

Human Motion Estimation with Everyday Wearables

Siqi Zhu^{1*} Yixuan Li^{1*} Junfu Li^{1,2*} Qi Wu^{1,2*}
Zan Wang^{1*} Haozhe Ma³ Wei Liang^{1,2}✉

* Equal contributors ¹ Beijing Institute of Technology

² Yangtze Delta Region Academy of Beijing Institute of Technology, Jiaxing ³ Shenzhen MSU-BIT University

<https://pie-lab.cn/EveryWear/>



Figure 1. **EveryWear** is a novel human motion capture approach based on a series of lightweight everyday wearables: a smartphone, smartwatch, earbuds, and smart glasses equipped with one forward-facing and two downward-facing cameras.

Abstract

While on-body device-based human motion estimation is crucial for applications such as XR interaction, existing methods often suffer from poor wearability, expensive hardware, and cumbersome calibration, which hinder their adoption in daily life. To address these challenges, we present **EveryWear**, a lightweight and practical human motion capture approach based entirely on everyday wearables: a smartphone, smartwatch, earbuds, and smart glasses equipped with one forward-facing and two downward-facing cameras, requiring no explicit calibration before use. We introduce **Ego-Elec**, a 9-hour real-world dataset covering 56 daily activities across 17 diverse indoor and outdoor environments, with ground-truth 3D annotations provided by the motion capture (MoCap), to facilitate robust research and benchmarking in this direction. Our

approach employs a multimodal teacher-student framework that integrates visual cues from egocentric cameras with inertial signals from consumer devices. By training directly on real-world data rather than synthetic data, our model effectively eliminates the sim-to-real gap that constrains prior work. Experiments demonstrate that our method outperforms baseline models, validating its effectiveness for practical full-body motion estimation.

1. Introduction

Human motion estimation is crucial for natural interaction and realistic avatar animation in XR applications [15, 18, 40, 42, 44, 49], as well as expressive humanoid whole-body control [8, 17, 22, 48]. Prior work on estimating human motion follows two main directions. Camera-based methods [1, 2, 4, 30, 33, 35, 37, 43, 50] suffer from self-occlusion

and environmental occlusions, limiting real-world robustness. IMU-based approaches [16, 27, 42, 46] mitigate occlusion issues but struggle with cumulative drift and sensor instability, degrading long-term accuracy.

To address these limitations, recent work [5, 9, 14, 20, 47] has explored fusing visual and inertial information to improve accuracy. Although these approaches improve accuracy, several practical factors still limit their deployment: (i) **Limited Wearability**, as they often rely on bulky head-mounted cameras or dense IMU configurations; (ii) **Non-trivial Calibration**, requiring careful alignment across multiple heterogeneous sensors; and (iii) **High Hardware Cost**, since professional-grade IMUs remain expensive and inaccessible. To this end, we address these issues by **enabling human motion estimation using only lightweight, calibration-free everyday wearables**.

In this paper, we propose a novel human motion estimation framework, **EveryWear**, that leverages RGB images from glasses with three cameras (one forward, two downward) and inertial data from typical wearables: a smartphone, a smartwatch, and earbuds, all carried in everyday configurations. This setup is lightweight and practical for daily use. However, these sparse consumer-grade sensors also introduce unique challenges for accurate motion estimation: (i) **Low sensor stability**. Consumer IMUs are inherently less stable than professional hardware, due to lower sensor quality (data loss, signal noise) and loose everyday wear configurations (phones in pockets, watches with play). The complex instability makes realistic sensor modeling difficult, leading to the failure of previous synthetic-data-based training approaches built on idealized assumptions. (ii) **Limited motion observability**. The sparse sensors setup provides only partial body motion observations compared to dense professional IMU configurations, especially in camera occluded scenarios, leading to incomplete constraints for accurate full-body motion estimation.

To address these challenges, we propose a teacher-student distillation approach. The teacher policy leverages multimodal inputs: RGB images from three cameras and IMU data from two wrists, two legs, and the head obtained from MoCap system [41], to learn accurate motion estimation. We then distill this knowledge into a student policy that uses only sparse, noisy IMU measurements from consumer devices while maintaining the same camera inputs. This distillation process constrains the student to learn from the teacher’s knowledge obtained from precise sensor observations, while simultaneously adapting to noisy consumer-grade IMU data.

Importantly, while sparse sensors provide incomplete motion constraints individually, the multimodal design enables cross-modal compensation: when one modality becomes unreliable (*e.g.*, cameras are occluded or IMU drift occurs), other modalities compensate to maintain robust

motion estimation. Furthermore, we integrate an off-the-shelf SLAM module using the forward-facing camera to provide global localization and compensate for drift in both position and head pose estimation.

To facilitate further research in this direction, we introduce **Ego-Elec**, a large-scale dataset comprising RGB images from three cameras and IMU data from three consumer devices collected using the prototype of **EveryWear**. The dataset features 9 hours of real-world human motion across 56 types of daily activities in 17 diverse environments, with ground-truth 3D body poses and global translations annotated using the motion capture (MoCap) system. By training directly on real-world data, our approach achieves robust generalization to real-world scenarios without the sim-to-real gap that limits synthetic-data-based methods. Moreover, the diverse activities and environments make our dataset valuable for future research and broadly applicable across various applications.

We conduct comprehensive experiments to validate the effectiveness of our approach. Our method achieves 8.459 cm MPJPE and 10.627 cm MPJVE, significantly outperforming baselines by 3.345 cm and 4.289 cm, demonstrating that distillation from accurate teacher observations helps the model adapt to noisy consumer-grade sensors. Ablation studies further validate our design choices, confirming that (i) multimodal fusion enables robust cross-modal compensation, and (ii) our approach maintains performance under occlusion scenarios where single-modality methods fail.

In summary, our contributions are as follows:

- A practical motion capture framework using only everyday consumer devices (smartphone, smartwatch, earbuds, smart glasses) with no calibration required, achieving robust motion estimation through teacher-student distillation and multimodal fusion.
- Ego-Elec, a comprehensive real-world dataset with 9 hours of egocentric motion across 56 activities and 17 environments, the first large-scale dataset combining egocentric vision with sparse consumer IMU measurements.
- State-of-the-art results demonstrating that our approach outperforms existing methods while using only consumer devices, with robust performance maintained under occlusions through cross-modal compensation.

2. Related Work

2.1. On-Body Device-Based Motion Capture

On-body device-based human motion capture systems fall into three categories: IMU-only, camera-only, and multimodal approaches. **IMU-only methods** use sparse sensors attached to body locations such as wrists, knees, or other joints [16, 34, 39, 52]. Early work used professional IMUs, which are costly and inaccessible for daily use. Recent approaches [7, 27, 42] leverage consumer electronics (smart-

phones, smartwatches), but sparse consumer IMU configurations provide limited body coverage, constraining motion estimation accuracy. **Camera-only methods** employ egocentric cameras mounted on the head [3, 30, 33, 36, 50] or chest [29]. While cameras provide rich visual observations, they struggle with occlusions caused by furniture, walls, or body parts blocking the camera’s view. **Multimodal methods** combine head-mounted cameras, chest cameras, and inertial sensors for comprehensive observations and cross-modal compensation [9, 12, 25, 38]. However, these systems burden users with multiple body-mounted sensors and complex calibration procedures, limiting practical deployment. In contrast, we adopt a lightweight configuration centered on glasses-mounted cameras, complemented by everyday consumer devices. This design preserves comprehensive egocentric observations through multimodal fusion while avoiding the hardware burden and calibration requirements of existing systems, making it practical for daily use.

