0.660

Table 3. Results for egocentric gaze estimation on the full RLR-CHAT golden subset. We leverage EgoGazeViT with different initializations. Best results are in bold.

6. Results

6.1. Egocentric Gaze Estimation Baselines

We trained and evaluated the baselines listed in Section 5.2 on the RLR-CHAT golden subset. Specifically, each model is trained on the train split of the golden subset for egocentric gaze estimation and evaluated then evaluated on its test set. The results are presented in Table 2.

Despite operating on a single image frame, both transformer and CNN-based models outperform the heuristic baselines by a significant margin, highlighting their ability to incorporate human priors and scene saliency for accurate egocentric gaze estimation. EgoGazeViT achieves the highest performance among the image-only models, and the best overall distance score. Therefore, we select this model for all subsequent experiments.

Interestingly, MAV-Gaze achieves the highest LAH F1-score. The incorporation of spatial audio helps the model identify the speaking person, which can serve as a strong cue for gaze target prediction. Exploring the role of spatial audio in egocentric gaze estimation remains an exciting avenue for future work.

6.2. Learning Exo Gaze Representations

We leverage the entire RLR-CHAT dataset by training our models on the designated train split and evaluating them on

Initialization	Number of People		
	Full	≥ 3	≥ 4
Standard Training	0.650	0.630	0.576
<i>SSL Approaches</i>			
Synchronization	0.656	0.634	0.578
Implicit matching	0.650	0.626	0.553
Explicit matching	0.660	0.640	0.587

Table 4. LAH F1-scores for different splits of the RLR-CHAT golden subset based on the number of participants in the sessions. Best results are in bold.

the full test split (the golden subset). Specifically, we compare the performance of EgoGazeViT when initialized with weights from standard egocentric gaze estimation training versus weights obtained via our proposed ego-exo alignment methods. Results (Table 3) indicate that performance on the Distance metric for all methods is comparable. However, Synchronization and Explicit Matching yield some improvements over Standard Training for the LAH metric, with Explicit Matching having the best performance.

As seen in Table 4, Explicit Matching consistently improves over Standard Training across evaluations on different splits of the RLR-CHAT golden subset based on the number of participants in the sessions (full results in supplementary). This suggests that our method captures exocentric gaze behaviors beyond shared attention. Shared attention in RLR-CHAT is predominantly observed when gaze is directed towards other people—a scenario naturally limited in sessions with only two participants, which constitute the majority of the dataset. However, the performance gap implies that the model also learns other behaviors, such as eye contact, that help disambiguate the egocentric gaze target.

Probing for Exocentric Gaze. To assess whether training with our proposed method effectively enables the learning of exocentric gaze representations, we probe the trained encoder by evaluating its performance on exocentric gaze prediction. The probing architecture is illustrated in Figure 4. Specifically, we freeze the trained encoder and train a new exocentric decoder D_{exo} to predict LAH labels for RLR-CHAT. D_{exo} is a 2 layer MLP that operates on ROI-aligned features corresponding to individuals visible within the egocentric field of view (\mathbf{G}_{exo}). Specifically, it processes their concatenated features and predicts pairwise LAH following the formulation of [11, 13]. For instance, to predict whether person B is looking at person C within person A’s FoV, the model proceeds as follows:

$$\text{LAH}_{B \rightarrow C} = D_{\text{exo}}(\mathbf{G}_{\text{exo}}^B, \mathbf{G}_{\text{exo}}^C) \quad (11)$$

Note that the order of individuals supplied to the decoder is crucial, as the LAH prediction is directional.

We present the results in Table 5. Unlike egocentric gaze prediction, where discrete LAH labels enable direct preci-

Initialization/Model	Distance \downarrow		Heatmap \uparrow		
	Mean	Median	Prec	Recall	F1
<i>Cross-Dataset Evaluation</i>					
Standard training	0.174	0.151	0.260	0.520	0.347
Synchronization	0.169	0.144	0.286	0.506	0.365
Implicit Matching	0.183	0.154	0.259	0.537	0.349
Explicit Matching	0.170	0.147	0.274	0.493	0.352
<i>Within-Dataset Evaluation</i>					
GBVS [14]	-	-	0.111	0.472	0.180
Attention Transition [19]	-	-	0.295	0.476	0.364
I3D-R50 [9]	-	-	0.292	0.525	0.375
MViT [28]	-	-	0.317	0.574	0.409
GLC [28]	0.156	0.123	0.347	0.570	0.431
EgoGazeViT (Explicit Matching init)	0.163	0.131	0.315	0.562	0.404

Table 6. Results for egocentric gaze estimation on the Ego4D dataset. Best results are in bold.

sion and recall calculations, the predicted LAH values in this setting are continuous. While applying a threshold can yield discrete values, the precision and recall scores can vary significantly depending on that choice. Therefore, we report the average precision (AP) score, which provides a threshold-independent evaluation.

