

# Sequence-Based Identification of First-Person Camera Wearers in Third-Person Views

Anonymous CVPR submission

Paper ID XXXXX

## Abstract

001 *As immersive education and collaborative environments*  
 002 *continue to advance, understanding complex student in-*  
 003 *teractions within these shared spaces has become increas-*  
 004 *ingly important. While large-scale datasets like Ego4D and*  
 005 *Ego-Exo4D have advanced egocentric vision research, they*  
 006 *lack the rich, multi-user interactions critical for collabora-*  
 007 *tive learning and robotics. To address this gap, we intro-*  
 008 *duce TF2025, an expanded dataset featuring synchronized*  
 009 *first- and third-person views of actors, enhanced synchron-*  
 010 *ization, and multiple train-test splits. We also propose a*  
 011 *sequence-based approach for identifying first-person cam-*  
 012 *era wearers in broader third-person views. By leverag-*  
 013 *ing motion cues and person re-identification, our method*  
 014 *improves robustness and significantly outperforms state-of-*  
 015 *the-art approaches. This work advances the analysis of*  
 016 *multi-camera interactions in embodied vision and educa-*  
 017 *tion. The code and dataset will be made publicly available*  
 018 *upon acceptance.*

## 019 1. Introduction

020 Rapid advancements in virtual reality (VR), augmented re-  
 021 ality (AR) and narrative-centered learning technologies [6,  
 022 36, 46], have enabled the development of highly interac-  
 023 tive applications in immersive education [25, 40]. In these  
 024 settings, first-person (egocentric) cameras, such as body-  
 025 worn or head-mounted devices, capture dynamic, close-  
 026 range perspectives directly from the students' viewpoints.  
 027 In contrast, traditional third-person (exocentric) cameras,  
 028 such as those monitoring a classroom, provide a broader,  
 029 contextual view that complements this egocentric footage.  
 030 Combining information from these two types of views of-  
 031 fers significant potential for analyzing complex student-  
 032 student or student-teacher interactions in collaborative en-  
 033 vironments [29, 69].

034 Datasets like **Ego4D** [30] and **Ego-Exo4D** [32] have  
 035 been instrumental in advancing egocentric vision. Ego4D

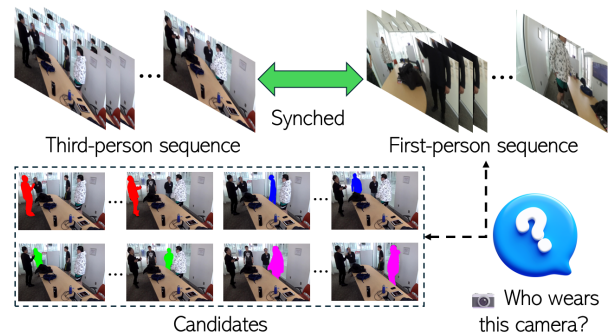


Figure 1. **Problem definition.**

offers a large-scale collection of first-person videos capturing 036  
 037 daily activities. Building on this, Ego-Exo4D pairs  
 038 egocentric footage with third-person (exocentric) views, 039  
 040 creating a multi-view dataset that enables researchers to  
 041 explore cross-view relationships—such as how robots or  
 042 first-person camera wearers can use widely available third-

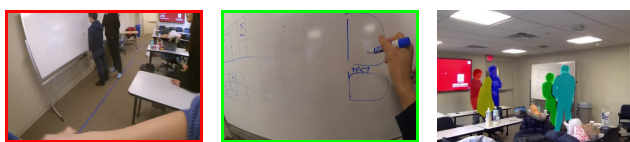
043 person videos to enhance learning [39].  
 044 However, Ego-Exo4D is limited to scenarios featuring  
 045 only a single camera wearer, restricting its applicability  
 046 in collaborative educational environments—such as multi-  
 047 user VR classrooms and immersive training [2, 77]—where  
 048 multiple students frequently interact and operate egocentric  
 049 devices simultaneously.

050 In this work, we study the identification of camera wears  
 051 in third-person views. In multi-user immersive learning en-  
 052 vironments, systems typically maintain basic identity in-  
 053 formation for students. However, simply knowing a stu-  
 054 dent's identity leaves a critical disconnect between their  
 055 first-person perspective and the overarching third-person  
 056 view. For example, an ID tag does not reveal how a stu-  
 057 dent's egocentric focus correlates with their body actions  
 058 and surroundings, which is crucial for activity recognition  
 059 and skill assessment. To create this essential link, we define  
 060 our task to identify the student's precise segmentation mask  
 061 in the third-person view using their first-person video as a  
 062 clue, as shown in Fig. 1. This mask-level cross-view match-

062 ing provides significantly more value for fine-grained action  
063 assessment, though it presents a major technical challenge.

064 Several methods [26, 79, 83] have been proposed for this  
065 task. However, these approaches are frame-based, mak-  
066 ing predictions using only frame-level information. Due to  
067 the inherent challenges of first-person cameras, such as oc-  
068 clusions and motion blur, relying on single frames reduces  
069 reliability. Additionally, since camera wearers frequently  
070 move, the first-person view can sometimes become unin-  
071 formative. For instance, when the camera points at blank  
072 walls, white boards or other featureless surfaces.

073 In this work, we introduce a sequence-based approach  
074 for identifying first-person camera wearers in third-person  
075 views. Compared to frame-based methods, our approach of-  
076 fers several key advantages. First, by using multiple frames,  
077 the model reduces reliance on any single uninformative  
078 frame. Second, our experiments demonstrate that matching  
079 the motion of camera wearers in third-person views with op-  
080 tical flow in first-person views provides strong and reliable  
081 evidence about identity. Lastly, using a sequence of frames  
082 increases the chance of camera wearers seeing each other,  
083 which leads to an improved identification clue.



(a) First-person view 1 (b) First-person view 2 (c) Third-person view

Figure 2. **Data sample from TF2025.** The TF2025 dataset aims to establish a link between students’ first-person views and their third-person segmentation masks in classrooms.

084 To support this objective, we introduce **TF2025 (Third-**  
085 **First 2025)**, an expanded dataset designed to benchmark al-  
086 gorithms that facilitate these collaborative environments, as  
087 shown in Fig. 2. TF2025 integrates videos from three source  
088 datasets: **TF2023** [83], **IUShareView** [79], and **Ego4D-**  
089 **TF**, a subset of Ego4D enhanced with our additional an-  
090 notated ground truth masks. TF2025 is 2.3 times the size  
091 of TF2023, the current largest dataset for this task, and in-  
092 cludes more challenging frames. It also features three new  
093 train/test splits designed to evaluate multiple levels of gen-  
094 eralizability, directly translating to diverse educational and  
095 training applications. Furthermore, we conducted manual  
096 synchronization to improve temporal alignment between  
097 first- and third-person views.

098 To summarize, our main contributions are as follows:

- 099 1. We introduce TF2025, a larger and more challenging  
100 dataset featuring multiple new evaluation benchmarks  
101 tailored for complex multi-user interactions.
- 102 2. We propose *motion matching*, a novel approach for first-  
103 and third-person matching.

3. We present a new framework termed *Motion Appearance  
Fusion (MAF)* that adaptively fuses information from  
multiple cues, outperforming state-of-the-art methods.

## 2. Related Works

**Egocentric Video Understanding** has advanced computer  
vision by capturing first-person experiences for applica-  
tions like activity recognition [17, 28, 30, 55, 57, 70], ob-  
ject interaction [30, 50, 54, 63], and navigation. Recent  
studies leverage self-supervised learning to enhance action  
recognition [4, 49, 53, 60, 78], while others focus on gaze  
prediction [37, 38, 47, 58, 75] and hand-object interac-  
tions [74, 80] to infer human intent. Egocentric videos also  
aid robotics in task and imitation learning [41–43]. In so-  
cial analysis, they help detect engagement, infer relation-  
ships, and identify conversational partners [3, 5, 11, 30, 83].  
Large-scale datasets like Ego4D [30] have further driven  
research, supporting a wide range of applications from  
human-robot collaboration to social behavior modeling.

**Cross-View Matching** aligns corresponding elements  
across different views to enhance spatial understanding.  
Recent advancements leverage Graph Neural Networks  
(GNNs) and attention mechanisms for improved feature  
matching. SuperGlue [64] uses GNNs with positional en-  
codings to refine keypoint descriptors, while LoFTR [72]  
employs self- and cross-attention for feature alignment.  
ASpanFormer [15] introduces a hierarchical attention with  
adaptive attention span adjustment. Other approaches en-  
hance efficiency [14, 66] and unify dense and sparse match-  
ing methods [24, 76]. GlueStick [61] refines correspon-  
dence estimation by jointly matching points and lines.

**Joint First- and Third-person Video Understanding**  
Recent research has made significant strides in connect-  
ing first-person (egocentric) and third-person (exocentric)  
video understanding. Approaches include cross-view ac-  
tion recognition [8, 67, 71], generating egocentric video  
from third-person footage, aligning frames, and learning  
viewpoint-invariant representations [48]. Cross-view data  
has improved pose estimation [1, 16, 54]. Other work fo-  
cuses on summarizing egocentric videos with third-person  
views [20, 35], refining egocentric models using third-  
person data [54, 62], and co-segmentation via attention pre-  
diction. Datasets with both views include Ego-Exo4D [31],  
Assembly101 [65], and Charades-Ego [68]. Other datasets  
like Epic-Kitchens [17–19] also combine these views, sup-  
porting gaze prediction and activity modeling.

Person identification has been explored in the context  
of first- and third-person cross-view understanding. While  
some work has focused on top-view cameras [7, 9, 10, 81],  
our work aligns with the use of side-view cameras as third-  
person views [26, 79, 83]. Although side-view cameras pro-  
vide a narrower field of view compared to top-view cam-  
eras, they offer richer appearance details for the camera

156 wearer candidates. Recent advancements, such as multi-  
157 branch deep networks and semi-Siamese architectures, have  
158 improved joint embeddings for aligning first- and third-  
159 person views [26, 79], while PEN [83] introduces a dual-  
160 branch framework, combining a personal branch for rela-  
161 tional analysis and an environment branch for geometric  
162 matching. Building on these foundations, our work focuses  
163 on identifying camera wearers in third-person views by in-  
164 corporating temporal motion cues and multi-cue informa-  
165 tion fusion to achieve better performance.

### 166 3. Methodology

#### 167 3.1. Problem Definition

168 We define the task of identifying the camera wearer in third-  
169 person views as follows: Given a  $t$ -frame sequence of first-  
170 person frames  $Ego : \{Ego_1, Ego_2, \dots, Ego_t\}$ , a sequence  
171 of third-person frames  $Exo : \{Exo_1, Exo_2, \dots, Exo_t\}$ ,  
172 and  $N$  mask sequences in the third-person view represent-  
173 ing  $N$  candidate people,

174 Candidate 1:  $\{Mask_{11}, Mask_{12}, \dots, Mask_{1t}\}$   
175 Candidate 2:  $\{Mask_{21}, Mask_{22}, \dots, Mask_{2t}\}$   
176  $\vdots$   
177 Candidate N:  $\{Mask_{N1}, Mask_{N2}, \dots, Mask_{Nt}\}$ ,

178 the goal is to identify the mask sequence which represents  
179 the correct candidate that corresponds to the first-person  
180 frame sequence  $Ego : \{Ego_1, Ego_2, \dots, Ego_t\}$ . In our  
181 setting, we choose  $t = 30$  as a balance between hav-  
182 ing enough frames for meaningful information and avoid-  
183 ing excessive computational cost; however, our introduced  
184 method can be applied to sequences of arbitrary length.

185 This task presents unique and significant challenges. Un-  
186 like traditional person identification or re-identification [12,  
187 82, 84], we cannot rely on the appearance of the cam-  
188 era wearers themselves, which is typically the most robust  
189 source of information. As egocentric camera wearers are  
190 often at the center of interactions, their view presents addi-  
191 tional challenges: limited scene coverage due to the cam-  
192 era’s closeness to the center of interaction, frequent occlu-  
193 sions from nearby objects or people, and motion blur from  
194 quick movements. These issues can substantially affect the  
195 quality and reliability of the first-person view, complicating  
196 information extraction.