2.2. On-Body-Device-Based Motion Dataset

Corresponding to the three types of motion capture methods discussed above, existing datasets for human motion estimation can also be categorized in the same way: camera-only [1, 3], IMU-only [7, 26, 27, 42], and multimodal approaches [1, 11–13, 23, 38, 45]. These datasets have significantly advanced human motion estimation research by providing diverse activities and motion patterns. However, most of these works rely on synthetic data to avoid the labor-intensive process of collecting and annotating real-world data. While synthetic data enables easy scaling with different sensor modalities, it introduces a significant sim-to-real gap that limits real-world deployment. First, real-world IMU measurements are inherently noisy and affected by environmental disturbances (magnetic interference, temperature drift) that are difficult to accurately simulate. Second, everyday wear configurations, such as phones loosely placed in pockets or watches with some wrist play, introduce relative movement between sensors and the body that is challenging to model in simulation. Recent efforts [3, 15, 23, 32, 38, 43] have attempted to combine synthetic and real data, but they still inherit the fundamental limitations of synthetic sensor modeling. To address this gap, we introduce the first large-scale real-world dataset that combines egocentric vision with sparse consumer IMU measurements. By training directly on real-world data, our approach eliminates the sim-to-real gap and enables robust deployment in practical scenarios.

2.3. Egocentric Motion Estimation

Previous single-modality approaches [3, 6, 31] process camera or IMU inputs independently, either through end-to-end regression or via 2D heatmap as intermediate representations. Multimodal methods [10] process modalities

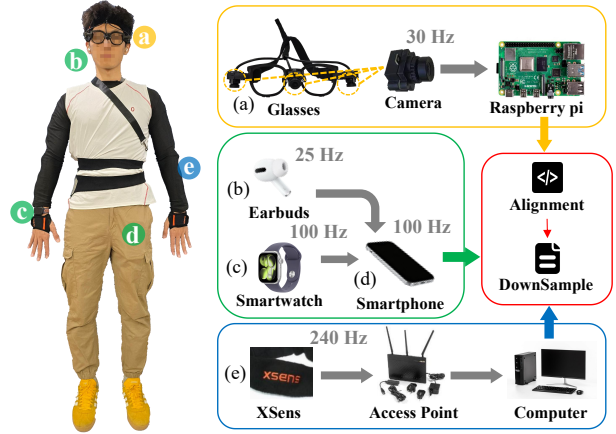


Figure 2. **System illustration.** The system comprises smart glasses (with onboard Raspberry Pi), smartphone, smartwatch, and earbuds. An XSens MoCap system provides annotations.

separately and fuse pose estimate results directly, failing to exploit cross-modal feature interactions. Critically, most existing methods train on synthetic IMU data [1, 3] and train directly on sparse sensors, limiting their applicability to real-world scenarios with sparse, noisy consumer devices. Our teacher-student distillation framework addresses these limitations by (i) training on real-world data to capture true noise characteristics, and (ii) learning to estimate poses from sparse consumer sensors through knowledge transfer from a teacher trained on dense, accurate sensors. Our multimodal fusion enables robust cross-modal compensation when individual sensors become unreliable.

3. Setup and Dataset

3.1. Hardware Setup

As shown in Fig. 2, our system comprises everyday consumer devices: a smartphone, a smartwatch, wireless earbuds, and custom smart glasses equipped with three cameras (one forward-facing, two downward-facing). All cameras are connected to an onboard Raspberry Pi 5, capturing synchronized RGB images at a resolution of 1080×720 and a frame rate of 30 Hz. IMU data from the smartphone, smartwatch, and earbuds are recorded at 100, 100, and 25 Hz, respectively.

During data collection, users were equipped with the system in typical everyday configurations: smart glasses on the head, a smartwatch on either wrist, earbuds in both ears, and a smartphone loosely placed in a pants pocket (left or right). To obtain ground-truth annotations, participants simultaneously wear the MoCap system [41], which provides 3D joint positions and global translations at 240 Hz using 17 body-mounted IMU sensors.

Table 1. Comparison with existing on-body device motion estimation datasets. **Type**: synthetic or real-world; **I/O**: indoor (**I**) or outdoor (**O**); **Cams**: number of cameras; **IMUs**: number of IMUs; **Acts**: number of activity types; **Frames**: number of frames.

Datasets	Type	I/O	Cams	IMUs	Acts	Frames
M2C2 [43]	Syn.	-	1	-	3K	530K
xR-EP [33]	Syn.	-	1	-	9	380K
UnrealEgo [1]	Syn.	-	2	-	30	900K
xw EgoCap [30]	Real	I/O	2	-	-	75K
EgoGlass [50]	Real	I	2	-	6	173K
EgoPW [36]	Real	I/O	1	-	20	318K
DIP-IMU [16]	Real	I	-	6	15	330K
IMUPoser [27]	Real	I	-	3	36	-
FRAME [4]	Real	I	2	-	50	1.6M
MEV-R [15]	Real	O	6	-	35	3.12M
Ego4View-R [3]	Real	I	4	-	-	930K
ColossusEgo [21]	Real	I	2	-	-	2.8M
Ego-VIP [5]	Real	I	4	4	-	38K
EMHI [9]	Real	I	2	5	39	3.07M
Ego-Elec	Real	I/O	3	3	56	2.88M

3.2. Multimodality Data Alignment

Since the cameras, IMU sensors, and motion capture system operate independently, their data streams are neither temporally synchronized nor uniformly sampled. To address this, participants perform predefined calibration gestures at the start and end of each recording session, which we use to temporally align all modalities to a standard reference frame (details in *supplementary material*). We then downsample all streams to 25 Hz, the earbuds’ native rate, and the lowest in our system, using nearest-neighbor sampling.

3.3. Ego-Elec Dataset

Our dataset was collected using the hardware setup described in Sec. 3.1. Totally, our dataset provides 9 **hours** of real-world human motion across 56 **types** of daily activities in 17 **diverse environments**. To contextualize its contributions, we compare it with existing datasets in Tab. 1. As shown, we provide the most diverse and comprehensive real-world data across modalities and environments, comparable in scale to prior work. As visualized in *supplementary material*, each recording session consists of 5-10 scripted daily activities performed continuously in one of 17 distinct indoor and outdoor environments, with an average duration of 4 minutes per session. Activities are drawn from a taxonomy of 56 distinct daily activity types (details in *supplementary material*), including sports, social interactions, household tasks, and work-related actions. This design captures naturalistic motion sequences suitable for both motion estimation and long-horizon egocentric tasks.

4. Method

4.1. Overview

Our objective is to estimate human motion $\mathcal{M} = \{\mathbf{P}_t, \Theta_t\}_{t=0}^T$, where $\mathbf{P}_t \in \mathbb{R}^3$ is global root translation and

$\Theta_t = \{\theta_j\}_{j=0}^J \in \mathbb{R}^{24 \times 6}$ contains $J = 24$ joint rotations in 6D representation [51] for SMPL [24] at time t over sequence length T .

We propose a teacher-student model (detailed in Sec. 4.3) that takes multi-view images $\mathcal{I} = \{\mathbf{I}_t^f, \mathbf{I}_t^l, \mathbf{I}_t^r\}_{t=0}^T$ and IMU signals $\mathcal{V}_{t=0}^T$ from five body-worn sensors as input to estimate motion sequence \mathcal{M} . Additionally, we employ an off-the-shelf SLAM module (Sec. 4.2) to estimate global head position, complementing the earbud’s orientation measurements. $\mathbf{I}_t^f, \mathbf{I}_t^l, \mathbf{I}_t^r$ denote frames from forward-facing, left downward-facing, and right downward-facing cameras at time t . IMU signals $\mathcal{V} \subseteq \{\mathbf{V}^h, \mathbf{V}^{lw}, \mathbf{V}^{rw}, \mathbf{V}^{lh}, \mathbf{V}^{rh}\}$ come from: head (earbuds), left wrist (smartwatch), right wrist (smartwatch), left hip (smartphone), and right hip (smartphone), respectively. Each sensor provides $\mathbf{V} = \{\mathbf{R}_t, \mathbf{A}_t\}_{t=0}^T$, where \mathbf{R}_t and \mathbf{A}_t represent orientation and acceleration respectively. The overall architecture of our model is illustrated in Fig. 3.