We find that all self-supervised approaches outperform the baseline, indicating that they successfully capture exocentric gaze information. Notably, the Synchronization approach significantly improves over the other self-supervised approaches. This may be related to the more general nature of the alignment task, which allows the exocentric features to encode global social gaze cues because it does not rely on head crops. The Implicit Matching approach also surpasses Explicit Matching, however, this may be a result of overfitting to scene geometry when learning head-identity correspondences. This interpretation is supported by its lower cross-dataset performance as discussed in the next section.

6.3. Evaluation on Ego4D

We provide cross-dataset evaluation results for our RLR-CHAT trained models on Ego4D in Table 6. Overall all methods have a marked drop in distance score, highlighting the domain gap between the two datasets.

Despite this gap, all self-supervised approaches except Implicit Matching improve over Standard Training in cross-dataset generalization, illustrating another benefit of ego-exo alignment. The Synchronization approach has the best overall performance, following results from exocentric probing. Interestingly, this trend is not followed for Implicit Matching, which suggests that it may be overfitting to scene geometry in order to learn head-identity correspondences.

For comparison with state-of-the-art methods, we additionally train one of our models—EgoGazeViT initialized with Explicit Matching—on Ego4D. Despite relying solely on single frames, it attains strong performance and even surpasses some temporal models, highlighting the potential of

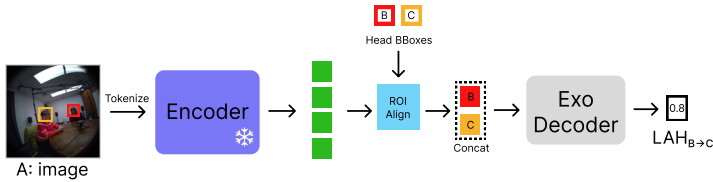


Figure 4. Architecture for probing learned exocentric gaze representations. We freeze the encoder, which was initially trained for egocentric gaze estimation, and train a 2-layer MLP probe to predict Looking at Heads (LAH).

Initialization	LAH AP \uparrow
Random init	0.178
Standard training	0.262
<i>SSL Approaches</i>	
Synchronization	0.498
Implicit Matching	0.371
Explicit Matching	0.304

Table 5. Results for exocentric gaze probing on the full RLR-CHAT golden subset. Best results are in bold.



Figure 5. Qualitative results on RLR-CHAT for egocentric (top) and exocentric (bottom) gaze prediction using the encoder initialized with our Explicit Matching based self-supervised approach. The predicted egocentric gaze heatmap is overlaid on the image, with the ground truth target marked by a green dot. Predicted exocentric gaze targets are indicated by the person ID following the 'LAH' prefix.

single-frame approaches in this domain. We observe that different initializations of EgoGazeViT do not yield significant performance variations on Ego4D, likely due to the pronounced domain shift between datasets.

6.4. Qualitative Results

We present qualitative results for egocentric and exocentric gaze prediction in Figure 5, using the encoder initialized with Explicit Matching. For egocentric predictions, we directly overlay the predicted heatmap onto the image. For exocentric predictions of a given person B , the target is determined by the argmax over LAH pairs $(B, *)$, and is visualized if the corresponding value exceeds 0.1.

We observe that the model accurately identifies the egocentric gaze target (columns 2-4). Generally, it tends to focus on salient items such as faces, whereas human gaze can sometimes be directed toward background individuals (column 1) or be in transition during a gaze shift (column 5), which is challenging for a static model to capture. Additionally, we observe that the model effectively leverages exocentric cues to resolve ambiguities (columns 1-4). How-

ever, in scenarios where exocentric gaze information is less informative (column 5), the model exhibits greater uncertainty, as reflected in the multimodal heatmap. Additional qualitative results are in the supplementary material.

7. Conclusion

In this work, we introduced novel self-supervised learning approaches for egocentric gaze estimation that leverage ego-exo alignment to learn exocentric gaze representations. Our methods improve egocentric gaze prediction in challenging single-frame setting across RLR-CHAT and Ego4D by leveraging the learned exocentric gaze representations. Furthermore, our probing analysis confirms that training with our method enhances the encoder's ability to learn these exocentric gaze representations.