197 Prior work [83] shows the importance of multi-source in-  
198 tegration for this task, introducing a frame-based approach  
199 that utilizes geometric cues from image matching and per-  
200 son detection in both views. However, image matching is  
201 often too unreliable, as methods like GlueStick [61] rely on  
202 significant overlap between views—an assumption that fre-  
203 quently breaks between first-person and third-person views.

204 This limitation highlights the need for more robust tech-  
205 niques to address the differences in these views.

206 In this section, we present a novel sequence-based  
207 framework, *Motion Appearance Fusion (MAF)*. First, we  
208 introduce a new matching method based on camera motion,  
209 termed *motion matching*. Next, we describe our *appear-*  
210 *ance matching* model, which computes appearance simi-  
211 larity between individuals in first-person and third-person  
212 views. Finally, we propose the *Confidence-Based Adaptive*  
213 *Fusing (CBAF)* module to integrate information from these  
214 two sources. The structure of *MAF* is illustrated in Fig. 3.

#### 215 3.2. Motion Matching

216 First, we establish a connection between the body motion  
217 of camera wearers in the third-person view and the optical  
218 flow in the first-person view. While optical flow has been  
219 widely used in egocentric vision analysis [22, 59, 79], to  
220 our knowledge, we are the first to explicitly link it to third-  
221 person body motion for first-third cross-view matching.

222 We start by discussing the concept of motion within the  
223 context of this paper. Some egocentric datasets, such as  
224 **Ego4D** and **Ego-Exo4D**, utilize head-mounted and eyewear  
225 cameras, while others like **IUShareView** and **TF2023** use  
226 chest-worn cameras. Our approach does not restrict the pre-  
227 diction of motion to any predefined body part. Instead, we  
228 seek to predict movements from any body part that could  
229 induce changes (optical flow) in the first-person view. This  
230 design choice enhances the generalizability of our model, as  
231 body-worn cameras in real-world scenarios can be mounted  
232 on various body parts, such as the shoulder, wrist, or waist.

233 In addition, we define two values to model the camera  
234 motion:  $T$  and  $R$ . Camera motion can be represented by  
235 translational and rotational motions, which have distinct im-  
236 pacts on the optical flow in the first-person view. For trans-  
237 lational motion, due to motion parallax, objects closer to  
238 the camera exhibit faster apparent motion across the field of  
239 view compared to those farther away, resulting in a depth-  
240 dependent optical flow pattern. In contrast, for rotational  
241 motion, the optical flow is depth-independent.

242 To model this, we apply a sequence-based backbone,  
243 *MViT* [27], to a candidate in the third-person view, which is  
244 cropped using the third-person frames and the candidate’s  
245 segmentation mask. We then use a Multi-Layer Perceptron  
246 (MLP) to predict two values,  $T_{exo}$  and  $R_{exo}$ , which together  
247 form a system that models the egocentric camera motion.

248 For the first-person view, we first compute two values,  
249  $T_i$  and  $R_i$  for each frame  $i$ ,

$$250 T_i = \operatorname{median}_{j \in \text{frame}_i} \left( \operatorname{depth}(p_{ij}) \sqrt{fl_x(p_{ij})^2 + fl_y(p_{ij})^2} \right) \quad (1)$$

$$252 R_i = \operatorname{median}_{j \in \text{frame}_i} \left( \sqrt{fl_x(p_{ij})^2 + fl_y(p_{ij})^2} \right), \quad (2)$$

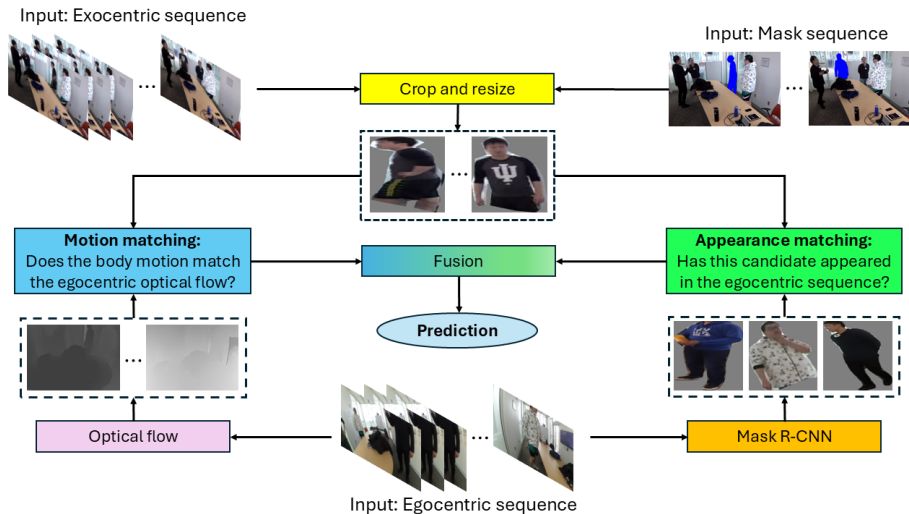


Figure 3. **Overview of the structure of our framework (MAF)** We utilize two models: *motion matching*, which matches the motion of the candidate to the egocentric optical flow, and *appearance matching* that checks if the candidate appeared in the first-person view. Then we integrates information from both sides to make the final prediction.

253 where  $p_{ij}$  denotes the  $j$ -th pixel in frame  $i$ , and  $depth(p_{ij})$ ,  
 254  $fl_x(p_{ij})$ , and  $fl_y(p_{ij})$  represent the depth, x-axis optical  
 255 flow, and y-axis optical flow, respectively.

256 We then compute the total motion of the first-person  
 257 view by aggregating the two motion values  $T$  and  $R$  across  
 258 all frames, represented as a tuple of the cumulative motion,  
 259 as shown in Eq. 3,

$$260 \text{motion}(\{Ego_1, \dots, Ego_t\}) = \left( \sum_{i=1}^{t-1} T_i, \sum_{i=1}^{t-1} R_i \right). \quad (3)$$

261 In our model, we utilize *Raft* [73] to compute optical  
 262 flow, and *ZoeDepth* [13], a monocular depth estimation  
 263 model, to predict absolute depth. Additionally, we employ  
 264 *Mask R-CNN* [33] to filter out foreground pixels where peo-  
 265 ple are detected, as their presence can introduce indepen-  
 266 dent motion not caused by the camera wearer.

267 For training, we randomly select 16 consecutive frames  
 268 from a 30-frame sequence. Using our motion model, we  
 269 predict the motion values  $T_{exo}$  and  $R_{exo}$  for the selected  
 270 sequence. The loss function is computed as the sum of  
 271 squared errors between the predicted motion values and the  
 272 corresponding ground truth values derived from the respec-  
 273 tive first-person sequence.

274 During inference, we employ a sliding window approach  
 275 to evaluate motion similarity across overlapping sequences.  
 276 Specifically, we calculate the sum of squared errors for the  
 277 motion values across three overlapping sequences: frames  
 278 1–16, 8–23, and 15–30. This sum of squared errors is de-  
 279 noted as “motion score”. A detailed illustration of the *mo-*  
 280 *tion matching* model is provided in the supplementary ma-  
 281 terials.

### 3.3. Appearance Matching 282

283 The presence of individuals in the first-person view pro-  
 284 vides additional cues for matching, as camera wear-  
 285 ers should not see themselves in their own first-person  
 286 views. This information has been utilized in prior work  
 287 on third/first-person matching. *PEN* [83] introduces a per-  
 288 sonal branch that implicitly learns this information through  
 289 cross-attention between third- and first-person views. Their  
 290 model takes as input the third-person candidate, the indi-  
 291 viduals visible in the corresponding first-person view, and a  
 292 binary training label indicating whether the pair is a match.

293 However, we empirically discovered that this approach  
 294 is suboptimal, as the visibility of the camera wearer in the  
 295 first-person view does not strictly correlate with their iden-  
 296 tity. For instance, even if a camera wearer is not visible in  
 297 their own first-person view, it does not necessarily imply  
 298 that the view belongs to them. In this work, we demonstrate  
 299 empirically that training a person re-identification (re-ID)  
 300 model directly is more effective for this task.

301 We implement our person re-ID model, denoted *appear-*  
 302 *ance matching*, using a *Vision Transformer (ViT)* [23] back-  
 303 bone. During training, the backbone outputs a feature vec-  
 304 tor, and we attach an MLP head to generate classification la-  
 305 bels. During inference, we use the feature vector to compute  
 306 the L2 distance between pairs of individuals. Following ex-  
 307 isting person re-ID works [34], we use three data augmenta-  
 308 tion techniques: random horizontal flip, random cropping,  
 309 and random patching, each applied with a 50% probability.  
 310 We apply two loss functions: cross-entropy loss for classi-  
 311 fication and triplet loss for metric learning.

312 While it is natural to use the camera wearer as the anchor

313 and another person in their view as the negative sample, this  
 314 approach may result in non-camera wearers only appearing  
 315 as negatives during training. To address this imbalance, we  
 316 introduced “pseudo” first/third-person pairs by using non-  
 317 camera wearers as anchors and the camera wearer as the  
 318 negative sample, as depicted in Fig. 4.

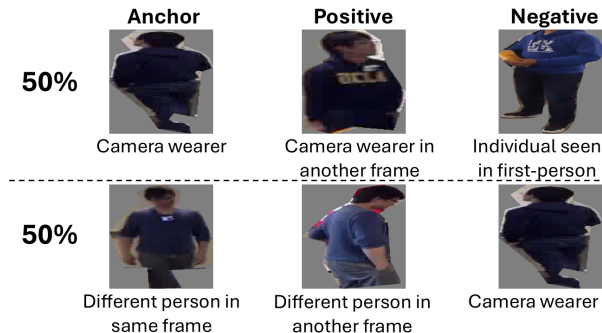


Figure 4. Triplet selection.

319 During inference, we select frames 1, 16, and 30 from  
 320 the first-person sequence and apply *Mask R-CNN* to crop  
 321 out all detected individuals. Each detected individual and  
 322 third-person candidates are passed through the same *ViT*  
 323 backbone, and we compute its average L2 distance from the  
 324 third-person candidate across all frames, denoted as “ap-  
 325 pearance score”. These three frames were chosen as a bal-  
 326 ance between capturing most individuals appearing in the  
 327 first-person sequence and avoiding excessive detections of  
 328 the same person, which could complicate the fusion pro-  
 329 cess. We provide a detailed illustration of the *appearance*  
 330 *matching* model in the supplementary materials.

### 331 3.4. Confidence-Based Adaptive Fusing

332 After *motion matching* and *appearance matching*, we have  
 333 the following results for  $N$  candidates in third-person and  
 334  $M$  detected individuals in the first-person sequence, where  
 335 motion scores and appearance scores represent the results  
 336 from the two models:

337 Motion scores:  $\{Motion_1, \dots, Motion_N\}$   
 338 Appear. scores 1:  $\{Appearance_{11}, \dots, Appearance_{1N}\}$   
 339  $\vdots$   
 340 Appear. scores  $M$ :  $\{Appearance_{M1}, \dots, Appearance_{MN}\}$

341 The goal of this task is to identify the correct candidate  
 342 among  $\{1, \dots, N\}$  that matches the first-person sequence.  
 343 This poses several challenges for the fusion method:

- 344 1. Both  $M$  and  $N$  are different for each query.
- 345 2. The same individual can be detected multiple times in  
 346 the 3 chosen frames.
- 347 3. Motion scores and appearance scores behave differently.  
 348 A lower motion score reflects better matching between

the third/first-person motion, while a lower appearance  
 score means the candidate is seen in the first-person  
 view, so unlikely to be the camera wearer.

We propose *Confidence-Based Adaptive Fusing (CBAF)*,  
 a novel module designed to integrate results from the two  
 preceding modules. *CBAF* adaptively selects which source  
 of information to use for candidate elimination or final pre-  
 diction, requires no training for different datasets, and can  
 handle arbitrary values of  $M$  and  $N$ .