4.2. Monocular SLAM Module

We employ the off-the-shelf SLAM module MAST3R-SLAM [28] to estimate global camera poses $\mathbf{C} = \{\mathbf{C}_t\}_{t=0}^T$ from the forward camera sequence $\{\mathbf{I}_t^f\}_{t=0}^T$. We then apply a rigid transformation to obtain head poses $\mathcal{H} = \{\mathbf{H}_t\}_{t=0}^T$ from the camera poses \mathbf{C} . Although the global translation scale differs between SLAM and IMU coordinate systems, our teacher-student model implicitly learns to align these frames during training.

4.3. Teacher-Student Model

4.3.1. Teacher Model

We formulate motion estimation as a sequence-to-sequence problem using a sliding window of length N . The model takes visual and IMU observations as input, extracts features using visual and IMU feature encoders, and then fuses the multimodal features through a temporal fusion module to estimate human motion, as illustrated in Fig. 3.

Visual Feature Encoder We employ a ResNet-18 backbone pre-trained on ImageNet to extract visual features $\mathbf{F}_t^f, \mathbf{F}_t^l$, and \mathbf{F}_t^r from input images $\mathbf{I}_t^f, \mathbf{I}_t^l, \mathbf{I}_t^r$ at each time step t .

IMU Feature Encoder An MLP-based encoder processes all IMU signals $\mathcal{V}_t^{tea} = \{\mathbf{V}_t^h, \mathbf{V}_t^{lw}, \mathbf{V}_t^{rw}, \mathbf{V}_t^{lh}, \mathbf{V}_t^{rh}\}$ into a unified latent representation \mathbf{U}_t . Note that during teacher training, we use dense IMU data from the motion capture system (5 IMUs) rather than the sparse consumer sensors (3 IMUs), enabling the teacher to learn from more accurate and complete motion observations.

Temporal Fusion A bidirectional LSTM processes the concatenated features $\{\mathbf{F}_t^f, \mathbf{F}_t^l, \mathbf{F}_t^r, \mathbf{U}_t, \mathbf{H}_t\}_{t=t-N}^{t+N}$ over the sliding window to estimate motion sequence \mathcal{M} , where \mathbf{H}_t denotes the SLAM-derived head pose detailed in Sec. 4.2.

Loss Function The teacher model is jointly optimized with two objectives: the local pose loss \mathcal{L}_{pose} and the global

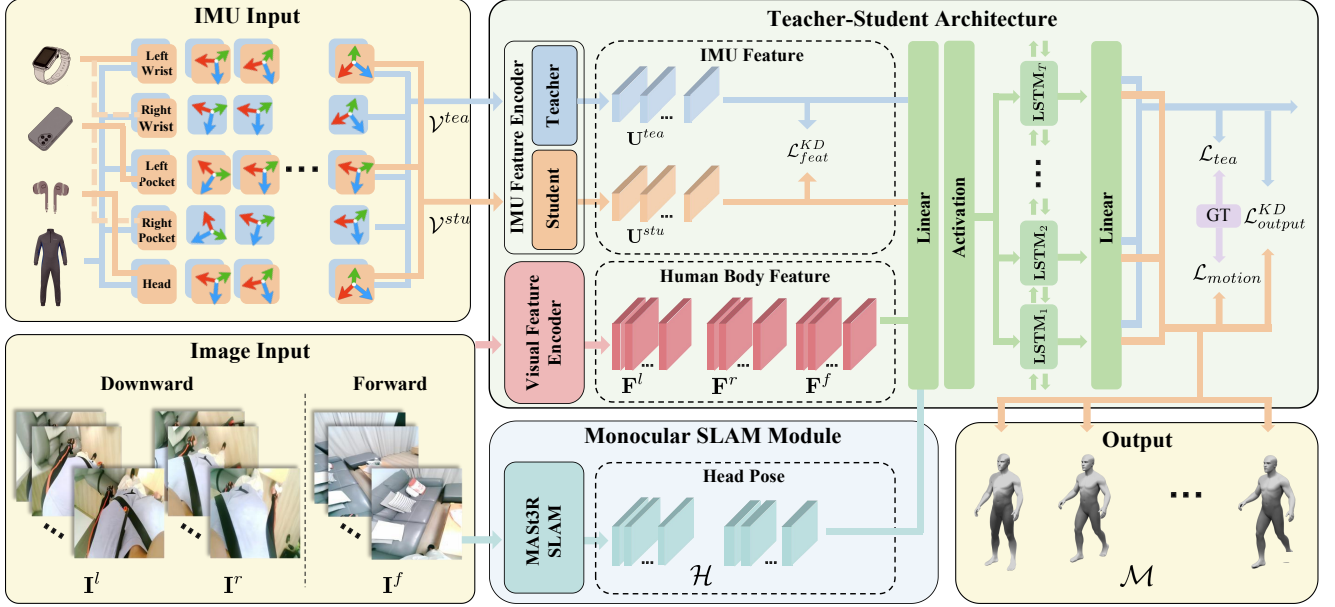


Figure 3. **Pipeline.** Our method takes egocentric images (three cameras) and IMU signals (everyday wearables) as input. We employ SLAM (Sec. 4.2) for head pose estimation, then use a teacher-student framework (Sec. 4.3) with shared visual feature encoder and separate IMU feature encoders, followed by bidirectional LSTM for temporal fusion and motion prediction.

translation loss \mathcal{L}_{trans} . Following Kendall et al. [19], we use learned task-dependent uncertainty weighting. The total teacher loss is:

$$\mathcal{L}_{tea} = \frac{1}{\sigma_{pose}^2} \mathcal{L}_{pose} + \log \sigma_{pose}^2 + \frac{1}{\sigma_{trans}^2} \mathcal{L}_{trans} + \log \sigma_{trans}^2, \quad (1)$$

where σ_{pose}^2 and σ_{trans}^2 are learnable uncertainty parameters that adaptively weight the losses during training.

Local Pose Loss The pose loss is a weighted mean squared error over the 6D rotation representations of all $J = 24$ body joints:

$$\mathcal{L}_{pose} = \frac{1}{J} \sum_{j=1}^J w_j \|\theta_j - \hat{\theta}_j\|_2^2, \quad (2)$$

where $\theta, \hat{\theta} \in \mathbb{R}^6$ are ground-truth and predicted rotations, and $w_j > 0$ are per-joint weights (Refer to *supplementary material* for detailed weights design.).

Global Translation Loss We implement the translation loss as the L2 norm between the predicted and ground-truth root positions: $\mathcal{L}_{trans} = \|\mathbf{P} - \hat{\mathbf{P}}\|_2^2$, where $\mathbf{P}, \hat{\mathbf{P}} \in \mathbb{R}^3$ are ground-truth and predicted root positions.

4.3.2. Student Model

The student model shares the teacher’s architecture but processes sparse IMU signals from consumer devices. While the teacher uses dense IMU data from 5 body locations, the student uses only 3 physical sensors: head (earbuds), wrist (smartwatch), and hips (smartphone): $\mathcal{V}_t^{stu} =$

$\{\mathbf{V}_t^{head}, \mathbf{V}_t^w, \mathbf{V}_t^{hip}\}$. This reflects typical daily usage with one watch and one phone.

The student model’s parameters are initialized with the teacher’s pre-trained weights and fine-tuned with a small learning rate to adapt to noisy consumer IMU signals. Training details are provided in the *supplementary material*.