Future research could explore integrating spatial audio cues to further refine gaze estimation, particularly in social settings where auditory information plays a key role in attention and interaction. Additionally, investigating the generalizability of these self-supervised techniques in temporal settings could be another interesting direction of research.

References

- [1] S.O. Ba and J.-M. Odobez. Recognizing human visual focus of attention from head pose in natural meetings. *IEEE Trans. on System, Man and Cybernetics: part B, Cybernetics*, 39(1): 16–34, 2009. 1
- [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. 2
- [3] 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. 2
- [4] 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. In *Proceedings of the International Conference on Learning Representations*, 2021. 14
- [5] Jakob Engel, Kiran Somasundaram, Michael Goesele, Albert Sun, Alexander Gamino, Andrew Turner, Arjang Talattof, Arnie Yuan, Bilal Souti, Brighid Meredith, et al. Project aria: A new tool for egocentric multi-modal ai research. *arXiv preprint arXiv:2308.13561*, 2023. 3
- [6] Chenyou Fan, Jangwon Lee, Mingze Xu, Krishna Kumar Singh, Yong Jae Lee, David J Crandall, and Michael S Ryoo. Identifying first-person camera wearers in third-person videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5125–5133, 2017. 3
- [7] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6824–6835, 2021. 13
- [8] 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. 2
- [9] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6202–6211, 2019. 7
- [10] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012, 2022. 5
- [11] Anshul Gupta, Samy Tafasca, Arya Farkhondeh, Pierre Vuillecard, and Jean-marc Odobez. Mtgs: A novel framework for multi-person temporal gaze following and social gaze prediction. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*. 2, 6, 7
- [12] Anshul Gupta, Samy Tafasca, and Jean-Marc Odobez. 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*, pages 5041–5050, 2022. 2
- [13] Anshul Gupta, Samy Tafasca, Naravich Chutisilp, and Jean-Marc Odobez. A unified model for gaze following and social gaze prediction. In *2024 IEEE 18th International Conference on Automatic Face and Gesture Recognition (FG)*, pages 1–9. IEEE, 2024. 7
- [14] Jonathan Harel, Christof Koch, and Pietro Perona. Graph-based visual saliency. *Advances in neural information processing systems*, 19, 2006. 7
- [15] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 4
- [16] 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. 14
- [17] Seongsil Heo, Calvin Murdock, Michael Proulx, and Christi Miller. Gaze-enhanced multimodal turn-taking prediction in triadic conversations. In *Proceedings of Interspeech*, 2025. 3
- [18] Yuqi Hou, Zhongqun Zhang, Nora Horanyi, Jaewon Moon, Yihua Cheng, and Hyung Jin Chang. Multi-modal gaze following in conversational scenarios. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1186–1195, 2024. 2
- [19] Yifei Huang, Minjie Cai, Zhenqiang Li, and Yoichi Sato. Predicting gaze in egocentric video by learning task-dependent attention transition. In *Proceedings of the European conference on computer vision (ECCV)*, pages 754–769, 2018. 1, 2, 7
- [20] Yifei Huang, Minjie Cai, Zhenqiang Li, Feng Lu, and Yoichi Sato. Mutual context network for jointly estimating egocentric gaze and action. *IEEE Transactions on Image Processing*, 29:7795–7806, 2020. 2
- [21] Apple Inc. Apple vision pro. <https://www.apple.com/apple-vision-pro/>, . Accessed: 2025-01-13. 1
- [22] Snap Inc. Spectacles '24. <https://www.spectacles.com/spectacles-24>, . Accessed: 2025-01-13. 1
- [23] Wenqi Jia, Miao Liu, Hao Jiang, Ishwarya Ananthabhotla, James M Rehg, Vamsi Krishna Ithapu, and Ruohan Gao. The audio-visual conversational graph: From an egocentric-exocentric perspective. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 3
- [24] Hao Jiang, Calvin Murdock, and Vamsi Krishna Ithapu. Egocentric deep multi-channel audio-visual active speaker localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10544–10552, 2022. 6, 12
- [25] Tianlei Jin, Zheyuan Lin, Shiqiang Zhu, Wen Wang, and Shunda Hu. Multi-person gaze-following with numerical co-