We first define the concept of confidence in the context  
 of this paper:

$$\text{Confidence} = \lambda_{\text{Re-ID}} \cdot \alpha_{\text{mask}} \cdot \frac{x_{(2)}}{x_{(1)}}$$

where:

- $x_{(1)}$  is the smallest value among  $N$  values from the same  
 source,
- $x_{(2)}$  is the second smallest value among  $N$  values from  
 the same source,
- $\lambda_{\text{Re-ID}}$  is a constant representing our trust in the *appear-  
 ance matching* model, which depends on prior knowledge  
 about the dataset (e.g., how clearly individuals are visible  
 in the first-person view). For motion scores,  $\lambda_{\text{Re-ID}} = 1$ ,
- $\alpha_{\text{mask}}$  is the confidence from *Mask R-CNN* when detecting  
 this individual. For motion scores,  $\alpha_{\text{mask}} = 1$ .

Essentially, this calculates the ratio between the two  
 smallest items from each score source. A higher confi-  
 dence value indicates a clearer distinction between the most  
 prominent candidate and the runner-up, while a lower confi-  
 dence value suggests ambiguity between the two candi-  
 dates. This ratio helps us evaluate the reliability of each  
 source of information, guiding how much trust we place in  
 their respective results.

The *CBAF* algorithm is presented in Algorithm 1. It  
 compares the confidence of motion scores against that of  
 all appearance confidence scores. If the motion scores ex-  
 hibit the highest confidence, the candidate with the lowest  
 motion score is selected as the prediction. Conversely, if an  
 appearance source yields the highest confidence, the corre-  
 sponding candidate is removed from all sources, and that  
 source is subsequently eliminated. Examples illustrating  
 this algorithm are provided in the supplementary materials.

## 4. Experiments

### 4.1. Dataset

We introduce **TF2025 (Third-First 2025)**, an expanded  
 dataset built on the foundation of three datasets: **TF2023**,  
**IUShareView**, and **Ego4D**.

We begin by merging TF2023 with IUShareView, as  
 these two datasets were both recorded in classroom settings,  
 utilize the same camera models, and feature overlapping ac-  
 tors. IUShareView [26, 79] was designed for the task of

**Algorithm 1** Confidence-Based Adaptive Fusing (CBAF)

---

**Require:** Motion scores  $\{Motion_1, \dots, Motion_N\}$ , appearance scores  $\{Appearance_{i1}, \dots, Appearance_{iN}\}$  for  $i = 1, \dots, M$ , Candidates  $\{1, \dots, N\}$

- 1: **while**  $\text{len}(Candidates) > 1$  **do**
- 2:   Compute motion confidence:  $C_{\text{motion}} \leftarrow \frac{Motion_{(2)}}{Motion_{(1)}}$
- 3:   Initialize best source  $s \leftarrow 0$ , best confidence  $C_{\text{best}} \leftarrow C_{\text{motion}}$
- 4:   **for** each appearance source  $i = 1, \dots, M$  **do**
- 5:     Compute appearance confidence:  $C_{\text{Re-ID}_i} \leftarrow \lambda_{\text{Re-ID}} \cdot \alpha_i \cdot \frac{Appearance_{i(2)}}{Appearance_{i(1)}}$
- 6:     **if**  $C_{\text{Re-ID}_i} > C_{\text{best}}$  **then**
- 7:       Update  $s \leftarrow i$ ,  $C_{\text{best}} \leftarrow C_{\text{Re-ID}_i}$
- 8:     **end if**
- 9:   **end for**
- 10:  **if**  $s = 0$  **then**
- 11:    **return**  $Candidates[\arg \min(Motion)]$
- 12:  **else**
- 13:    Remove the candidate with index  $\arg \min(s)$  from all sources and  $Candidates$
- 14:    Remove  $\{Appearance_{s1}, \dots, Appearance_{sN}\}$  from all scores
- 15:  **end if**
- 16: **end while**
- 17: **return**  $Candidates[0]$

---

399 identifying first-person camera wearers, containing 9 sets  
400 of frame sequences. Later, TF2023 [83] was introduced and  
401 is currently the largest dataset featuring synchronized first-  
402 and third-person views, along with segmentation masks for  
403 every individual in third-person views. Both datasets were  
404 designed to study multi-camera wearer interactions in so-  
405 cial scenarios, such as presentations, discussions and snack  
406 breaks. To support this focus, the official splits of both  
407 datasets include only those frames where actors are actively  
408 engaged in interactions. In TF2025, we address a more  
409 realistic scenario by extracting all frames from the source  
410 videos of these two datasets and constructing continuous  
411 30-frame sequences. The newly incorporated frames intro-  
412 duce significantly greater challenges compared to existing  
413 ones, as examples shown in Fig. 5.

414 Ego4D [30], launched in 2021, is a large-scale egocentric  
415 dataset with over 3,670 hours of first-person video captured  
416 using wearable cameras. Ego4D does not inherently sup-  
417 port this task, as it primarily consists of egocentric videos  
418 featuring solo camera wearers. To address this, we created  
419 and annotated **Ego4D-TF**, a subset for cross-view match-  
420 ing. We identified scenes in Ego4D containing at least three  
421 people and at least two synchronized egocentric videos. We  
422 then sampled one of every six frames to match the frame  
423 rate of TF2023 and IUShareview and selected one of the



Figure 5. Examples of highly challenging first-person frames including camera wearers looking at a whiteboard, ceiling, or being occluded by their own arms.

first-person views as a “pseudo-third-person” view. Using  
YOLOv11 [45], we extracted segmentation masks for visi-  
ble individuals and tracked them across frames for each 30-  
frame sequence. We provide a more comprehensive expla-  
nation and examples of our annotation method in the sup-  
plementary materials.

**Manual Synchronization:** We performed manual cali-  
bration to improve the synchronization accuracy of TF2023  
and IUShareView using a technique inspired by the manual  
synchronization procedure used in **Ego-Exo4D**. By iden-  
tifying clear single-frame actions visible across multiple  
cameras, such as moments when a hand touches specific ob-  
jects, we aligned their timestamps across different camera  
views. This process reduced the temporal alignment error  
by up to 4 frames (0.8 seconds).

**Cross-dataset Evaluation:** Prior work [83] identified a  
phenomenon in this task where some methods may unin-  
tentionally exploit dataset biases in third-person views. For  
instance, in TF2023, camera wearers often appear more fre-  
quently near the center of frames, leading models to learn  
visual cues that are not transferable to other datasets or ap-  
plications. To address this issue, we constructed a test split  
using Ego4D-TF, enabling cross-dataset evaluation to better  
assess generalizability.

**Dataset Splits:** Following prior works in first- and third-  
person learning [32, 54], we created three splits designed to  
evaluate methods at different levels of generalizability. We  
believe this is crucial for evaluating diverse potential ap-  
plications. For instance, in an immersive learning environ-  
ment where multiple students engage in VR/AR group activi-  
ties, the model should be trainable using data from the class-  
room and the appearances of the participants. In contrast, in pub-  
lic settings such as detecting unauthorized camera wearers,  
there may be no prior knowledge of the location or even the  
camera models used.

We designed three levels of splits, as outlined below:

1. **Seen:** We allocated the first 80% of each source video  
from TF2023 and IUShareView for training and the re-  
maining 20% for testing.
2. **Unseen:** We split the videos from TF2023 and IUShare-  
View such that neither the camera wearers nor the loca-  
tions overlap between training and testing.
3. **Cross-Dataset:** We used data from TF2023 and  
IUShareView for training and Ego4D-TF for testing.

| Methodology               | TF2023                  | Seen(Acc $\uparrow$ ) | TF2025                  | Cross-dataset(Acc $\uparrow$ ) |
|---------------------------|-------------------------|-----------------------|-------------------------|--------------------------------|
|                           | Unseen(Acc $\uparrow$ ) |                       | Unseen(Acc $\uparrow$ ) |                                |
| Random guess              | 32.5%                   | 28.7%                 | 33.1%                   | 46.4%                          |
| Sequence baseline         | 34.7%                   | 39.4%                 | 34.1%                   | 47.0%                          |
| ThirdFirstNet [79]        | 35.4%                   | 31.1%                 | 33.8%                   | 47.4%                          |
| PEN [83]                  | 67.2%                   | 85.9%                 | 58.8%                   | 52.9%                          |
| Motion matching(ours)     | 84.6%                   | 88.8%                 | 80.5%                   | 67.7%                          |
| Appearance matching(ours) | 79.5%                   | 66.1%                 | 63.4%                   | 72.2%                          |
| MAF(ours)                 | <b>90.3%</b>            | <b>93.8%</b>          | <b>84.4%</b>            | <b>75.7%</b>                   |

Table 1. **Experiment Results on TF2025 and TF2023.** Models were evaluated by matching a first-person sequence to the correct third-person candidate, given a third-person sequence and multiple candidate masks. Results are reported in accuracy. We converted the original TF2023 dataset to 30-frame sequences in order to compare with sequence-based methods.

468 This ensures no shared individuals, locations, or camera  
469 models between the training and testing sets.

470 A detailed comparison of the sizes of TF2023 and each  
471 split of TF2025 is provided in Tab. 2.

| Dataset             | TF2023       | TF2025 |        |               |
|---------------------|--------------|--------|--------|---------------|
| Split               | Unseen       | Seen   | Unseen | Cross-Dataset |
| Number of Frames    |              |        |        |               |
| Train               | 47335        | 104760 | 80490  | 131010        |
| Test                | 21689        | 26250  | 50520  | 32880         |
| Total               | <b>69024</b> | 131010 | 131010 | <b>160890</b> |
| Number of Sequences |              |        |        |               |
| Train               | N/A          | 3492   | 2683   | 4367          |
| Test                | N/A          | 875    | 1684   | 1096          |
| Total               | <b>N/A</b>   | 4367   | 4367   | <b>5363</b>   |

Table 2. **Detailed Comparison of dataset sizes.** TF2025 is approximately 2.3 times larger than TF2023, offering multiple evaluation splits and providing full 30-frame sequences rather than sporadic frame sets.

## 472 4.2. Implementation Details and Results

473 We implemented three baseline methods for comparison:  
474 two frame-based methods and one sequence-based method.

475 **PEN** [83] is the current state-of-the-art method for this  
476 task. We applied it to each of the 30 frames and summed its  
477 prediction scores to generate a sequence-level prediction.

478 **ThirdFirstNet** [79] originally requires segmentation  
479 masks from the previous frame, making it unsuitable for  
480 this task. Following the adaptation in [83], we modified it to  
481 use the current frame’s masks and the model’s identification  
482 branch. Frame-level predictions were summed to generate  
483 sequence-level results.

484 Since no sequence-based methods exist for this task, we  
485 implemented a *two-stream* baseline using two sequence-

based backbones (*MViT*): one for third-person and one for  
first-person sequences. The third-person backbone was  
modified to accept 4-channel inputs by concatenating the  
candidate mask sequence with the third-person sequence.  
Both sequences are fed into their respective backbones and  
trained with a contrastive loss: matching sequences mini-  
mize the L2 distance between outputs, while non-matching  
sequences maximize it. During inference, we also em-  
ployed the sliding window approach and selected the can-  
didate mask sequence with the minimal cumulative L2 dis-  
tance to the first-person sequence.

For our proposed method (*MAF*), we trained our *motion*  
*matching* and *appearance matching* models separately us-  
ing the AdamW [52] optimizer with a learning rate of 1e-4.  
Both models were trained to process input images and se-  
quences with a resolution of 224 $\times$ 224 pixels. The *motion*  
*matching* model was trained with a batch size of 32 for 50  
epochs, while the *appearance matching* model used a batch  
size of 64 for 100 epochs. The *MViT* backbone for *motion*  
*matching* was pretrained on Kinetics-400 [44], and the *ViT*  
backbone for *appearance matching* was pretrained on Im-  
ageNet [21]. All experiments were conducted on two Nvidia  
RTX A6000 GPUs. We set the hyperparameter  $\lambda_{\text{Re-ID}}$  to  
1.75, 1.75 and 3 for the three splits, supported by sensitivity  
study in supplement materials.

The experimental results are shown in Tab. 1. *MAF*  
significantly outperforms frame-based methods across all  
splits, particularly in the unseen and cross-dataset splits,  
where higher levels of generalizability are required. On  
the original TF2023 dataset, *PEN* achieved an accuracy of  
67.2%. However, its performance dropped to 58.8% on the  
unseen split of TF2025 due to the inclusion of more chal-  
lenging frames.