Loss Functions The student model is trained with a multi-objective loss that combines three components: (i) a motion estimation loss \mathcal{L}_{motion} , which is the same as the \mathcal{L}_{tea} . (ii) an distillation loss $\mathcal{L}_{output}^{KD}$ that makes the student’s output $\mathbf{M} = \{\mathbf{P}_t, \Theta_t\}_{t=0}^T$ as same as the teacher’s. (iii) a feature-level distillation loss \mathcal{L}_{feat}^{KD} that aligns the student’s IMU feature with those in the teacher model.

The loss $\mathcal{L}_{output}^{KD}$ is defined as:

$$\mathcal{L}_{output}^{KD} = \frac{1}{J} \sum_{j=1}^J \|\hat{\theta}_j^T - \hat{\theta}_j^S\|_2^2, \quad (3)$$

where $\hat{\theta}_j^T$ and $\hat{\theta}_j^S$ denote the predicted 6D joint rotations of joint j from the teacher and student models, respectively.

The loss \mathcal{L}_{feat}^{KD} is defined as: $\mathcal{L}_{feat}^{KD} = \frac{1}{d} \|\mathbf{U}^{tea} - \mathbf{U}^{stu}\|_2^2$, where \mathbf{U}^{tea} and \mathbf{U}^{stu} represent the extracted IMU features in teacher and student models, respectively.

Formally, the total student loss \mathcal{L}_{stu} is defined as:

$$\mathcal{L}_{stu} = \lambda_{motion} \mathcal{L}_{motion} + \lambda_{output} \mathcal{L}_{output}^{KD} + \lambda_{feat} \mathcal{L}_{feat}^{KD}, \quad (4)$$

where λ_* are the balanced weights. We provide detailed weight settings in the *supplementary material*.

5. Experiments

5.1. Implementation and Metrics

Full implementation details are provided in *supplementary material*. To evaluate the proposed framework, we segment the dataset into non-overlapping sequences with a window size N , then randomly divide them into three subsets: 712 sequences for training (80%), 89 for validation (10%), and 89 for testing (10%).

Following prior work [20, 27, 42], we evaluate our model using the following metrics: *Mean Per Joint Position Error (MPJPE)*: Mean Euclidean distance between predicted and ground-truth joint positions (cm) with the pelvis aligned; *Procrustes Aligned Mean Joint Position Error (PA-MPJPE)*: MPJPE after Procrustes alignment, measuring pose shape accuracy (cm); *Mean Per Joint Rotation Error (MPJRE)*: Mean angular error between predicted and ground-truth joint orientations (degrees) with the pelvis aligned; *Mean Per Joint Vertex Error (MPJVE)*: Mean per-vertex error across all vertices of the SMPL mesh (cm) with the pelvis aligned; *Root Position Error (Root PE)*: Root position error in global coordinates (cm).

5.2. Comparison

To demonstrate the effectiveness of our approach, we conduct a comprehensive comparison against two representative categories of baseline methods on our dataset:

IMU-Only We compare against IMUPoser [27] and MobilePoser [42], sparse IMU methods using smartphones, smartwatches, and earbuds.

Camera-Only We adapt Fish2Mesh [31] for two downward-facing cameras instead of one fisheye camera.

In Table 2, we report quantitative results, which demonstrate that our full model outperforms all baselines across all metrics. The qualitative results in Fig. 4 also confirm these findings. We also present results on challenging scenarios (sitting with occlusions, waving with body parts exiting the camera view), demonstrating our model’s robustness across diverse conditions.

Comparison with IMU-based Baselines We compare our full model against IMU-only baselines and our IMU-only ablation on our test set. As shown in the first section of Tab. 2, our approach **significantly outperforms all baselines across all metrics**. Compared to the best baseline IMUPoser, our method achieves improvements of 3.35 cm MPJPE, 2.17 cm PA-MPJPE, 2.40 deg MPJRE, 4.29 cm MPJVE, and 12.12 cm Root PE, demonstrating the effectiveness of multimodal fusion. Notably, our IMU-only variant also achieves improvement over IMUPoser despite using the same sensor configuration. These **demonstrate three key insights**: (i) the sparse 3-IMU setup with loose everyday attachment poses significant challenges for pure inertial-based motion estimation. (ii) the addition of ego-

centric vision provides substantial benefits, enabling robust cross-modal compensation when individual sensors become unreliable due to noise or drift. (iii) our teacher-student distillation framework enables the student model to learn effective denoising and adaptation strategies for handling noisy and unstable consumer-grade IMU sensors.

Comparison with Camera-based Baselines We compare Fish2Mesh [31] against our full approach and our camera-only ablation. Our full model achieves 4.22 cm and 5.59 cm improvements over Fish2Mesh on MPJPE and MPJVE respectively, with similar gains across other metrics. This improvement can be attributed to the fundamental limitations of camera-only Methods: body parts leaving the field of view, environmental and self-occlusions, and ambiguity in head pose. Our approach addresses these by **fusing visual observations with IMU signals, which provide motion and orientation cues when cameras fail**. Importantly, our camera-only variant performs comparably to Fish2Mesh, demonstrating that both vision-only approaches share similar limitations.

Comparison in Occluded Scenario Table 3 evaluates performance on challenging scenarios: sitting (lower body occlusion) and waving (upper body exits view). We report region-specific Upper-MPJPE and Lower-MPJPE, which measure upper-body and lower-body joint errors, respectively. While baseline methods struggle with these challenging scenarios, achieving comparable errors to each other, our approach significantly outperforms them by approximately 4 cm on both upper and lower body metrics. This demonstrates effective multimodal fusion: IMU sensors compensate when cameras are occluded, while cameras track visible regions when body parts leave the frame.

Qualitative Comparisons Figure 4 presents qualitative comparisons across unoccluded and occluded scenarios. For the unoccluded scenes, rows 1-2 show camera-only methods failing when the upper body exits the field of view, while our method maintains accuracy via IMU data. Rows 2-3 demonstrate that IMU-only methods struggle to estimate leg movements with sparse sensors, a challenge our approach addresses by incorporating complementary visual information. For the occluded scenes, the last rows show scenarios where objects (boxes, books) block camera views. Camera-only methods fail under these occlusions, whereas our multimodal fusion enables robust tracking by leveraging IMU signals. This **cross-modal compensation and teacher-student distillation are key to handling real-world challenges**.

Failure Cases Figure 5 shows our method struggles when consumer IMU sensors become unstable due to loose attachment or sensor noise. Our model fails when the upper body exits the camera view while IMU signals are simultaneously unreliable, though the lower body maintains reasonable tracking. In severely occluded scenarios, our model

Table 2. Evaluation of **EveryWear** on **Ego-Elec**. The first section compares our model with baseline methods, while the second section presents ablation studies of our model design.

	Method	MPJPE (cm)	PA-MPJPE (cm)	MPJRE (deg)	MPJVE (cm)	Root PE (cm)
Baseline	(a) IMU-Only Method					
	IMUPoser [27]	11.804	7.647	12.300	14.916	24.499
	MobilePoser [42]	17.384	15.099	20.818	17.303	25.028
	(b) Camera-Only Method					
	Fish2Mesh [31]	12.677	7.049	11.036	16.213	-
Ablation	(a) Ablation on IMUs					
	Ours-w/o-IMUs	13.014	7.249	10.901	16.650	-
	Ours-w/o-Phone	9.447	6.183	9.968	11.786	14.774
	Ours-w/o-Watch	9.642	6.253	9.977	12.121	13.835
	Ours-w/o-Earbuds	10.709	6.164	10.247	13.615	14.642
	(b) Ablation on Cameras					
	Ours-w/o-Cams	11.284	6.927	11.219	14.405	14.441
	Ours-w/o-Cam.f.feature	9.691	6.249	10.080	12.231	13.211
	Ours-w/o-Cam.f.SLAM	9.024	5.966	9.680	11.299	16.937
	(c) Ablation on Teacher-Student Architecture					
	Ours-w/o-Teacher	9.193	6.258	9.896	11.631	13.055
	Ours	8.459	5.482	9.035	10.627	12.382

Table 3. Comparison of our methods with baseline models on a typical challenging scenario, including occluded (sitting) and the body exiting the camera view (waving).