- ordinate regression. In *2021 16th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2021)*, pages 01–08. IEEE, 2021. 2
- [26] 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. 2
- [27] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 13
- [28] Bolin Lai, Miao Liu, Fiona Ryan, and James M Rehg. In the eye of transformer: Global–local correlation for egocentric gaze estimation and beyond. *International Journal of Computer Vision*, pages 1–18, 2023. 1, 2, 5, 6, 7, 13
- [29] Bolin Lai, Fiona Ryan, Wenqi Jia, Miao Liu, and James M Rehg. Listen to look into the future: Audio-visual egocentric gaze anticipation. In *European Conference on Computer Vision*, pages 192–210. Springer, 2025. 2
- [30] Yin Li, Miao Liu, and James M Rehg. In the eye of beholder: Joint learning of gaze and actions in first person video. In *Proceedings of the European conference on computer vision (ECCV)*, pages 619–635, 2018. 5
- [31] Yin Li, Miao Liu, and James M Rehg. In the eye of the beholder: Gaze and actions in first person video. *IEEE transactions on pattern analysis and machine intelligence*, 45(6): 6731–6747, 2021. 2
- [32] Yanghao Li, Tushar Nagarajan, Bo Xiong, and Kristen Grauman. Ego-exo: Transferring visual representations from third-person to first-person videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6943–6953, 2021. 3
- [33] 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. 2
- [34] Zhaoyang Lv, Nicholas Charron, Pierre Moulon, Alexander Gamino, Cheng Peng, Chris Sweeney, Edward Miller, Huixuan Tang, Jeff Meissner, Jing Dong, Kiran Somasundaram, Luis Pesqueira, Mark Schwesinger, Omkar Parkhi, Qiao Gu, Renzo De Nardi, Shangyi Cheng, Steve Saarinen, Vijay Baiyya, Yuyang Zou, Richard Newcombe, Jakob Julian Engel, Xiaqing Pan, and Carl Ren. Aria everyday activities dataset, 2024. 3, 5
- [35] Qiaomu Miao, Alexandros Graikos, Jingwei Zhang, Sounak Mondal, Minh Hoai, and Dimitris Samaras. Diffusion-refined vqa annotations for semi-supervised gaze following. In *European Conference on Computer Vision*, pages 439–457. Springer, 2025. 2
- [36] Calvin Murdock, Ishwarya Ananthabhotla, Hao Lu, and Vamsi Krishna Ithapu. Self-motion as supervision for egocentric audiovisual localization. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7835–7839. IEEE, 2024. 3
- [37] Cheng Peng and Oya Celiktutan. Visual saliency guided gaze target estimation with limited labels. In *2024 IEEE 18th International Conference on Automatic Face and Gesture Recognition (FG)*, pages 1–9. IEEE, 2024. 2
- [38] Ray-Ban. Ray-ban meta smart glasses. <https://www.ray-ban.com/usa/ray-ban-meta-smart-glasses>. Accessed: 2025-01-13. 1
- [39] Adria Recasens, Aditya Khosla, Carl Vondrick, and Antonio Torralba. Where are they looking? *Advances in neural information processing systems*, 28, 2015. 2, 6
- [40] Federico Rossano. Gaze in conversation. *The handbook of conversation analysis*, pages 308–329, 2012. 1
- [41] Fiona Ryan, Ajay Bati, Sangmin Lee, Daniel Bolya, Judy Hoffman, and James M Rehg. Gaze-llc: Gaze target estimation via large-scale learned encoders. *arXiv preprint arXiv:2412.09586*, 2024. 2
- [42] Gunnar A Sigurdsson, Abhinav Gupta, Cordelia Schmid, Ali Farhadi, and Karteek Alahari. Actor and observer: Joint modeling of first and third-person videos. In *proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7396–7404, 2018. 3
- [43] Yuehao Song, Xinggang Wang, Jingfeng Yao, Wenyu Liu, Jinglin Zhang, and Xiangmin Xu. Vitgaze: Gaze following with interaction features in vision transformers. *arXiv preprint arXiv:2403.12778*, 2024. 2
- [44] 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. 2
- [45] Samy Tafasca, Anshul Gupta, and Jean-Marc Odobez. Child-play: A new benchmark for understanding children’s gaze behaviour. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20935–20946, 2023. 2, 6
- [46] 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. 2, 6
- [47] Hamed Rezazadegan Tavakoli, Esa Rahtu, Juho Kannala, and Ali Borji. Digging deeper into egocentric gaze prediction. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 273–282. IEEE, 2019. 1, 2
- [48] Sanket Kumar Thakur, Cigdem Beyan, Pietro Morerio, and Alessio Del Bue. Predicting gaze from egocentric social interaction videos and imu data. In *Proceedings of the 2021 International Conference on Multimodal Interaction*, pages 717–722, 2021. 2
- [49] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. *Advances in neural information processing systems*, 35:10078–10093, 2022. 13
- [50] Yangming Wen, Krishna Kumar Singh, Markham Anderson, Wei-Pang Jan, and Yong Jae Lee. Seeing the unseen: Predicting the first-person camera wearer’s location and pose in third-person scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3446–3455, 2021. 3
- [51] Mingze Xu, Chenyou Fan, Yuchen Wang, Michael S Ryoo, and David J Crandall. Joint person segmentation and identification in synchronized first-and third-person videos. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 637–652, 2018. 3