We observed that while *PEN* achieved significant im-  
provements over random guessing on both the seen and un-  
seen splits, it failed to outperform random guessing by a  
large margin on the cross-dataset split. We attribute this  
to *PEN*’s reliance on geometric matching, which requires a

| Dataset     | TF2025       |              |               |
|-------------|--------------|--------------|---------------|
|             | Seen         | Unseen       | Cross-dataset |
| $T$ only    | 85.1%        | 73.8%        | 64.2%         |
| $R$ only    | 83.0%        | 75.3%        | 66.7%         |
| Both (ours) | <b>88.8%</b> | <b>80.5%</b> | <b>67.7%</b>  |

Table 3. Ablation study on motion calculation.

524 portion of the view to be visible in both the first-person and  
 525 third-person views. Essentially, it assumes that the camera  
 526 wearer is facing away from the third-person camera. How-  
 527 ever, in the Ego4D-TF setting, where participants often face  
 528 each other in board games, this assumption does not hold,  
 529 showing its lack of robustness.

530 The introduced *two-stream* baseline achieved a slight  
 531 improvement over random guessing on the seen split but  
 532 did not surpass random guessing on the unseen and cross-  
 533 dataset splits by any notable margin. This demonstrates that  
 534 the significant improvement of our model stems from its  
 535 design of visual cues, rather than the sequence-based back-  
 536 bones themselves.

### 537 4.3. Ablation studies

538 We conducted several ablation studies to justify the choice  
 539 of our model design in each modules.

540 First, we validated our design to predict two motion val-  
 541 ues separately. We compared this approach to predicting  
 542 only one of the two motion values, i.e.,  $T$  or  $R$ . As shown  
 543 in Tab. 3, our method outperformed single-source predic-  
 544 tion across all three splits.

545 Second, we compared our *appearance matching* model  
 546 to the *personal branch* introduced in *PEN* [83], as both  
 547 methods utilize the appearances of individuals in first-  
 548 person views. The results are shown in Tab. 4. We observed  
 549 that the *personal branch* of *PEN* significantly outperformed  
 550 our method on the seen split. This is because the *personal*  
 551 *branch* model takes both the third-person candidate and the  
 552 individual in the first-person view as inputs, learning their  
 553 differences implicitly through cross-attention. As a result,  
 554 for the seen split, it can memorize the appearance of third-  
 555 person camera wearers and make predictions even when no  
 556 individual appears in the first-person sequence. In contrast,  
 557 our method takes in the L2 distance between individuals  
 558 from both views. If no individual appears in the first-person  
 559 sequence, our *appearance matching* method defaults to ran-  
 560 dom guessing. However, this does not negatively impact our  
 561 overall performance, as our fusion method make predictions  
 562 based on the *motion matching* results when this happens.

563 On the other two splits, where camera wearers did  
 564 not appear during training (unseen and cross-dataset), our  
 565 method significantly outperformed *PEN*, demonstrating  
 566 much better generalizability. This finding also underscores

| Dataset           | TF2025       |              |               |
|-------------------|--------------|--------------|---------------|
|                   | Seen         | Unseen       | Cross-dataset |
| PEN-Personal      | <b>83.0%</b> | 50.3%        | 48.3%         |
| Appearance (ours) | 66.1%        | <b>63.4%</b> | <b>72.2%</b>  |

Table 4. Ablation study on appearance matching.

| Methodology | Dataset | TF2025 |              |              |               |
|-------------|---------|--------|--------------|--------------|---------------|
|             |         | Frames | Seen         | Unseen       | Cross-dataset |
| No fusing   | N/A     |        | 88.8%        | 80.5%        | 67.7%         |
| Direct      | 16      |        | 77.6%        | 66.7%        | 67.3%         |
| CBAF(ours)  | 16      |        | 91.5%        | 82.2%        | 72.9%         |
| Subtraction | 1,16,30 |        | 89.9%        | 79.7%        | 73.1%         |
| CBAF(ours)  | 1,16,30 |        | <b>93.8%</b> | <b>84.4%</b> | <b>75.7%</b>  |

Table 5. Ablation study on fusing methods.

567 the benefit of having multiple levels of splits in the TF2025  
 568 dataset for evaluation.

569 Finally, we evaluated the performance of our fusion  
 570 method (*CBAF*). We compared it against two other meth-  
 571 ods: *subtraction fusing*, which subtracts appearance scores  
 572 from motion scores, and *direct fusing*, which uses individu-  
 573 als detected in the first-person view to eliminate candidates  
 574 based on the order of *Mask R-CNN* confidence scores and  
 575 then selects the best candidate among the remaining, us-  
 576 ing motion scores. Since *direct fusing* cannot handle cases  
 577 where the same person appears multiple times, we also  
 578 modified our *CBAF* method to use only frame 16 for com-  
 579 parison. The results are shown in Tab. 5.

## 580 5. Conclusion

581 We explore the task of identifying first-person camera  
 582 wearers in third-person views. We introduce a new  
 583 dataset, **TF2025**, which integrates two existing datasets  
 584 (**IUShareView** and **TF2023**) with additional frames, along  
 585 with **Ego4D-TF**, an annotated subset of **Ego4D**. The  
 586 dataset also includes three new splits to evaluate different  
 587 levels of generalizability. We propose a novel sequence-  
 588 based framework, *MAF*, which incorporates a new *motion*  
 589 *matching* method and an *appearance matching* model. Ad-  
 590 ditionally, we introduce a novel method to integrate results  
 591 from these two models. Experiments show that our ap-  
 592 proach significantly outperforms existing works, achieving  
 593 state-of-the-art performance. Future work could incorpo-  
 594 rate more egocentric datasets to increase the diversity of  
 595 camera wearer appearances and camera models. This could  
 596 also lead to improved models for measuring camera motions  
 597 and better generalization across various camera models.  
 598

599

## References

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

[1] Adnan Abdullah, Ruo Chen, Ioannis Rekleitis, and Md Jahidul Islam. Ego-to-exo: Interfacing third person visuals from egocentric views in real-time for improved rov teleoperation. *arXiv preprint arXiv:2407.00848*, 2024. 2

[2] Laure Abensur Vuillaume, Jonathan Goffoy, Nadège Dubois, Nathacha Almoynes, Cécile Bardet, Evelyne Dubreucq, Sophie Klenkenberg, Anne-Françoise Donneau, Camille Dib, Alexandre Ghuysen, et al. Collaborative virtual reality environment in disaster medicine: moving from single player to multiple learners. *BMC medical education*, 24(1):422, 2024. 1

[3] Maedeh Aghaei, Mariella Dimiccoli, and Petia Radeva. With whom do i interact? detecting social interactions in egocentric photo-streams. In *2016 23rd International Conference on Pattern Recognition (ICPR)*, pages 2959–2964. IEEE, 2016. 2

[4] Peri Akiva, Jing Huang, Kevin J Liang, Rama Kovvuri, Xingyu Chen, Matt Feiszli, Kristin Dana, and Tal Hassner. Self-supervised object detection from egocentric videos. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5225–5237, 2023. 2

[5] Stefano Alletto, Giuseppe Serra, Simone Calderara, and Rita Cucchiara. Understanding social relationships in egocentric vision. *Pattern Recognition*, 48(12):4082–4096, 2015. 2

[6] Asmaa Saeed Alqahtani, Lamyia Fouad Daghestani, and Lamiaa Fattouh Ibrahim. Environments and system types of virtual reality technology in stem: A survey. *International Journal of Advanced Computer Science and Applications (IJACSA)*, 8(6), 2017. 1

[7] Shervin Ardeshir and Ali Borji. Ego2top: Matching viewers in egocentric and top-view videos. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part V 14*, pages 253–268. Springer, 2016. 2

[8] Shervin Ardeshir and Ali Borji. An exocentric look at egocentric actions and vice versa. *Computer Vision and Image Understanding*, 171:61–68, 2018. 2

[9] Shervin Ardeshir and Ali Borji. Integrating egocentric videos in top-view surveillance videos: Joint identification and temporal alignment. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 285–300, 2018. 2

[10] Shervin Ardeshir, Sandesh Sharma, and Ali Broji. Egoreid: Cross-view self-identification and human re-identification in egocentric and surveillance videos. *arXiv preprint arXiv:1612.08153*, 2016. 2

[11] Sophia Bano, Tamas Suveges, Jianguo Zhang, and Stephen J Mckenna. Multimodal egocentric analysis of focused interactions. *IEEE Access*, 6:37493–37505, 2018. 2

[12] Apurva Bedagkar-Gala and Shishir K Shah. A survey of approaches and trends in person re-identification. *Image and vision computing*, 32(4):270–286, 2014. 3

[13] Shariq Farooq Bhat, Reiner Birkl, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023. 4

[14] Hongkai Chen, Zixin Luo, Jiahui Zhang, Lei Zhou, Xuyang 656

Bai, Zeyu Hu, Chiew-Lan Tai, and Long Quan. Learning 657

to match features with seeded graph matching network. In 658

*Proceedings of the IEEE/CVF international conference on 659**computer vision*, pages 6301–6310, 2021. 2 660

[15] Hongkai Chen, Zixin Luo, Lei Zhou, Yurun Tian, Mingmin 661

Zhen, Tian Fang, David McKinnon, Yanghai Tsin, and Long 662

Quan. Aspanformer: Detector-free image matching with 663

adaptive span transformer. *European Conference on Com- 664**puter Vision (ECCV)*, 2022. 2 665

[16] Feng Cheng, Mi Luo, Huiyu Wang, Alex Dimakis, Lorenzo 666

Torresani, Gedas Bertasius, and Kristen Grauman. 4diff: 3d- 667

aware diffusion model for third-to-first viewpoint translation. 668

In *European Conference on Computer Vision*, pages 409– 669

427. Springer, 2024. 2 670

[17] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, 671

Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Da- 672

vide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, 673

and Michael Wray. Scaling egocentric vision: The epic- 674

kitchens dataset. In *European Conference on Computer Vi- 675**sion (ECCV)*, 2018. 2 676

[18] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, 677

Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide 678

Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and 679

Michael Wray. The epic-kitchens dataset: Collection, chal- 680

lenges and baselines. *IEEE Transactions on Pattern Analy- 681**sis and Machine Intelligence (TPAMI)*, 43(11):4125–4141, 682

2021. 683

[19] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, 684

Antonino Furnari, Jian Ma, Evangelos Kazakos, Davide 685

Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and 686

Michael Wray. Rescaling egocentric vision: Collection, 687

pipeline and challenges for epic-kitchens-100. *International 688**Journal of Computer Vision (IJCV)*, 130:33–55, 2022. 2 689

[20] Ana Garcia Del Molino, Cheston Tan, Joo-Hwee Lim, and 690

Ah-Hwee Tan. Summarization of egocentric videos: A com- 691

prehensive survey. *IEEE Transactions on Human-Machine 692**Systems*, 47(1):65–76, 2016. 2 693

[21] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, 694

and Li Fei-Fei. Imagenet: A large-scale hierarchical image 695

database. In *2009 IEEE conference on computer vision and 696**pattern recognition*, pages 248–255. Ieee, 2009. 7 697

[22] Ameya Dhamanaskar, Mariella Dimiccoli, Enric Corona, Al- 698

bert Pumarola, and Francesc Moreno-Noguer. Enhancing 699

egocentric 3d pose estimation with third person views. *Pat- 700**tern Recognition*, 138:109358, 2023. 3 701

[23] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, 702

Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, 703

Mostafa Dehghani, Matthias Minderer, Georg Heigold, Syl- 704

vain Gelly, et al. An image is worth 16x16 words: Trans- 705

formers for image recognition at scale. *arXiv preprint 706**arXiv:2010.11929*, 2020. 4 707

[24] Johan Edstedt, Ioannis Athanasiadis, Márten Wadenbäck, 708

and Michael Felsberg. Dkm: Dense kernelized feature 709

matching for geometry estimation. In *Proceedings of the 710**IEEE/CVF Conference on Computer Vision and Pattern 711**Recognition*, pages 17765–17775, 2023. 2 712