Method	Sitting		Waving	
	Upper PE	Lower PE	Upper PE	Lower PE
IMUPoser [27]	9.230	8.946	21.197	12.472
Fish2Mesh [31]	9.369	6.058	22.582	12.052
Ours	5.315	4.541	18.465	9.674

exhibits reduced accuracy for occluded body parts while maintaining overall motion structure.

5.3. Ablation study

We conducted systematic ablations to assess the impact of various model configurations on human motion estimation performance, as shown in the second section of Tab. 2.

Ablation on IMUs We ablate IMU sensor configurations, and the results show that removing all IMUs degrades performance but remains comparable to the camera-only baseline, validating our model’s adaptability. Individual sensor ablations reveal that earbuds contribute most: removing earbuds causes 2.25 cm MPJPE and 2.99 cm MPJVE degradation versus 1-2 cm for others. Notably, MPJRE, which measures the local pose, degrades minimally while MPJPE and MPJVE increase > 2 cm, indicating earbuds’ critical role in global orientation estimation.

Ablation on Cameras We conduct camera ablation studies with three configurations: removing all cameras

(w/o-Cams), only excluding the forward-facing camera from motion estimation (w/o-Cam.f.feature), or only excluding it from SLAM (w/o-Cam.f.SLAM).

Individual Camera Contributions Results show expected patterns based on camera viewpoints. Removing the forward-facing camera from SLAM degrades Root PE (global localization) most severely, increasing error by 4.56 cm. This is expected, as SLAM provides environmental context critical for global positioning. Conversely, removing the forward-facing camera from motion estimation degrades MPJPE by 1.23 cm, as this camera directly observes the upper body.

All Cameras Removed (w/o-Cams) When all cameras are removed (IMU-only configuration), MPJPE increases by 2.83 cm and MPJVE by 3.78 cm, demonstrating substantial degradation. This validates the importance of multi-modal fusion: visual information significantly outperforms IMU-only estimation.

Interesting Finding We observe a counterintuitive result: removing *all* cameras (Root PE = 14.4 cm) achieves better global localization than removing only the forward camera (Root PE = 16.9 cm). This suggests that downward-facing cameras, while beneficial for body pose estimation, may introduce conflicting signals for global localization when environmental context from the forward camera is absent. The model handles pure IMU-based global tracking more robustly than partial visual information lacking environmental context.

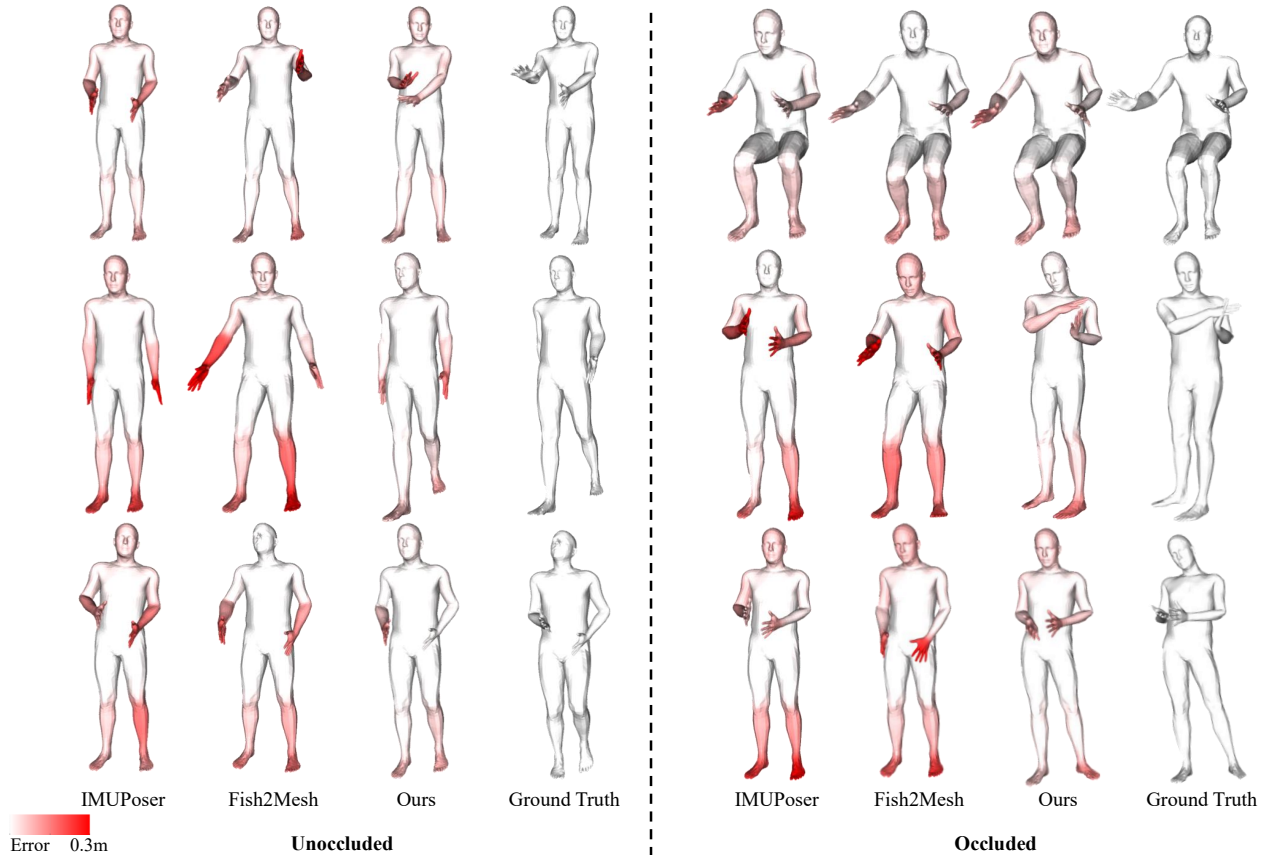


Figure 4. **Qualitative comparison.** This figure presents a comparison against baseline methods on Ego-Elec, where IMUPoser serves as the IMU-only baseline and Fish2Mesh represents the camera-only baseline. We visualize per-vertex SMPL error using a color map ranging from 0 to 0.3 m (white: low error, red: high error). Our method achieves consistently lower errors across diverse activities and maintains robustness under challenging occlusion scenarios.

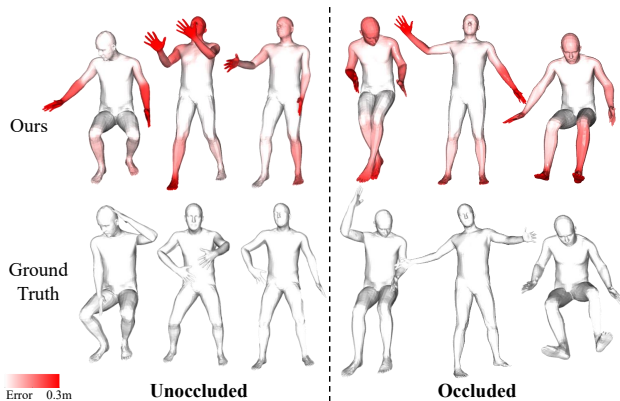


Figure 5. **Failure cases.** We show rare failure scenarios where loose attachment of consumer devices (phone shifting in pocket, watch with play) introduces sensor noise and instability.