- [52] Zihui Sherry Xue and Kristen Grauman. Learning fine-grained view-invariant representations from unpaired ego-exo videos via temporal alignment. *Advances in Neural Information Processing Systems*, 36:53688–53710, 2023. [3](#)
- [53] Huangyue Yu, Minjie Cai, Yunfei Liu, and Feng Lu. First-and third-person video co-analysis by learning spatial-temporal joint attention. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(6):6631–6646, 2020. [3](#), [4](#)
- [54] Heeseung Yun, Ruohan Gao, Ishwarya Ananthabhotla, Anurag Kumar, Jacob Donley, Chao Li, Gunhee Kim, Vamsi Krishna Ithapu, and Calvin Murdock. Spherical world-locking for audio-visual localization in egocentric videos. In *European Conference on Computer Vision*, pages 256–274. Springer, 2024. [3](#)
- [55] Mengmi Zhang, Keng Teck Ma, Joo Hwee Lim, Qi Zhao, and Jiashi Feng. Deep future gaze: Gaze anticipation on egocentric videos using adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4372–4381, 2017. [2](#)

Supplementary Material

A. Head Bounding Box Identification

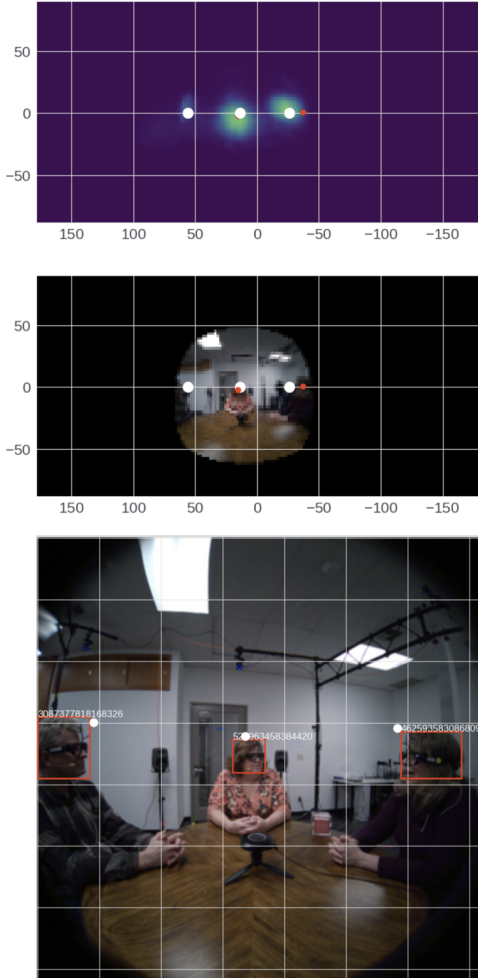


Figure S1. Overall pipeline for identifying head bounding boxes. We first obtain the average locations of other participants (white dots) in the conversation in 360 degrees world-locked coordinates. These are then mapped to the head-locked FoV coordinates and matched to the nearest head bounding box within a threshold.

To reliably identify the head bounding boxes of individuals visible in a participant’s field of view (FoV), we leverage a pre-trained MAV-ASL model [24] to obtain active speaker heatmaps for each egocentric image frame. The MAV-ASL model produces two types of heatmaps: one that indicates the direction of the active speaker over a full 360-degree span, and another that provides the 2D location of the active speaker when present in the FoV. Both heatmaps are initially computed in a head-locked coordinate system.

We begin by utilizing the directional heatmap and converting it into world-locked coordinates with the help of SLAM data. By processing 5000 frames per participant, we compute the average world-locked location of the other participants in the conversation. For each frame, these average locations are then transformed back into the head-locked FoV coordinate system.

Subsequently, we match the detected head bounding boxes to these averaged locations by selecting the nearest match within a threshold of 200 pixels. Although this thresholding process means that not all head bounding boxes are assigned an identity, the matches that are made have been verified to be of high quality. The overall pipeline is illustrated in Figure S1.