- 713 [25] Noel Enyedy, Joshua A Danish, and David DeLiema. Con- 771  
714 structing liminal blends in a collaborative augmented-reality 772  
715 learning environment. *International Journal of Computer- 773*  
716 *Supported Collaborative Learning*, 10(1):7–34, 2015. 1 774  
717 [26] Chenyou Fan, Jangwon Lee, Mingze Xu, Krishna Ku- 775  
718 mar Singh, Yong Jae Lee, David J Crandall, and Michael S 776  
719 Ryoo. Identifying first-person camera wearers in third- 777  
720 person videos. In *Proceedings of the IEEE Conference 778*  
721 *on Computer Vision and Pattern Recognition*, pages 5125– 779  
722 5133, 2017. 2, 3, 5 780  
723 [27] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, 781  
724 Zhicheng Yan, Jitendra Malik, and Christoph Feichten- 782  
725 hofer. Multiscale vision transformers. In *Proceedings of 783*  
726 *the IEEE/CVF international conference on computer vision*, 784  
727 pages 6824–6835, 2021. 3 785  
728 [28] Alireza Fathi, Ali Farhadi, and James M Rehg. Under- 786  
729 standing egocentric activities. In *2011 international conference 787*  
730 *on computer vision*, pages 407–414. IEEE, 2011. 2 788  
731 [29] Joyce Horn Fonteles, Clayton Cohn, Efrat Ayalon, Mengxi 789  
732 Zhou, Ashwin TS, Eduardo Davalos, Zhijian Li, Surya Ray- 790  
733 ala, Divya Mereddy, Austin Coursey, et al. Analyzing em- 791  
734 bodied learning in classroom settings: A human-in-the-loop 792  
735 ai approach for multimodal learning analytics. *Learning and 793*  
736 *Instruction*, 103:102274, 2026. 1 794  
737 [30] Kristen Grauman, Andrew Westbury, Eugene Byrne, 795  
738 Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson 796  
739 Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: 797  
740 Around the world in 3,000 hours of egocentric video. In 798  
741 *Proceedings of the IEEE/CVF conference on computer vi- 799*  
742 *sion and pattern recognition*, pages 18995–19012, 2022. 1, 800  
743 2, 6 801  
744 [31] Kristen Grauman, Andrew Westbury, Lorenzo Torresani, 802  
745 Kris Kitani, Jitendra Malik, Triantafyllos Afouras, Kumar 803  
746 Ashutosh, Vijay Baiyya, Siddhant Bansal, Bikram Boote, 804  
747 Eugene Byrne, Zach Chavis, Joya Chen, Feng Cheng, Fu- 805  
748 Jen Chu, Sean Crane, Avijit Dasgupta, Jing Dong, Maria 806  
749 Escobar, Cristhian Forigua, Abraham Gebreselasie, Sanjay 807  
750 Haresh, Jing Huang, Md Mohaiminul Islam, Suyog Jain, 808  
751 Rawal Khirodkar, Devansh Kukreja, Kevin J Liang, Jia- 809  
752 Wei Liu, Sagnik Majumder, Yongsen Mao, Miguel Martin, 810  
753 Effrosyni Mavroudi, Tushar Nagarajan, Francesco Ragusa, 811  
754 Santhosh Kumar Ramakrishnan, Luigi Seminara, Arjun So- 812  
755 mayazulu, Yale Song, Shan Su, Zihui Xue, Edward Zhang, 813  
756 Jinxu Zhang, Angela Castillo, Changan Chen, Xinzhu Fu, 814  
757 Ryosuke Furuta, Cristina Gonzalez, Prince Gupta, Jiabo 815  
758 Hu, Yifei Huang, Yiming Huang, Weslie Khoo, Anush Ku- 816  
759 mar, Robert Kuo, Sach Lakhavani, Miao Liu, Mi Luo, 817  
760 Zhengyi Luo, Brigid Meredith, Austin Miller, Oluwatumin- 818  
761 inu Oguntola, Xiaqing Pan, Penny Peng, Shraman Praman- 819  
762 ick, Merrey Ramazanov, Fiona Ryan, Wei Shan, Kiran So- 820  
763 masundaram, Chenan Song, Audrey Southerland, Masatoshi 821  
764 Tateno, Huiyu Wang, Yuchen Wang, Takuma Yagi, Mingfei 822  
765 Yan, Xitong Yang, Zecheng Yu, Shengxin Cindy Zha, Chen 823  
766 Zhao, Ziwei Zhao, Zhifan Zhu, Jeff Zhuo, Pablo Arbe- 824  
767 laez, Gedas Bertasius, Dima Damen, Jakob Engel, Gio- 825  
768 vanni Maria Farinella, Antonino Furnari, Bernard Ghanem, 826  
769 Judy Hoffman, C.V. Jawahar, Richard Newcombe, Hyun Soo 827  
770 Park, James M. Rehg, Yoichi Sato, Manolis Savva, Jianbo 828  
Shi, Mike Zheng Shou, and Michael Wray. Ego-exo4d: 771  
Understanding skilled human activity from first- and third- 772  
person perspectives. In *Proceedings of the IEEE/CVF 773*  
*Conference on Computer Vision and Pattern Recognition 774*  
*(CVPR)*, pages 19383–19400, 2024. 2 775  
[32] Kristen Grauman, Andrew Westbury, Lorenzo Torresani, 776  
Kris Kitani, Jitendra Malik, Triantafyllos Afouras, Kumar 777  
Ashutosh, Vijay Baiyya, Siddhant Bansal, Bikram Boote, 778  
et al. Ego-exo4d: Understanding skilled human activity 779  
from first-and third-person perspectives. In *Proceedings of 780*  
*the IEEE/CVF Conference on Computer Vision and Pattern 781*  
*Recognition*, pages 19383–19400, 2024. 1, 6 782  
[33] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Gir- 783  
shick. Mask r-cnn. In *Proceedings of the IEEE international 784*  
*conference on computer vision*, pages 2961–2969, 2017. 4 785  
[34] Lingxiao He, Xingyu Liao, Wu Liu, Xinchun Liu, Peng 786  
Cheng, and Tao Mei. Fastred: A pytorch toolbox for general 787  
instance re-identification. In *Proceedings of the 31st ACM 788*  
*International Conference on Multimedia*, pages 9664–9667, 789  
2023. 4 790  
[35] Hsuan-I Ho, Wei-Chen Chiu, and Yu-Chiang Frank Wang. 791  
Summarizing first-person videos from third persons’ points 792  
of view. In *Proceedings of the European conference on com- 793*  
*puter vision (ECCV)*, pages 70–85, 2018. 2 794  
[36] Miao Hu, Xianzhuo Luo, Jiawen Chen, Young Choon Lee, 795  
Yipeng Zhou, and Di Wu. Virtual reality: A survey of en- 796  
abling technologies and its applications in iot. *Journal of 797*  
*Network and Computer Applications*, 178:102970, 2021. 1 798  
[37] Yifei Huang, Minjie Cai, Zhenqiang Li, and Yoichi Sato. 799  
Predicting gaze in egocentric video by learning task- 800  
dependent attention transition. In *Proceedings of the Eu- 801*  
*ropean conference on computer vision (ECCV)*, pages 754– 802  
769, 2018. 2 803  
[38] Yifei Huang, Minjie Cai, Zhenqiang Li, Feng Lu, and Yoichi 804  
Sato. Mutual context network for jointly estimating egocen- 805  
tric gaze and action. *IEEE Transactions on Image Process- 806*  
*ing*, 29:7795–7806, 2020. 2 807  
[39] Yifei Huang, Guo Chen, Jilan Xu, Mingfang Zhang, Li- 808  
jin Yang, Baoqi Pei, Hongjie Zhang, Lu Dong, Yali Wang, 809  
Limin Wang, et al. Egoexolearn: A dataset for bridging asyn- 810  
chronous ego-and exo-centric view of procedural activities 811  
in real world. In *Proceedings of the IEEE/CVF Conference 812*  
*on Computer Vision and Pattern Recognition*, pages 22072– 813  
22086, 2024. 1 814  
[40] Lasse Jensen and Flemming Konradsen. A review of the 815  
use of virtual reality head-mounted displays in education 816  
and training. *Education and information technologies*, 23 817  
(4):1515–1529, 2018. 1 818  
[41] Baoxiong Jia, Ting Lei, Song-Chun Zhu, and Siyuan Huang. 819  
Egotaskqa: Understanding human tasks in egocentric videos. 820  
*Advances in Neural Information Processing Systems*, 35: 821  
3343–3360, 2022. 2 822  
[42] Simar Kareer, Dhruv Patel, Ryan Punamiya, Pranay Mathur, 823  
Shuo Cheng, Chen Wang, Judy Hoffman, and Danfei Xu. 824  
Egomimic: Scaling imitation learning via egocentric video. 825  
*arXiv preprint arXiv:2410.24221*, 2024. 826  
[43] Haresh Karnan, Garrett Warnell, Xuesu Xiao, and Peter 827  
Stone. Voila: Visual-observation-only imitation learning 828