Teacher-Student Model We ablate the teacher-student framework by training without distillation. Performance degrades by 0.8 cm across all metrics (MPJPE, PA-MPJPE, MPJVE), confirming that the distillation enables effective adaptation to noisy consumer sensors.

6. Discussions

Conclusions We present Ego-Elec, a lightweight human motion capture method using everyday wearables without calibration, enabled by multimodal fusion and teacher-student distillation. Through this distillation paradigm, our model learns to adapt to noisy real-world consumer sensors. We also introduce a large-scale real-world dataset spanning diverse daily activities and environments. By achieving state-of-the-art performance and providing comprehensive real-world data, our work establishes a foundation for future research on full-body motion estimation with consumer devices, with applications in XR gaming, telepresence, and humanoid teleoperation.

Limitations Our current approach operates offline, but can be adapted to real-time use for widespread real-world applications. The IMU configuration currently supports only fixed sensor placements and requires prior specification, which can be improved by enabling adaptive handling of variable sensor configurations. Future work will further incorporate richer interaction modalities to enhance robustness in complex real-world environments.

References

- [1] Hiroyasu Akada, Jian Wang, Soshi Shimada, Masaki Takahashi, Christian Theobalt, and Vladislav Golyanik. Unrealego: A new dataset for robust egocentric 3d human motion capture. In *Proceedings of European Conference on Computer Vision (ECCV)*, 2022. 1, 3, 4
- [2] Hiroyasu Akada, Jian Wang, Vladislav Golyanik, and Christian Theobalt. 3d human pose perception from egocentric stereo videos. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 1
- [3] Hiroyasu Akada, Jian Wang, Vladislav Golyanik, and Christian Theobalt. Bring your rear cameras for egocentric 3d human pose estimation. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2025. 3, 4
- [4] Andrea Boscolo Camiletto, Jian Wang, Eduardo Alvarado, Rishabh Dabral, Thabo Beeler, Marc Habermann, and Christian Theobalt. Frame: Floor-aligned representation for avatar motion from egocentric video. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025. 1, 4
- [5] Young-Woon Cha, Husam Shaik, Qian Zhang, Fan Feng, Andrei State, Adrian Ilie, and Henry Fuchs. Mobile. egocentric human body motion reconstruction using only eyeglasses-mounted cameras and a few body-worn inertial sensors. In *IEEE Virtual Reality and 3D User Interfaces (VR)*, 2021. 2, 4
- [6] Hanz Cuevas-Velasquez, Charlie Hewitt, Sadegh Aliakbarian, and Tadas Baltrušaitis. Simpleego: Predicting probabilistic body pose from egocentric cameras. In *International Conference on 3D Vision (3DV)*, 2024. 3
- [7] Nathan DeVrio, Vimal Mollyn, and Chris Harrison. Smartposer: Arm pose estimation with a smartphone and smartwatch using uwb and imu data. In *ACM Symposium on User Interface Software and Technology (UIST)*, 2023. 2, 3
- [8] Yushi Du, Yixuan Li, Baoxiong Jia, Yutang Lin, Pei Zhou, Wei Liang, Yanchao Yang, and Siyuan Huang. Learning human-humanoid coordination for collaborative object carrying, 2025. 1
- [9] Zhen Fan, Peng Dai, Zhuo Su, Xu Gao, Zheng Lv, Jiarui Zhang, Tianyuan Du, Guidong Wang, and Yang Zhang. Emhi: A multimodal egocentric human motion dataset with hmd and body-worn imus. In *Proceedings of AAAI Conference on Artificial Intelligence (AAAI)*, 2025. 2, 3, 4
- [10] Andrew Gilbert, Matthew Trumble, Charles Malleson, Adrian Hilton, and John Collomosse. Fusing visual and inertial sensors with semantics for 3d human pose estimation. *International Journal of Computer Vision (IJCV)*, 127(4):381–397, 2019. 3
- [11] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, and et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [12] Kristen Grauman, Andrew Westbury, Lorenzo Torresani, Kris Kitani, Jitendra Malik, and et al. Ego-exo4d: Understanding skilled human activity from first- and third-person perspectives. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 3
- [13] Vladimir Guzov, Aymen Mir, Torsten Sattler, and Gerard Pons-Moll. Human poseitioning system (hps): 3d human pose estimation and self-localization in large scenes from body-mounted sensors. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [14] Vladimir Guzov, Aymen Mir, Torsten Sattler, and Gerard Pons-Moll. Human poseitioning system (hps): 3d human pose estimation and self-localization in large scenes from body-mounted sensors. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [15] Dominik Hollidt, Paul Streli, Jiaxi Jiang, Yasaman Haghighi, Changlin Qian, Xintong Liu, and Christian Holz. Egosim: An egocentric multi-view simulator and real dataset for body-worn cameras during motion and activity. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2024. 1, 3, 4
- [16] Yinghao Huang, Manuel Kaufmann, Emre Aksan, Michael J Black, Otmar Hilliges, and Gerard Pons-Moll. Deep inertial poser: Learning to reconstruct human pose from sparse inertial measurements in real time. *ACM Transactions on Graphics (TOG)*, 37(6):1–15, 2018. 2, 4
- [17] Mazeyu Ji, Xuanbin Peng, Fangchen Liu, Jialong Li, Ge Yang, Xuxin Cheng, and Xiaolong Wang. Exbody2: Advanced expressive humanoid whole-body control. *arXiv preprint arXiv:2412.13196*, 2024. 1
- [18] Jiaxi Jiang, Paul Streli, Manuel Meier, and Christian Holz. Egoposer: Robust real-time egocentric pose estimation from sparse and intermittent observations everywhere. In *Proceedings of European Conference on Computer Vision (ECCV)*, 2024. 1
- [19] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 5
- [20] Jiye Lee and Hanbyul Joo. Mocap everyone everywhere: Lightweight motion capture with smartwatches and a head-mounted camera. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2, 6
- [21] Jihyun Lee, Weipeng Xu, Alexander Richard, Shih-En Wei, Shunsuke Saito, Shaojie Bai, Te-Li Wang, Minhyuk Sung, Tae-Kyun Kim, and Jason Saragih. Rewind: Real-time egocentric whole-body motion diffusion with exemplar-based identity conditioning. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025. 4
- [22] Yixuan Li, Yutang Lin, Jieming Cui, Tengyu Liu, Wei Liang, Yixin Zhu, and Siyuan Huang. Clone: Closed-loop whole-body humanoid teleoperation for long-horizon tasks. In *Conference on Robot Learning (CoRL)*, 2025. 1
- [23] Yunze Liu, Yun Liu, Che Jiang, Kangbo Lyu, Weikang Wan, Hao Shen, Boqiang Liang, Zhoujie Fu, He Wang, and Li Yi. Hoi4d: A 4d egocentric dataset for category-level human-object interaction. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [24] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, 2023. 4