B. Discussion

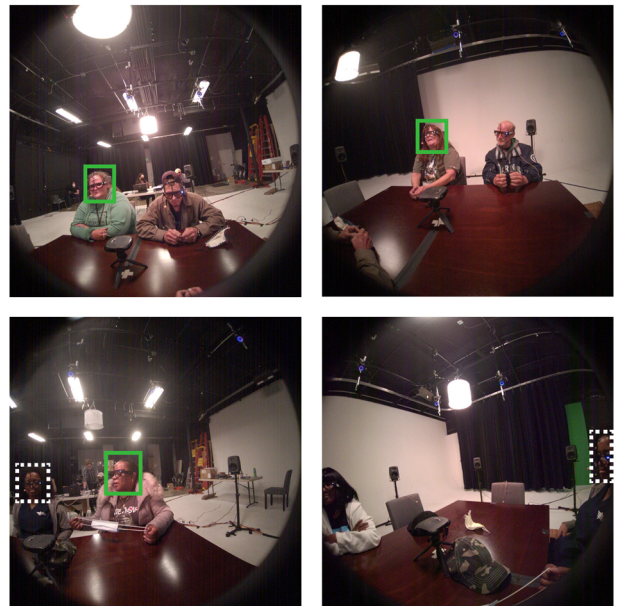


Figure S2. Egocentric and exocentric feature alignment results for the Implicit Matching approach. The correct exocentric person is highlighted with a green box. In the top row, these are correctly selected by the model. Incorrect selections made by the model are indicated with dotted white boxes in the bottom row.

Implicit Matching vs. Explicit Matching. In Figure S2, we examine how well the Implicit Matching method learns to align egocentric and exocentric gaze features. The model successfully matches features in several cases (top row), demonstrating its ability to capture meaningful ego-exo correspondences. However, it also exhibits failure cases

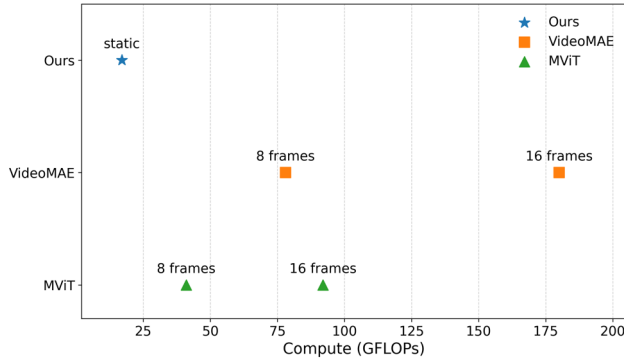


Figure S3. FLOPs comparison of our model against video-based models. Our model requires significantly fewer FLOPs while still achieving competitive performance.

(bottom row), where mismatches occur. In some cases, such as the bottom right example, failure is expected because the corresponding exocentric person is outside the field of view. However, the model also fails in other scenarios (e.g., bottom left), indicating that implicit matching alone may not always be sufficient for robust alignment.

Impact of number of people. Table S1 provides a detailed breakdown of egocentric gaze estimation performance on different splits of the RLR-CHAT Golden Subset, based on the number of participants in the included sessions. As expected, performance generally declines in sessions with a higher number of people due to the increased number of potential gaze targets and the resulting complexity of the task. Notably, there is an apparent spike in performance for sessions with 5 participants; however, since this split comprises only 2 sessions, the result is likely subject to high variance and may not be representative.

Impact of using exocentric views at inference. Since there is no feature sharing between the two branches of our siamese architecture, using exocentric views does not improve performance during inference. During training, the architecture enables self-supervised learning of exocentric representations through ego-exo alignment, but at inference each branch operates independently. The architecture was deliberately designed this way to avoid the need for additional views during inference.

Memory and Compute of Static vs. Temporal Models. As motivated in the paper introduction, static models provide benefits for practical applications as they require lower compute and memory compared to temporal models. First, temporal models require maintaining a frame buffer. For a model processing even just 8 frames, this multiplies storage requirements by $8\times$. They also incur substantially higher computational costs. In Figure S3, we compare the FLOPs consumed by our model, VideoMAE-based models [49], and MViT-based models [7].

We observe that all video-based models consume significantly more compute than our static model, which is reasonable as they process a larger number of tokens. MViT, used in [28], employs a more efficient attention mechanism and thus requires less compute than VideoMAE. Increasing the number of frames from 8 to 16 further raises compute requirements. Actual FPS would then depend on the hardware and platform on which these models are deployed.

C. Qualitative Comparisons

We provide qualitative comparisons of our models for egocentric gaze prediction in Figures S4 and S5. In Figure S4, all models perform similarly, as these cases involve either a single salient target or a target positioned near the image center, reducing task complexity.