- 829 for autonomous navigation. In *2022 International Confer-*  
830 *ence on Robotics and Automation (ICRA)*, pages 2497–2503.  
831 IEEE, 2022. 2
- 832 [44] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang,  
833 Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola,  
834 Tim Green, Trevor Back, Paul Natsev, et al. The kinetics hu-  
835 man action video dataset. *arXiv preprint arXiv:1705.06950*,  
836 2017. 7
- 837 [45] Rahima Khanam and Muhammad Hussain. Yolov11: An  
838 overview of the key architectural enhancements. *arXiv*  
839 *preprint arXiv:2410.17725*, 2024. 6, 15
- 840 [46] James C Lester, Hiller A Spires, John L Nietfeld, James  
841 Minogue, Bradford W Mott, and Eleni V Lobene. Design-  
842 ing game-based learning environments for elementary sci-  
843 ence education: A narrative-centered learning perspective.  
844 *Information Sciences*, 264:4–18, 2014. 1
- 845 [47] Yin Li, Alireza Fathi, and James M Rehg. Learning to predict  
846 gaze in egocentric video. In *Proceedings of the IEEE inter-*  
847 *national conference on computer vision*, pages 3216–3223,  
848 2013. 2
- 849 [48] Yanghao Li, Tushar Nagarajan, Bo Xiong, and Kristen Grau-  
850 man. Ego-exo: Transferring visual representations from  
851 third-person to first-person videos. In *Proceedings of the*  
852 *IEEE/CVF Conference on Computer Vision and Pattern*  
853 *Recognition*, pages 6943–6953, 2021. 2
- 854 [49] Kevin Qinghong Lin, Jinpeng Wang, Mattia Soldan, Michael  
855 Wray, Rui Yan, Eric Z Xu, Difei Gao, Rong-Cheng Tu, Wen-  
856 zhe Zhao, Weijie Kong, et al. Egocentric video-language  
857 pretraining. *Advances in Neural Information Processing Sys-*  
858 *tems*, 35:7575–7586, 2022. 2
- 859 [50] Yunze Liu, Yun Liu, Che Jiang, Kangbo Lyu, Weikang Wan,  
860 Hao Shen, Boqiang Liang, Zhoujie Fu, He Wang, and Li Yi.  
861 Hoi4d: A 4d egocentric dataset for category-level human-  
862 object interaction. In *Proceedings of the IEEE/CVF Con-*  
863 *ference on Computer Vision and Pattern Recognition*, pages  
864 21013–21022, 2022. 2
- 865 [51] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang,  
866 Stephen Lin, and Han Hu. Video swin transformer. In *Pro-*  
867 *ceedings of the IEEE/CVF conference on computer vision*  
868 *and pattern recognition*, pages 3202–3211, 2022. 13
- 869 [52] Ilya Loshchilov and Frank Hutter. Decoupled weight decay  
870 regularization. *arXiv preprint arXiv:1711.05101*, 2017. 7
- 871 [53] Minlong Lu, Danping Liao, and Ze-Nian Li. Learning spa-  
872 tiotemporal attention for egocentric action recognition. In  
873 *Proceedings of the IEEE/CVF International Conference on*  
874 *Computer Vision Workshops*, pages 0–0, 2019. 2
- 875 [54] Mi Luo, Zihui Xue, Alex Dimakis, and Kristen Grauman.  
876 Put myself in your shoes: Lifting the egocentric perspective  
877 from exocentric videos. In *European Conference on Com-*  
878 *puter Vision*, pages 407–425. Springer, 2024. 2, 6
- 879 [55] Minghuang Ma, Haoqi Fan, and Kris M Kitani. Going deeper  
880 into first-person activity recognition. In *Proceedings of the*  
881 *IEEE Conference on Computer Vision and Pattern Recogni-*  
882 *tion*, pages 1894–1903, 2016. 2
- 883 [56] TorchVision maintainers and contributors. Torchvision: Py-  
884 torch’s computer vision library. [https://github.com/](https://github.com/pytorch/vision)  
885 [pytorch/vision](https://github.com/pytorch/vision), 2016. 13
- [57] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra 886  
Malik. Egoschema: A diagnostic benchmark for very long- 887  
form video language understanding. *Advances in Neural In-* 888  
*formation Processing Systems*, 36:46212–46244, 2023. 2 889
- [58] Kyle Min and Jason J Corso. Integrating human gaze into at- 890  
tention for egocentric activity recognition. In *Proceedings of* 891  
*the IEEE/CVF Winter Conference on Applications of Com-* 892  
*puter Vision*, pages 1069–1078, 2021. 2 893
- [59] Keisuke Ogaki, Kris M Kitani, Yusuke Sugano, and Yoichi 894  
Sato. Coupling eye-motion and ego-motion features for first- 895  
person activity recognition. In *2012 IEEE Computer Soci-* 896  
*ety Conference on Computer Vision and Pattern Recognition* 897  
*Workshops*, pages 1–7. IEEE, 2012. 3 898
- [60] A Emin Orhan, Wentao Wang, Alex N Wang, Mengye Ren, 899  
and Brenden M Lake. Self-supervised learning of video 900  
representations from a child’s perspective. *arXiv preprint* 901  
*arXiv:2402.00300*, 2024. 2 902
- [61] Rémi Pautrat, Iago Suárez, Yifan Yu, Marc Pollefeys, and 903  
Viktor Larsson. Gluestick: Robust image matching by 904  
sticking points and lines together. In *Proceedings of the* 905  
*IEEE/CVF International Conference on Computer Vision*, 906  
pages 9706–9716, 2023. 2, 3 907
- [62] Chiara Plizzari, Gabriele Goletto, Antonino Furnari, Sid- 908  
dhanth Bansal, Francesco Ragusa, Giovanni Maria Farinella, 909  
Dima Damen, and Tatiana Tommasi. An outlook into the fu- 910  
ture of egocentric vision. *International Journal of Computer* 911  
*Vision*, 132(11):4880–4936, 2024. 2 912
- [63] Francesco Ragusa, Antonino Furnari, Salvatore Livatino, 913  
and Giovanni Maria Farinella. The meccano dataset: 914  
Understanding human-object interactions from egocentric 915  
videos in an industrial-like domain. In *Proceedings of the* 916  
*IEEE/CVF Winter Conference on Applications of Computer* 917  
*Vision*, pages 1569–1578, 2021. 2 918
- [64] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, 919  
and Andrew Rabinovich. Superglue: Learning feature 920  
matching with graph neural networks. In *Proceedings of* 921  
*the IEEE/CVF conference on computer vision and pattern* 922  
*recognition*, pages 4938–4947, 2020. 2 923
- [65] Fadime Sener, Dibiyadip Chatterjee, Daniel Shelepov, Kun 924  
He, Dipika Singhanian, Robert Wang, and Angela Yao. As- 925  
sembly101: A large-scale multi-view video dataset for un- 926  
derstanding procedural activities. In *Proceedings of the* 927  
*IEEE/CVF Conference on Computer Vision and Pattern* 928  
*Recognition*, pages 21096–21106, 2022. 2 929
- [66] Yan Shi, Jun-Xiong Cai, Yoli Shavit, Tai-Jiang Mu, Wensen 930  
Feng, and Kai Zhang. Clustergnn: Cluster-based coarse-to- 931  
fine graph neural network for efficient feature matching. In 932  
*Proceedings of the IEEE/CVF conference on computer vi-* 933  
*sion and pattern recognition*, pages 12517–12526, 2022. 2 934
- [67] Gunnar A Sigurdsson, Abhinav Gupta, Cordelia Schmid, Ali 935  
Farhadi, and Karteek Alahari. Actor and observer: Joint 936  
modeling of first and third-person videos. In *proceedings of* 937  
*the IEEE conference on computer vision and pattern recog-* 938  
*nition*, pages 7396–7404, 2018. 2 939
- [68] Gunnar A. Sigurdsson, Abhinav Gupta, Cordelia Schmid, 940  
Ali Farhadi, and Karteek Alahari. Charades-ego: A large- 941  
scale dataset of paired third and first person videos, 2018. 942  
2 943

- 944 [69] Nithin Sivakumaran, Chia-Yu Yang, Abhay Zala, Shoubin  
945 Yu, Daeun Hong, Xiaotian Zou, Elias Stengel-Eskin, Dan  
946 Carpenter, Wookhee Min, Cindy Hmelo-Silver, et al. A mul-  
947 timodal classroom video question-answering framework for  
948 automated understanding of collaborative learning. In *Pro-  
949 ceedings of the 27th International Conference on Multimodal  
950 Interaction*, pages 516–525, 2025. 1
- 951 [70] Sibong Song, Vijay Chandrasekhar, Bappaditya Mandal,  
952 Liyuan Li, Joo-Hwee Lim, Giduthuri Sateesh Babu, Phyto  
953 Phyto San, and Ngai-Man Cheung. Multimodal multi-stream  
954 deep learning for egocentric activity recognition. In *Proceed-  
955 ings of the IEEE conference on computer vision and pattern  
956 recognition workshops*, pages 24–31, 2016. 2
- 957 [71] Bilge Soran, Ali Farhadi, and Linda Shapiro. Action recog-  
958 nition in the presence of one egocentric and multiple static  
959 cameras. In *Asian Conference on Computer Vision*, pages  
960 178–193. Springer, 2014. 2
- 961 [72] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and  
962 Xiaowei Zhou. Loftr: Detector-free local feature matching  
963 with transformers. In *Proceedings of the IEEE/CVF con-  
964 ference on computer vision and pattern recognition*, pages  
965 8922–8931, 2021. 2
- 966 [73] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field  
967 transforms for optical flow. In *Computer Vision—ECCV  
968 2020: 16th European Conference, Glasgow, UK, August 23–  
969 28, 2020, Proceedings, Part II 16*, pages 402–419. Springer,  
970 2020. 4
- 971 [74] Bugra Tekin, Federica Bogo, and Marc Pollefeys. H+ o: Uni-  
972 fied egocentric recognition of 3d hand-object poses and in-  
973 teractions. In *Proceedings of the IEEE/CVF conference on  
974 computer vision and pattern recognition*, pages 4511–4520,  
975 2019. 2
- 976 [75] Sanket Kumar Thakur, Cigdem Beyan, Pietro Morerio, and  
977 Alessio Del Bue. Predicting gaze from egocentric social in-  
978 teraction videos and imu data. In *Proceedings of the 2021  
979 International Conference on Multimodal Interaction*, pages  
980 717–722, 2021. 2
- 981 [76] Prune Truong, Martin Danelljan, Radu Timofte, and Luc  
982 Van Gool. Pdc-net+: Enhanced probabilistic dense corre-  
983 spondence network. *IEEE Transactions on Pattern Analysis  
984 and Machine Intelligence*, 45(8):10247–10266, 2023. 2
- 985 [77] Nesse Van der Meer, Vivian van der Werf, Willem-Paul  
986 Brinkman, and Marcus Specht. Virtual reality and collabor-  
987 ative learning: A systematic literature review. *Frontiers in  
988 Virtual Reality*, 4:1159905, 2023. 1
- 989 [78] Xiaohan Wang, Linchao Zhu, Heng Wang, and Yi Yang. In-  
990 teractive prototype learning for egocentric action recogni-  
991 tion. In *Proceedings of the IEEE/CVF International Con-  
992 ference on Computer Vision*, pages 8168–8177, 2021. 2
- 993 [79] Mingze Xu, Chenyou Fan, Yuchen Wang, Michael S Ryoo,  
994 and David J Crandall. Joint person segmentation and iden-  
995 tification in synchronized first-and third-person videos. In  
996 *Proceedings of the European Conference on Computer Vi-  
997 sion (ECCV)*, pages 637–652, 2018. 2, 3, 5, 7
- 998 [80] Yue Xu, Yong-Lu Li, Zheming Huang, Michael Xu Liu,  
999 Cewu Lu, Yu-Wing Tai, and Chi-Keung Tang. Egopca: A  
1000 new framework for egocentric hand-object interaction un-  
derstanding. In *Proceedings of the IEEE/CVF International  
Conference on Computer Vision*, pages 5273–5284, 2023. 2
- [81] Liang Yang, Hao Jiang, Jizhong Xiao, and Zhouyuan Huo. Ego-downward and ambient video based person location association. *arXiv preprint arXiv:1812.00477*, 2018. 2
- [82] Mang Ye, Jianbing Shen, Gaojie Lin, Tao Xiang, Ling Shao, and Steven CH Hoi. Deep learning for person re-identification: A survey and outlook. *IEEE transactions on pattern analysis and machine intelligence*, 44(6):2872–2893, 2021. 3
- [83] Ziwei Zhao, Yuchen Wang, and Chuhua Wang. Fusing personal and environmental cues for identification and segmentation of first-person camera wearers in third-person views. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16477–16487, 2024. 2, 3, 4, 6, 7, 8
- [84] Liang Zheng, Yi Yang, and Alexander G Hauptmann. Person re-identification: Past, present and future. *arXiv preprint arXiv:1610.02984*, 2016. 3

# Sequence-Based Identification of First-Person Camera Wearers in Third-Person Views

## Supplementary Material

### 1020 6. Additional Visualization of Models

1021 We provide detailed visualization of our *motion matching*  
1022 and *appearance matching* methods in Fig. 6 and Fig. 7.

### 1023 7. Additional Ablation Studies

1024 We evaluated the choice of backbones and the effective-  
1025 ness of the sliding window approach. For this analysis,  
1026 we selected two transformer-based video methods, *MViT*  
1027 and *Swin3D* [51], using their default weights from torchvi-  
1028 sion [56], both pretrained on **Kinetics-400**. We compared  
1029 two input strategies: (1) direct prediction using the full 30-  
1030 frame sequence (frames 1 to 30), and (2) cumulative predic-  
1031 tions from three overlapping 16-frame sequences (frames  
1032 1–16, 7–23, and 14–30). This comparison was conducted  
1033 using *Swin3D*, as the pretrained *MViT* weights do not sup-  
1034 port 30-frame inputs. The results of this ablation study are  
1035 presented in Tab. 6.

| Dataset                   | TF2025       |              |               |
|---------------------------|--------------|--------------|---------------|
| Methodology               | Seen         | Unseen       | Cross-dataset |
| Swin3d w/o sliding window | 74.7%        | 63.2%        | 57.7%         |
| Swin3d w/ sliding window  | 82.2%        | 72.9%        | 64.5%         |
| MViT w/ sliding window    | <b>88.8%</b> | <b>80.5%</b> | <b>67.7%</b>  |

Table 6. Ablation study on backbones and sliding window.

1036 We validated our choice of optical flow computation.  
1037 Specifically, we compared our approach of calculating the  
1038 hypotenuse of the optical flow components along the x and  
1039 y axes against using optical flow from a single direction, as  
1040 shown in Tab. 7.

| Dataset               | TF2025       |              |               |
|-----------------------|--------------|--------------|---------------|
| Methodology           | Seen         | Unseen       | Cross-dataset |
| X axis only           | 83.1%        | 72.8%        | 63.3%         |
| Y axis only           | 74.7%        | 70.9%        | 63.5%         |
| Hypotenuse of x and y | <b>88.8%</b> | <b>80.5%</b> | <b>67.7%</b>  |

Table 7. Ablation study on optical flow directions.