- [25] Lingni Ma, Yuting Ye, Fangzhou Hong, Vladimir Guзов, Yifeng Jiang, Rowan Postyeni, Luis Pesqueira, Alexander Gamino, Vijay Baiyya, Hyo Jin Kim, Kevin Bailey, David Soriano Fosas, C. Karen Liu, Ziwei Liu, Jakob Engel, Renzo De Nardi, and Richard Newcombe. Nymeria: A massive collection of multimodal egocentric daily motion in the wild. In *ECCV*, 2024. 3
- [26] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. Amass: Archive of motion capture as surface shapes. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2019. 3
- [27] Vimal Mollyn, Riku Arakawa, Mayank Goel, Chris Harrison, and Karan Ahuja. Imuposer: Full-body pose estimation using imus in phones, watches, and earbuds. In *ACM Conference on Human Factors in Computing Systems (CHI)*, 2023. 2, 3, 4, 6, 7
- [28] Riku Murai, Eric Dexheimer, and Andrew J Davison. Mast3r-slam: Real-time dense slam with 3d reconstruction priors. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025. 4
- [29] Keishi Nishikawa, Muhammad Taufique Popal, and Jun Ohya. Egocentric pose estimation using the image acquired by omni-directional camera attached on chest and cubemap. In *International Conference on Control, Automation and Robotics (ICCAR)*, 2025. 3
- [30] Helge Rhodin, Christian Richardt, Dan Casas, Eldar Insafutdinov, Mohammad Shafiee, Hans-Peter Seidel, Bernt Schiele, and Christian Theobalt. Egocap: egocentric marker-less motion capture with two fisheye cameras. *ACM Transactions on Graphics (TOG)*, 35(6):1–11, 2016. 1, 3, 4
- [31] Tianma Shen, Aditya Puranik, James Vong, Vrushabh Degirirkar, Ryan Fell, Julianna Dietrich, Maria Kyrarini, Christopher Kitts, and David C Jeong. Fish2mesh transformer: 3d human mesh recovery from egocentric vision. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2025. 3, 6, 7
- [32] Lan Sun, Songpengcheng Xia, Junyuan Deng, Jiarui Yang, Zengyuan Lai, Qi Wu, and Ling Pei. Suite-in: Aggregating motion features from apple suite for robust inertial navigation. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2025. 3
- [33] Denis Tome, Patrick Peluse, Lourdes Agapito, and Hernan Badino. xr-egopose: Egocentric 3d human pose from an hmd camera. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2019. 1, 3, 4
- [34] Tom Van Wouwe, Seunghwan Lee, Antoine Falisse, Scott Delp, and C Karen Liu. Diffusionposer: Real-time human motion reconstruction from arbitrary sparse sensors using autoregressive diffusion. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2
- [35] Jian Wang, Lingjie Liu, Weipeng Xu, Kripasindhu Sarkar, and Christian Theobalt. Estimating egocentric 3d human pose in global space. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2021. 1
- [36] Jian Wang, Lingjie Liu, Weipeng Xu, Kripasindhu Sarkar, Diogo Luvizon, and Christian Theobalt. Estimating egocentric 3d human pose in the wild with external weak supervision. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3, 4
- [37] Jian Wang, Diogo Luvizon, Weipeng Xu, Lingjie Liu, Kripasindhu Sarkar, and Christian Theobalt. Scene-aware egocentric 3d human pose estimation. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 1
- [38] Xin Wang, Taein Kwon, Mahdi Rad, Bowen Pan, Ishani Chakraborty, Sean Andrist, Dan Bohus, Ashley Feniello, Bugra Tekin, Felipe Vieira Frujeri, Neel Joshi, and Marc Pollefeys. Holoassist: an egocentric human interaction dataset for interactive ai assistants in the real world. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2023. 3
- [39] Xuan Xiao, Jianjian Wang, Pingfa Feng, Ao Gong, Xiangyu Zhang, and Jianfu Zhang. Fast human motion reconstruction from sparse inertial measurement units considering the human shape. *Nature Communications*, 2024. 2
- [40] Xianghui Xie, Jan Eric Lenssen, and Gerard Pons-Moll. Intertrack: Tracking human object interaction without object templates. In *International Conference on 3D Vision (3DV)*, 2025. 1
- [41] Xsens Technologies B.V. Xsens mvn link motion capture system, 2015. Product, Enschede, The Netherlands. 2, 3
- [42] Vasco Xu, Chenfeng Gao, Henry Hoffmann, and Karan Ahuja. Mobileposer: Real-time full-body pose estimation and 3d human translation from imus in mobile consumer devices. In *ACM Symposium on User Interface Software and Technology (UIST)*, 2024. 1, 2, 3, 6, 7
- [43] Weipeng Xu, Avishek Chatterjee, Michael Zollhoefer, Helge Rhodin, Pascal Fua, Hans-Peter Seidel, and Christian Theobalt. Mo 2 cap 2: Real-time mobile 3d motion capture with a cap-mounted fisheye camera. *IEEE Transactions on Visualization and Computer Graph (TVCG)*, 25(5):2093–2101, 2019. 1, 3, 4
- [44] Pradyumna Yalandur Muralidhar, Yuxuan Xue, Xianghui Xie, Margaret Kostyrko, and Gerard Pons-Moll. Physic: Physically plausible 3d human-scene interaction and contact from a single image. In *ACM SIGGRAPH / Eurographics Symposium on Computer Animation (SCA)*, 2025. 1
- [45] Jingkan Yang, Shuai Liu, Hongming Guo, Yuhao Dong, Xiamengwei Zhang, Sicheng Zhang, Pengyun Wang, Zitang Zhou, Binzhu Xie, Ziyue Wang, et al. Egolife: Towards egocentric life assistant. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025. 3
- [46] Xinyu Yi, Yuxiao Zhou, and Feng Xu. Transpose: Real-time 3d human translation and pose estimation with six inertial sensors. *ACM Transactions on Graphics (TOG)*, 40(4):1–13, 2021. 2
- [47] Xinyu Yi, Yuxiao Zhou, Marc Habermann, Vladislav Golyanik, Shaohua Pan, Christian Theobalt, and Feng Xu. Egolocate: Real-time motion capture, localization, and mapping with sparse body-mounted sensors. *ACM Transactions on Graphics (TOG)*, 2023. 2
- [48] Shaofeng Yin, Yanjie Ze, Hong-Xing Yu, C Karen Liu, and Jiajun Wu. Visualmimic: Visual humanoid locomanipulation via motion tracking and generation. *arXiv preprint arXiv:2509.20322*, 2025. 1

- [49] Xiaohan Zhang, Bharat Lal Bhatnagar, Sebastian Starke, Ilya A. Petrov, Vladimir Guzov, Helisa Dhamo, Eduardo Pérez Pellitero, and Gerard Pons-Moll. Force: Dataset and method for intuitive physics guided human-object interaction. In *International Conference on 3D Vision (3DV)*, 2025. [1](#)
- [50] Dongxu Zhao, Zhen Wei, Jisan Mahmud, and Jan-Michael Frahm. Egoglass: Egocentric-view human pose estimation from an eyeglass frame. In *International Conference on 3D Vision (3DV)*, 2021. [1](#), [3](#), [4](#)
- [51] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. [4](#)
- [52] Zunjie Zhu, Yan Zhao, Yihan Hu, Guoxiang Wang, Hai Qiu, Bolun Zheng, Chenggang Yan, and Feng Xu. Progressive inertial poser: Progressive real-time kinematic chain estimation for 3d full-body pose from three imu sensors. *IEEE Transactions on Instrumentation and Measurement (TIM)*, 2025. [2](#)

Human Motion Estimation with Everyday Wearables

Supplementary Material

A. Ego-Elec Dataset

We provide statistical analysis of the **Ego-Elec** dataset, characterizing environment types, activity categories, activity durations, and their distributions.

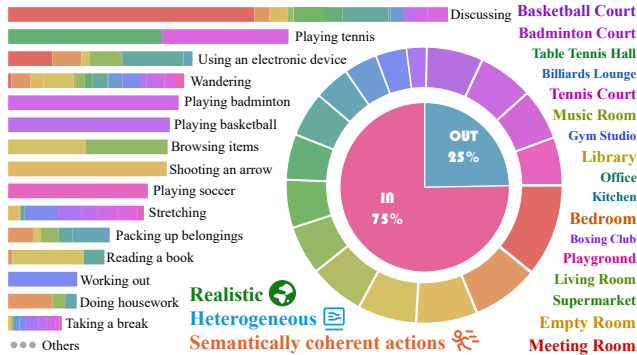


Figure A1. **Pie Chart: Distribution of Environments.** The dataset comprises a variety of environments, with their respective quantities illustrated by the proportions shown in the pie chart. **Bar Chart: Distribution of Activity Types.** Ego-Elec includes diverse activity categories. Here, we present the most frequent activities, along with the distribution of environments associated with each activity.