In contrast, Figure S5 highlights scenarios where our ego-exo alignment approaches improve egocentric gaze prediction. These cases are more ambiguous, featuring multiple salient targets near the center, making gaze estimation more challenging.

- *Row 1:* The Standard Training and Implicit Matching approaches incorrectly identify the target person.
- *Row 2:* The Synchronization and Implicit Matching approaches struggle to differentiate between two people. The Standard Training approach confidently selects the wrong target, whereas the Explicit Matching approach correctly identifies the target with high confidence.
- *Row 3:* The Standard Training and Synchronization approaches misidentify the target, while the Implicit Matching and Explicit Matching approaches show uncertainty between two possible targets.
- *Row 4:* The Standard Training approach produces a diffused heatmap due to confusion, whereas the ego-exo alignment approaches correctly select the target. The Explicit Matching approach still exhibits some uncertainty.

These results suggest that ego-exo alignment helps disambiguate complex scenarios by leveraging exocentric gaze cues, leading to more precise egocentric gaze predictions.

D. Training and Evaluation

During training, we randomly sample two people A and B for each timestamp. Models are trained for 20 epochs using the Adam optimizer [27] with a learning rate of 2×10^{-5} , and with a batch size of 512. We employ standard augmentations, namely center cropping, flipping, and color jittering. The method was trained on a distributed system with two nodes, each equipped with eight H100 GPUs.

Subset	Initialization	Distance		LAH		
		Mean	Median	Precision	Recall	F1
Full	Standard Training	0.102	0.057	0.538	0.819	0.650
	Synchronization	0.100	0.055	0.536	0.843	0.656
	Implicit Matching	0.101	0.056	0.533	0.833	0.650
	Explicit Matching	0.101	0.055	0.545	0.836	0.660
≥ 3 people	Standard Training	0.111	0.067	0.524	0.790	0.630
	Synchronization	0.110	0.064	0.519	0.815	0.634
	Implicit Matching	0.111	0.065	0.512	0.805	0.626
	Explicit Matching	0.110	0.065	0.532	0.803	0.640
≥ 4 people	Standard Training	0.110	0.074	0.466	0.754	0.576
	Synchronization	0.107	0.069	0.461	0.771	0.578
	Implicit Matching	0.111	0.074	0.438	0.750	0.553
	Explicit Matching	0.106	0.069	0.473	0.773	0.587
≥ 5 people	Standard Training	0.096	0.054	0.542	0.820	0.653
	Synchronization	0.090	0.049	0.546	0.849	0.664
	Implicit Matching	0.101	0.058	0.524	0.796	0.632
	Explicit Matching	0.092	0.052	0.554	0.827	0.664

Table S1. Evaluation results on different splits of the RLR-CHAT Golden Subset based on the number of people in the session. Best results for each split are given in bold.

During evaluation, we leverage only one of the branches (since both share weights), referred to as EgoGazeViT, to assess egocentric gaze estimation performance. Specifically, EgoGazeViT can be initialized with weights from either the self-supervised training or standard egocentric gaze estimation training.

E. Implementation Details

The model uses a ViT-B encoder initialized with masked autoencoder (MAE) pretraining [16]. The Prediction Module consists of four transformer layers, each with a token dimension of 384—half the dimension of the ViT tokens. Input images are processed at a resolution of 224×224 , and the predicted gaze heatmap is generated at the same resolution, following the MAE architecture. The ground truth gaze heatmap \mathbf{H}_{gt} is constructed by placing a Gaussian centered at the gaze point, with a standard deviation of 9.35 pixels. It is converted to two channels ($\mathbf{H}_{gt}, 1 - \mathbf{H}_{gt}$) to facilitate training using the cross-entropy loss detailed in Section 4.4.

F. Glossary

We provide definitions for key terms used in the paper.

- **Ego view:** For a person A wearing augmented reality glasses, this refers to the scene as observed from their own perspective.
- **Exo view:** This refers to the third-person observation of person A from another person’s (e.g., person B’s) perspective.
- **Looking at heads:** Defined for a pair of people. A person’s gaze point falls within another person’s head bounding box.
- **Eye contact:** Defined for a pair of people. A symmetric version of looking at heads, where the gaze points of both people fall within each other’s head box.
- **Shared attention:** Defined for two or more people. Occurs when multiple people look at the same object or person. In RLR-CHAT, this typically refers to multiple people looking at the same other person.
- **CLS token:** Standard transformer architectures, particularly the Vision Transformer (ViT) [4] used in our work, include an additional learnable token called the CLS token. This token is appended to the extracted image tokens to capture global information and is often used for tasks, such as classification, that require holistic image understanding.