1041 We also performed a sensitivity study on the hyperpa-  
1042 rameter  $\lambda_{\text{Re-ID}}$  in CBAF, which balances between motion  
1043 and appearance matching. As shown in Tab. 8, optimal val-  
1044 ues are 1.75 for the seen split, 1.75 for the unseen split, and

3 for the cross-dataset split, with minimal differences from  
non-optimal values. This demonstrates the adaptability of  
our method to different scenarios. For example, the higher  
value in the cross-dataset scenario reflects the greater reli-  
ability of appearance matching compared to motion match-  
ing due to the difference in camera models. In addition,  
CBAF exhibits robustness to changes in this hyperparam-  
eter, with performance varying only by  $\sim 1\%$  in the seen  
split for values ranging from 1.25 to 2.5.

| Dataset                  | TF2025       |              |               |
|--------------------------|--------------|--------------|---------------|
| $\lambda_{\text{Re-ID}}$ | Seen         | Unseen       | Cross-dataset |
| 0(motion)                | 88.8%        | 80.5%        | 67.7%         |
| 0.25                     | 88.9%        | 80.6%        | 67.7%         |
| 0.5                      | 89.6%        | 81.5%        | 68.0%         |
| 0.75                     | 90.5%        | 82.5%        | 69.0%         |
| 1                        | 91.4%        | 82.6%        | 71.4%         |
| 1.25                     | 93.0%        | 83.6%        | 73.0%         |
| 1.5                      | 93.5%        | 84.3%        | 73.6%         |
| 1.75                     | <b>93.8%</b> | <b>84.4%</b> | 73.7%         |
| 2                        | 93.3%        | 83.9%        | 74.7%         |
| 2.25                     | 92.9%        | 83.6%        | 74.8%         |
| 2.5                      | 93.3%        | 83.2%        | 75.4%         |
| 3                        | 92.3%        | 82.4%        | <b>75.7%</b>  |
| 3.5                      | 91.9%        | 81.8%        | 75.5%         |
| 4                        | 91.9%        | 81.4%        | 75.2%         |

Table 8. Sensitivity analysis for the hyperparameter of CBAF.

To further demonstrate the robustness of our model to  
the selection of  $\lambda$ , we conducted validation by splitting each  
test set into three subsets. In each run, we used one subset  
to select  $\lambda$ , and the remaining two for testing, averaging the  
results across all three combinations. As shown in Tab. 9,  
this results in only a negligible reduction ( $\leq 0.5\%$ ) while  
we significantly outperform state-of-the-art ( $\sim 20\%$ ).

| Selection of $\lambda$ | Seen  | Unseen | Cross |
|------------------------|-------|--------|-------|
| w/ validation          | 93.7% | 84.3%  | 75.2% |
| w/o validation         | 93.8% | 84.4%  | 75.7% |

Table 9. Validation for the selection of  $\lambda$

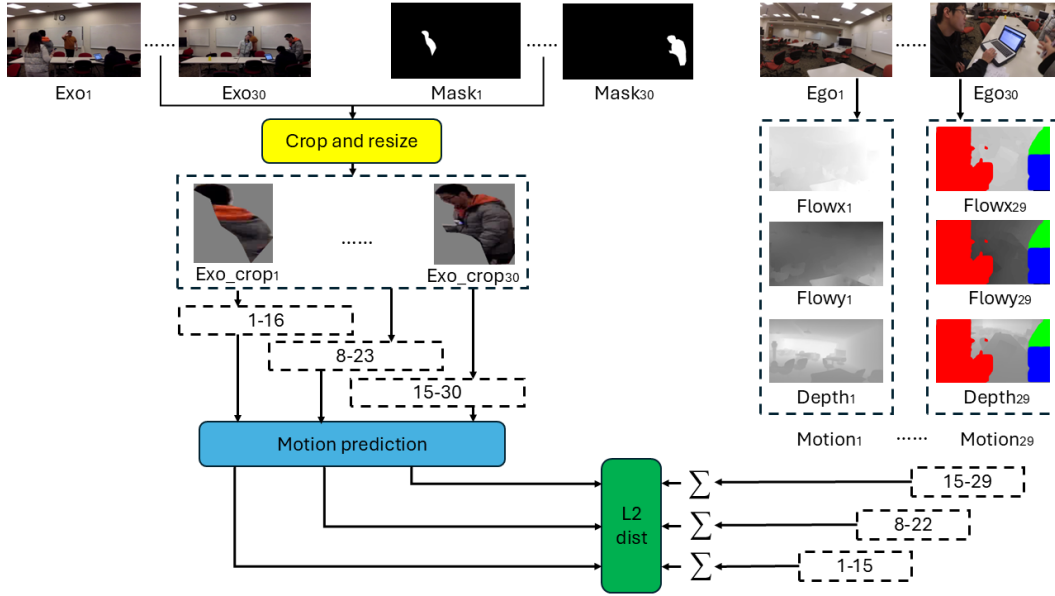


Figure 6. **Overview of the motion matching model.** We employ a sequence-based backbone (*MViT*) on the cropped third-person candidate and align it with the motion estimated from the first-person view. For visualization purposes, we rescale the optical flows such that black represents negative flow, white represents positive flow, and colored masks indicate foreground pixels excluded from motion estimation due to person detection. There are only 29 frames for first-person because optical flow requires 2 frames to calculate.

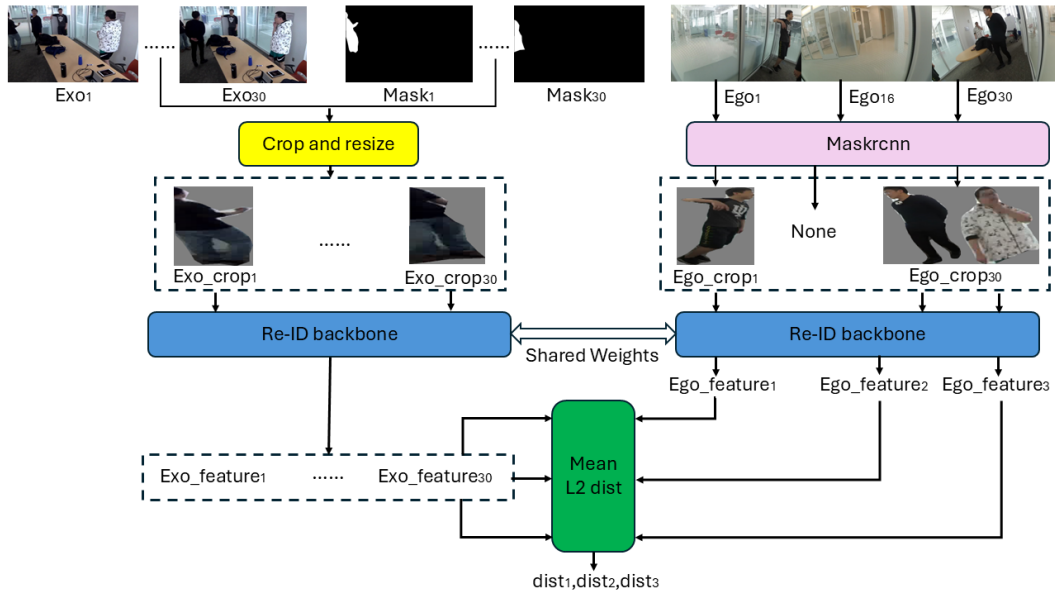


Figure 7. **Overview of appearance matching.** We apply *Mask R-CNN* on three selected frames from the first-person view and pass each of the  $M$  detected individuals, along with the third-person candidate, through the re-ID backbone. We then compute the average L2 distance between each detected individual and the third-person candidate across all 30 frames. This process results in  $M$  distance scores.

## 1061 8. Examples of Fusing Algorithm

1062 To demonstrate our fusion method (**CBAF**), we provide  
1063 step-by-step examples in Figs. 8–11. For clarity in visu-  
1064 alization, we set  $\lambda_{\text{Re-ID}} = 1$  and  $\alpha_{\text{mask}} = 1$  for each mask,  
1065 enabling us to omit these terms during confidence calcula-

tions.

## 1066 9. Annotation Details and Data Samples

1067 To create **Ego4D-TF**, a subset of **Ego4D** annotated to sup-  
1068 port our task, we used the following annotation pipeline: 1069

| Candidates:                 | 1    | 2    | 3    | 4   |
|-----------------------------|------|------|------|-----|
| Motion scores:              | 12.4 | 22.3 | 1.2  | 5.8 |
| Re-ID scores:<br>(Frame 1)  |      |      | None |     |
| Re-ID scores:<br>(Frame 16) |      |      | None |     |
| Re-ID scores:<br>(Frame 30) |      |      | None |     |

Figure 8. **Example 1:** No individuals are detected in the first-person view, leading to the prediction of candidate 3 (with the smallest motion score).

| Candidates:                 | 1    | 2    | 3    | 4    |                    |
|-----------------------------|------|------|------|------|--------------------|
| Motion scores:              | 3.7  | 16.7 | 21.4 | 15.8 | Conf.=15.8/3.7=4.3 |
| Re-ID scores:<br>(Frame 1)  | 12.4 | 11.6 | 5.8  | 14.7 | Conf.=11.6/5.8=2   |
| Re-ID scores:<br>(Frame 16) | 14.3 | 7.4  | 13.8 | 16.1 | Conf.=13.8/7.4=1.9 |
| Re-ID scores:<br>(Frame 30) | 11.9 | 9.2  | 12.3 | 17.5 | Conf.=11.9/9.2=1.3 |

Figure 9. **Example 2:** Detected several individuals in first-person view, but motion has the highest confidence, predicting candidate 1.

1070 First, we identified group activity videos in Ego4D contain-  
 1071 ing at least three people, with at least two videos already  
 1072 synchronized. From these, we selected one video with a  
 1073 broad enough view to frequently capture at least two other  
 1074 people. We then sampled one out of every six frames to  
 1075 align the frame rate with **IUShareView** and **TF2023**, and  
 1076 applied YOLO-v11 [45] to 30-frame sequences to detect  
 1077 and track individuals. For sequences where more than two  
 1078 people were detected, we manually verified the masks of  
 1079 the first and last frames to ensure quality. Sequences were  
 1080 rejected if the mask quality was low, or the masks corre-  
 1081 sponded to different individuals in the first and last frames  
 1082 (lost tracking). The likelihood of losing tracking during a  
 1083 sequence was minimal, as all five groups of videos we used  
 1084 involved board games, where participants rarely made large  
 1085 movements. We show two examples of accepted annota-  
 1086 tions in Fig. 12 and two examples of rejected annotations in  
 1087 Fig. 13.

1088 We also present sample data from the three source  
 1089 datasets used in **TF2025**: **TF2023**, **IUShareView**, and  
 1090 **Ego4D-TF**, illustrating the similarities and differences in  
 1091 their settings, as visualized in Fig. 14.

| Candidates:    | 1    | 2    | 3    |                    |
|----------------|------|------|------|--------------------|
| Motion scores: | 2.4  | 18.7 | 1.9  | Conf.=2.4/1.9=1.3  |
| Re-ID scores:  | 11.3 | 11.7 | 3.6  | Conf.=11.3/3.6=3.1 |
| (Frame 1)      | 17.1 | 6.7  | 10.7 | Conf.=10.7/6.7=1.6 |
| Re-ID scores:  | 16.5 | 4.2  | 19.4 | Conf.=16.5/4.2=3.9 |
| (Frame 16)     | 14.0 | 13.3 | 9.8  | Conf.=13.3/9.8=1.6 |
| Re-ID scores:  | 11.4 | 9.9  | 13.1 | Conf.=11.4/9.9=1.2 |
| (Frame 30)     |      |      |      |                    |

(a) Calculate confidence score for each source.

| Candidates:    | 1    | 3    |                            |
|----------------|------|------|----------------------------|
| Motion scores: | 2.4  | 1.9  | Conf.=2.4/1.9=1.3          |
| Re-ID scores:  | 11.3 | 3.6  | Conf.=11.3/3.6=3.1         |
| (Frame 1)      | 17.1 | 10.7 | <b>Conf.=17.1/10.7=1.6</b> |
| Re-ID scores:  | 14.0 | 9.8  | <b>Conf.=14.0/9.8=1.4</b>  |
| (Frame 16)     | 11.4 | 13.1 | <b>Conf.=13.1/11.4=1.1</b> |
| (Frame 30)     |      |      |                            |

(c) Recalculate confidence scores. The updated scores are highlighted in **bold**.

| Candidates:    | 1    | 3    |
|----------------|------|------|
| Motion scores: | 2.4  | 1.9  |
| Re-ID scores:  | 11.3 | 3.6  |
| (Frame 1)      | 17.1 | 10.7 |
| Re-ID scores:  | 14.0 | 9.8  |
| (Frame 16)     | 11.4 | 13.1 |
| (Frame 30)     |      |      |

(b) Remove the appearance source with highest confidence, and the candidate with the smallest score.

| Candidates:    | 1    | 3 |
|----------------|------|---|
| Motion scores: | 2.4  |   |
| Re-ID scores:  | 17.1 |   |
| (Frame 1)      | 14.0 |   |
| Re-ID scores:  | 11.4 |   |
| (Frame 16)     |      |   |
| (Frame 30)     |      |   |

(d) Only 1 candidate 1 remaining, predicting candidate 1.