Environment Diversity As shown in the pie chart in Fig. A1, our dataset spans 17 distinct environment types across indoor and outdoor settings, with their relative proportions indicated. This diversity reflects the variety of real-world scenarios encountered in daily life.

Activity Coverage The bar chart in Fig. A1 shows the 56 most common daily activities in our dataset. Activity durations range from 50 to several thousand frames (Fig. A2), capturing both short actions and extended activities.

Data Collection Protocol We recorded full sessions containing 5–10 activities arranged in contextually meaningful sequences, then segmented each session into individual activities.

B. System Setup

B.1. Glasses Prototype

We developed a glasses prototype (Fig. A3) with three RGB cameras (1280×720 @ 30 FPS, 120° FOV): one forward-facing camera for environmental observation, and two downward-facing cameras (left and right) for egocentric body observation.

We will release all hardware designs, assembly instructions, calibration procedures, and data collection code to facilitate future research.

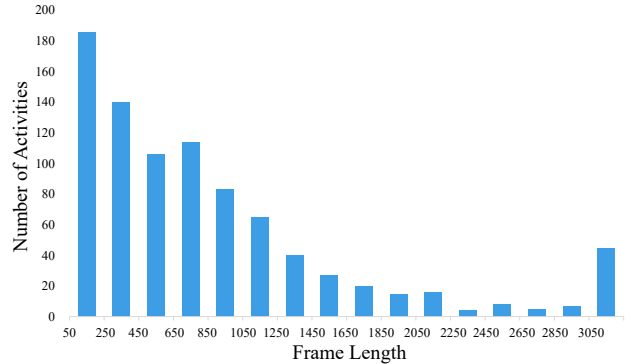


Figure A2. **Frame Length Distribution across Activities.** This figure shows the distribution of frame lengths across different activities, ranging from about 50 frames to several thousand.



Figure A3. **Glasses Prototype.** The glasses prototype is equipped with one forward-facing camera and two downward-facing cameras, and it runs on a Strawberry 5 platform.

B.1.1. Camera Synchronization

To achieve frame-level synchronization across the three cameras, we implement a software-based synchronization mechanism built around three core components: a **barrier-based capture trigger**, a **global frame index system**, and a **asynchronous storage pipeline**. Each camera operates in its own capture thread, coordinated by a centralized controller that manages shared synchronization primitives.

Barrier-Based Capture Trigger Each camera thread runs continuously but must wait at a shared barrier before capturing the next frame. The barrier releases only when all three threads have reached it, triggering simultaneous frame capture. This mechanism eliminates temporal drift caused by asynchronous execution or uneven driver latency.

Global Frame Index System A global frame counter ensures consistent labeling across all camera views. In each synchronized capture cycle, the first thread that exits the barrier assigns the next global frame index, and the remaining threads reuse that same index. Immediately after capture, each thread records a timestamp and pushes it into a camera-specific queue. This guarantees temporally consistent timestamps across all cameras.

Asynchronous Storage Pipeline To prevent disk I/O latency from interrupting capture operations, all frames are placed into per-camera queues and saved asynchronously by dedicated writer threads. Timestamps are flushed only after all cameras have submitted the same global frame index, ensuring completeness. This design maintains high capture throughput and prevents queue backpressure from degrading synchronization accuracy.

C. Implementation Details

C.1. Data Alignment

Since the cameras, everyday wearables, and MoCap system run on different devices, their timestamps are not inherently synchronized. To achieve temporal alignment across all modalities, we introduce two simple alignment gestures performed at the start of each recording session.

Gesture 1 (Wearables–MoCap Alignment) The user extends their right arm with the palm facing upward, places the smartphone on the palm, and raises the arm vertically three times while keeping the phone horizontal. Using this gesture, we compute the z-axis acceleration from both the smartphone IMU and the MoCap right-hand sensor. By temporally aligning the peak responses from the three repeated arm raises, we compute the time offset needed to synchronize the smartphone IMU with the MoCap system. Since the smartphone, smartwatch, and earbuds share a synchronized clock (all connected to the same mobile device), aligning the smartphone with the MoCap system automatically synchronizes all consumer wearable devices.

Gesture 2 (Cameras–MoCap Alignment) The user stands upright, quickly turns their head to the right, and then returns to a forward-facing position. From the camera streams, we select the frame where the camera’s field of view reaches its farthest rightward extent, corresponding to the moment of maximum head rotation. We then match this frame to the extremum of the z-axis orientation recorded by the MoCap head sensor, yielding the time offset between the cameras and the MoCap system for alignment.

Together, these two gestures provide accurate temporal synchronization across the cameras, wearables, and the MoCap ground-truth system.

C.2. Data Preprocess

Input images are resized to 224×224 for visual feature encoding, while the SLAM module operates on the original camera resolution (1280×720) to retain sufficient spatial detail for reliable tracking. To reduce high-frequency noise in the IMU signals while preserving natural motion dynamics, we apply a sliding-window smoothing filter with a window size of 5. All ground-truth global translations are expressed in a head-relative coordinate frame, computed as the relative transformation with respect to the first frame of

Table A1. Per-joint weights used in \mathcal{L}_{pose} .

Joint Name	Weight	Joint Name	Weight
Pelvis	1.0	Neck	0.1
Left Hip	0.2	Left Collar	0.3
Right Hip	0.2	Right Collar	0.3
Spine1	0.1	Head	0.3
Left Knee	0.3	Left Shoulder	0.2
Right Knee	0.3	Right Shoulder	0.2
Spine2	0.1	Left Elbow	0.3
Left Ankle	0.3	Right Elbow	0.3
Right Ankle	0.3	Left Wrist	0.4
Spine3	0.1	Right Wrist	0.4
Left Foot	0.3	Left Hand	0.4
Right Foot	0.3	Right Hand	0.4

each sequence.

C.3. Network Architecture

The teacher and student models share the same network architecture. The visual feature encoder is implemented using a ResNet backbone pre-trained on ImageNet, producing a 512-dimensional feature vector. The IMU feature encoder is an MLP with layer sizes (60, 256, 128). The resulting latent features are then concatenated with the head pose estimated from the SLAM module. Finally, the fused representation is passed through a bidirectional LSTM, which outputs the predicted motion representation \mathcal{M} .

C.4. Training Details

We train both teacher and student models for 500 epochs with batch size 256 on an NVIDIA RTX 4080 GPU using the Adam optimizer with window size $N = 50$ frames.

Teacher Training We optimize \mathcal{L}_{tea} with learning rate 1×10^{-3} . For \mathcal{L}_{pose} , we employ per-joint weights w_j (Eq. (2)) to emphasize challenging joints: the global pelvis rotation and limb joint rotations, which are often difficult to estimate accurately from egocentric observations, the loss weights are shown in Tab. A1.

Student Training We initialize the student with the teacher’s pre-trained weights, except for the student IMU encoder, which is randomly initialized. The MLP-based IMU encoder is trained with a learning rate 1×10^{-3} to learn representations that compensate for missing dense IMU observations. The remaining network parameters are fine-tuned at a lower learning rate of 1×10^{-4} to preserve their learned capacity for mapping multi-modal features to motion predictions.

Student Loss Weighting We initialize the loss coefficients in \mathcal{L}_{stu} (Eq. (4)) as: $\lambda_{motion} = 1.0$, $\lambda_{output} = 0.5$, and $\lambda_{feat} = 0.5$. The latter two are decayed by a factor of 0.8 every 10 epochs to gradually shift emphasis toward direct motion prediction.