Figure S4. Qualitative results on images from RLR-CHAT where all approaches yield similar predictions. In these cases, either a single salient target dominates the scene, making gaze estimation straightforward, or the gaze target is near the image center, reducing ambiguity across models.

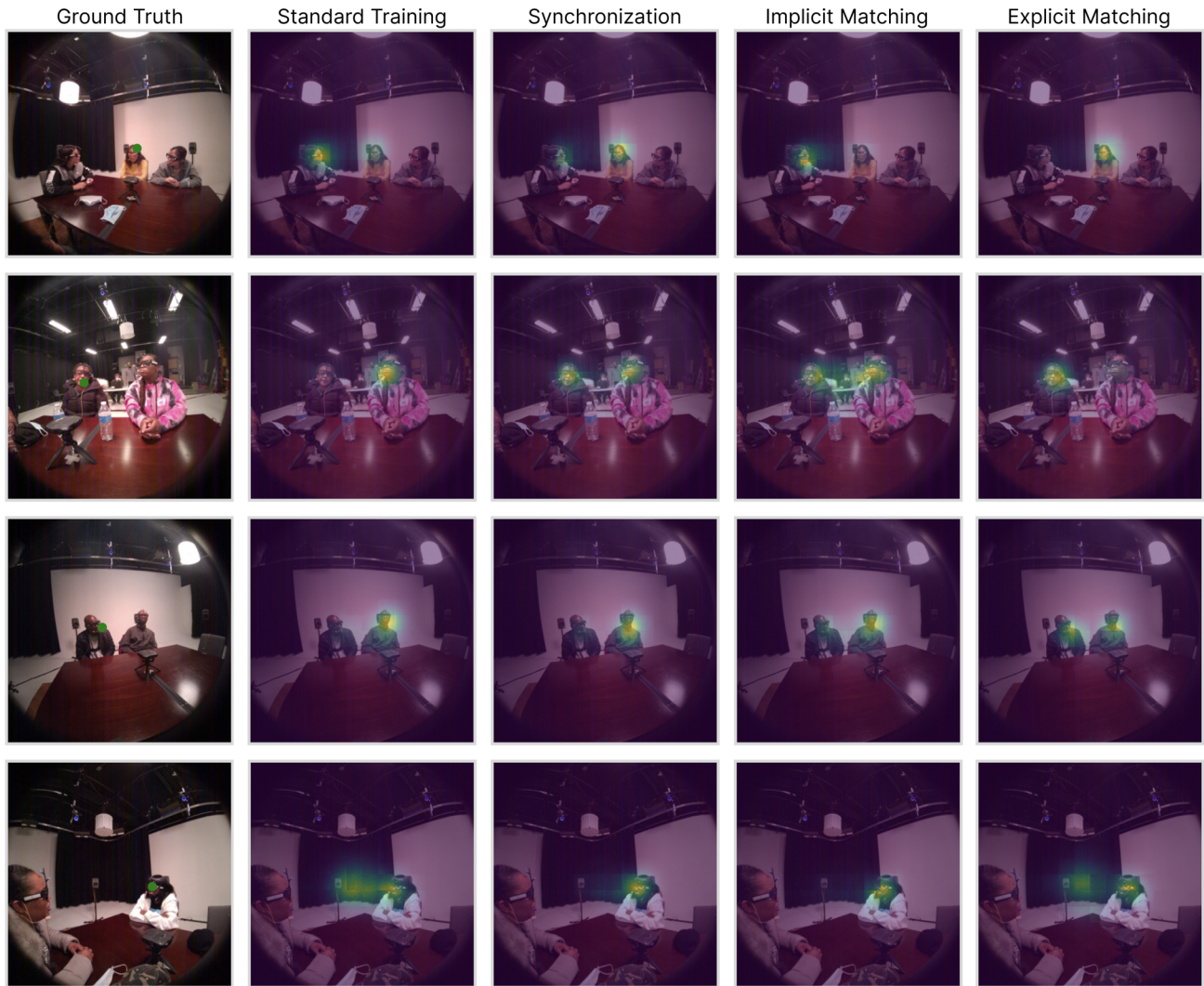


Figure S5. Qualitative results on images from RLR-CHAT where our proposed ego-exo alignment approaches improve egocentric gaze prediction compared to standard training. These cases involve greater ambiguity, with multiple salient targets positioned near the center, making gaze estimation more challenging.