Figure 10. **Example 3:** Two candidates (1 and 3) have similar motion scores, a scenario that can occur when both candidates remain idle during the sequence. Our method resolves this by leveraging appearance information to eliminate candidates.

| Candidates:    | 1    | 2    | 3    | 4    |                    |
|----------------|------|------|------|------|--------------------|
| Motion scores: | 12.3 | 1.5  | 22.7 | 2.4  | Conf.=2.4/1.5=1.6  |
| Re-ID scores:  |      | None |      |      |                    |
| (Frame 1)      |      |      |      |      |                    |
| Re-ID scores:  | 2.1  | 11.4 | 13.2 | 12.7 | Conf.=11.4/2.1=5.4 |
| (Frame 16)     | 19.6 | 6.8  | 17.5 | 16.3 | Conf.=16.3/6.8=2.4 |
| Re-ID scores:  | 13.1 | 12.7 | 7.3  | 15.0 | Conf.=12.7/7.3=1.7 |
| (Frame 30)     | 1.5  | 16.7 | 20.1 | 14.2 | Conf.=14.2/1.5=9.5 |

(a) Calculate confidence score for each source.

| Candidates:    | 2    | 3    | 4    |                            |
|----------------|------|------|------|----------------------------|
| Motion scores: | 1.5  | 22.7 | 2.4  | Conf.=2.4/1.5=1.6          |
| Re-ID scores:  | 11.4 | 13.2 | 12.7 | <b>Conf.=12.7/11.4=1.1</b> |
| (Frame 1)      | 6.8  | 17.5 | 16.3 | Conf.=16.3/6.8=2.4         |
| Re-ID scores:  | 12.7 | 7.3  | 15.0 | Conf.=12.7/7.3=1.7         |
| (Frame 16)     |      |      |      |                            |
| (Frame 30)     |      |      |      |                            |

(c) Recalculate confidence scores. The updated scores are highlighted in **bold**.

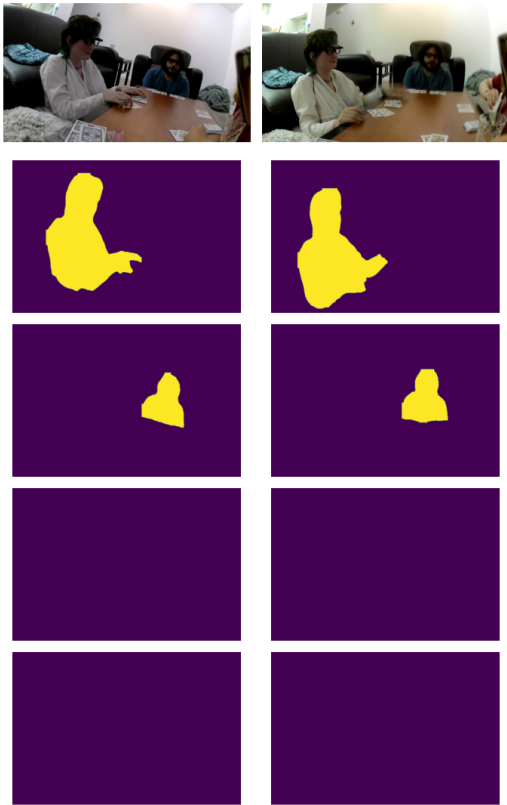
| Candidates:    | 2    | 3    | 4    |
|----------------|------|------|------|
| Motion scores: | 1.5  | 22.7 | 2.4  |
| Re-ID scores:  | 11.4 | 13.2 | 12.7 |
| (Frame 1)      | 6.8  | 17.5 | 16.3 |
| Re-ID scores:  | 12.7 | 7.3  | 15.0 |
| (Frame 16)     |      |      |      |
| (Frame 30)     |      |      |      |

(b) Remove the appearance source with highest confidence, and the candidate with the smallest score.

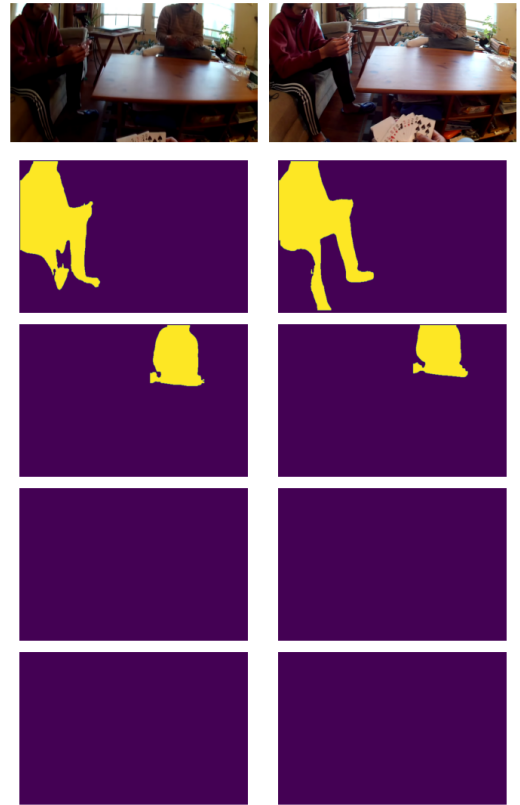
| Candidates:    | 3    | 4    |                            |
|----------------|------|------|----------------------------|
| Motion scores: | 22.7 | 2.4  | <b>Conf.=22.7/2.4=9.5</b>  |
| Re-ID scores:  | 17.5 | 16.3 | <b>Conf.=17.5/16.3=1.1</b> |
| (Frame 1)      | 7.3  | 15.0 | <b>Conf.=15.0/7.3=2.1</b>  |
| (Frame 16)     |      |      |                            |
| (Frame 30)     |      |      |                            |

(d) Motion score has the highest confidence, predicting candidate 4.

Figure 11. **Example 4:** A more complex scenario with two candidates (2 and 4) with similar motion scores. After eliminating two candidates, the updated motion scores exhibit the highest confidence, leading to the final prediction. This example demonstrates how our method adaptively selects which source to trust, effectively integrating information from both sources.



(a) Accepted annotation.

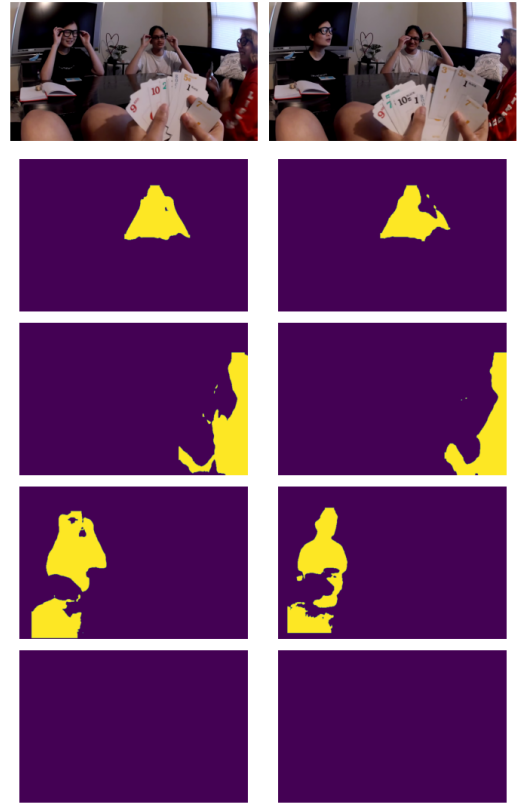


(b) Accepted annotation.)

Figure 12. Examples of accepted annotations for **Ego4D-TF**.



(a) Rejected due to tracking lost.



(b) Rejected due to low quality masks. (The individual to the left was merged with the legs of the camera wearer.)

Figure 13. Examples of rejected annotations for **Ego4D-TF**.

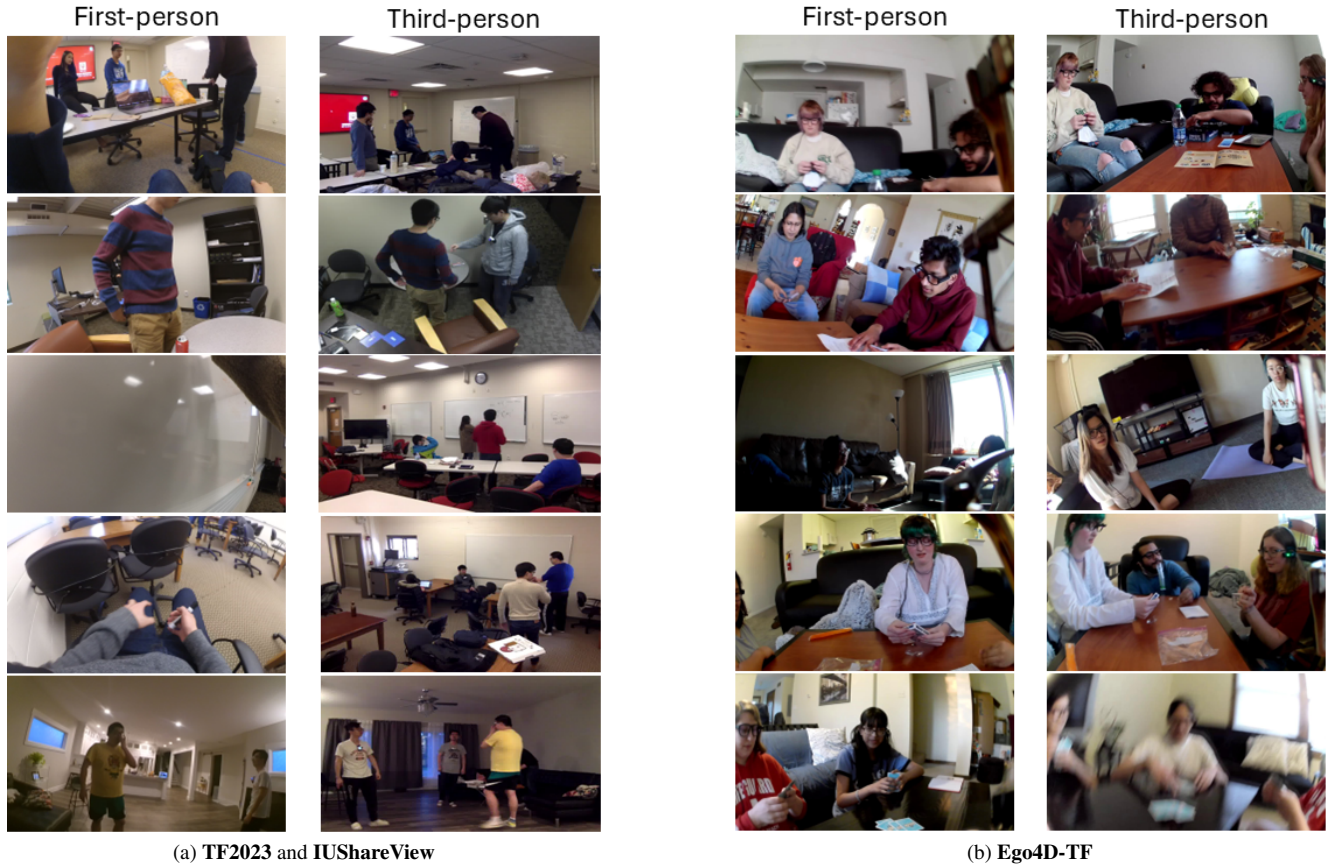


Figure 14. **Data Comparison** We present sample data from **TF2023**, **IUShareView** (which share actors and settings), and our labeled subset of **Ego4D**, referred to as **Ego4D-TF**. A key difference between these datasets is that in **Ego4D-TF**, we treat first-person views as "pseudo" third-person views. Additionally, **Ego4D-TF** utilizes head-mounted cameras, whereas **TF2023** and **IUShareView** use chest-mounted cameras.