

Fish2Mesh Transformer: 3D Human Mesh Recovery from Egocentric Vision

Tianma Shen¹ Aditya Puranik¹ James Vong¹ Vrushabh Deogirikar¹ Ryan Fell¹
 Julianna Dietrich¹ Maria Kyrarini¹ Christopher Kitts¹
 David C. Jeong¹

¹Santa Clara University, 500 El Camino Real, Santa Clara, CA 95053

{tshen2, apuranik, jvong, vdeogirikar, rfell, jdietrich, mkyrarini, ckitts, dcjeong}@scu.edu

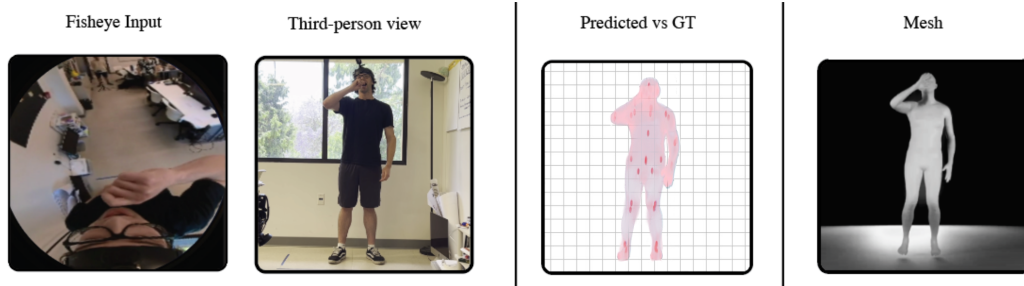


Figure 1. The Fish2Mesh pipeline enables accurate 3D mesh recovery from egocentric fisheye perspectives. From left to right: (1) Fisheye Input with a wide field of view; (2) Third-person view (shown for context, not used as input); (3) Predicted (blue) vs. Ground Truth (red) vertices demonstrating near-complete overlap; and (4) Reconstructed 3D human mesh model from the fisheye input.

Abstract

*Egocentric human body estimation allows for the inference of user body pose and shape from a wearable camera’s first-person perspective. Although pose estimation techniques have been used to overcome self-occlusions and image distortions caused by head-mounted fisheye images, similar advances in 3D human mesh recovery (HMR) techniques have been limited. We address this gap with **Fish2Mesh**, a fisheye-aware transformer-based model designed for 3D egocentric human mesh recovery. We propose an egocentric position embedding block to generate an ego-specific position table for the Swin Transformer to reduce fisheye image distortion. Our model utilizes multi-task heads for SMPL parametric regression and camera translations, estimating 3D and 2D joints as auxiliary loss to support model training. Further, we augment egocentric camera data with a training dataset by employing the pre-trained 4D-Human model and third-person cameras for weak supervision. Our experiments demonstrate that Fish2Mesh outperforms state-of-the-art 3D HMR models. Code and data are available on our [website](#).*

1. Introduction

Egocentric human body estimation is an emerging subfield of computer vision [2, 6, 7] that focuses on estimating the

human body’s pose and shape from a first-person point of view through wearable cameras [5, 16, 17, 25, 28], typically mounted on the head [28]. Egocentric estimation uses captured images from this unique perspective, enabling applications in assistive robotics [12, 20, 30], augmented reality [3, 10, 27], and virtual reality [21].

While forward-facing cameras [9, 13] excel at egocentric views ideal for third-person reconstruction but are limited for head- and arm-level scenes for first-person view, downward-facing cameras [17, 28] capture a broader, ground-oriented perspective of the human body and surroundings. This downward-facing perspective enriches contextual detail and boosts the processing potential. Building on these advantages, our work employs head-mounted cameras that capture downward views for human body estimation through key joint mapping, enabling accurate modeling of the body’s structure and movements. These benefits motivate the exploration of advanced modeling techniques, offering promising avenues for capturing precise body structures from egocentric views.

Recent research in human body estimation has focused on pose estimation, which involves predicting the positions of key body joints [26]. In contrast, our current work employs **human mesh recovery** (HMR), or human mesh reconstruction, to enhance spatial precision by mapping egocentric image data onto humanoid 3D mesh models. In addition to meeting the precision demands of realistic simu-

lations and capturing body shape and volume, HMR is better suited to address inconsistencies arising from varying keypoint standards across pose estimation datasets [16, 25, 28]. For instance, datasets such as COCO 2017 (32 keypoints) [14] and Human3.6M (17 keypoints) [8] use different sets of key points to define human poses.

Egocentric HMR faces several unique challenges. First, diverse and annotated egocentric datasets are scarce, as capturing such data requires wearable cameras [17]. Second, head-mounted cameras often use fisheye lenses [5, 16, 17, 25, 28], which introduce distortions, particularly near the edges of the frame. The amount of distortion is dependent on the unique fisheye lens used [7], which also adds inconsistent standards between datasets and complicates the model training of depth estimation needed for accurate body meshes. Third, the egocentric perspective introduces self-occlusions [32], where farther body parts such as legs can often be covered up. For example, hands and arms frequently block the view of the torso, while the head naturally limits the visibility of the lower body. Additionally, since egocentric cameras are typically worn on the head, they have a restricted field of view, often resulting in partial or complete exclusion of body parts from the frame. Legs, feet, and even parts of the arms may be out of view, especially during dynamic activities such as walking, running, or interacting with objects. Together, these challenges underscore the need for innovative approaches in egocentric HMR to effectively manage limited data, occlusions, head movements, and varying lens distortions.

In order to address these challenges, we propose a fundamentally different approach in Fish2Mesh, a transformer-based architecture that advances the handling of fisheye distortions in egocentric vision through principled geometric understanding, addressing key limitations in existing methods. Rather than treating these distortions as aberrations that must be corrected, our architecture implements a parameterized **Egocentric Position Embedding** (EPE), which employs a learnable table that encodes equirectangular geometric information directly. The values in this table represent relative 3D spherical information which is then embedded in the 2D distorted fisheye images.

While standard transformer position embeddings typically use fixed sinusoidal functions or learnable embeddings that assume a regular grid [18] and thus do not account for the non-linear distortions introduced by egocentric fisheye lens images, our EPE is designed to capture these subtle, non-linear distortions, thereby facilitating a more accurate mapping between 2D distorted inputs and their corresponding 3D representations. Although previous literature on position embeddings does not directly address this specific challenge, our approach is motivated by the need to bridge this gap and is supported by our empirical improvements in egocentric human mesh recovery.

In addition to our novel position embedding, we introduce a multi-headed architecture that jointly optimizes SMPL parameters and camera transformations. By implementing auxiliary 3D-to-2D consistency objectives alongside parameter regression, we create a geometric framework that maintains spatial coherence across both 2D and 3D dimensions, which is critical for egocentric perspectives.

To address the limitations of existing datasets, we augment an existing dataset by introducing a prompt-based collection system that captures natural human motion. Moving beyond scripted actions, our approach generates a diverse range of movements that include realistic head motions and self-occlusions that are characteristic of daily activities. This improved data foundation, combined with our geometric-aware architecture, enables robust performance in challenging real-world scenarios.

In this paper, we propose the following novel contributions in the pursuit of egocentric HMR:

- **Fish2Mesh:** A fisheye-aware, transformer-based architecture that takes egocentric fisheye image inputs and reconstructs 3D human meshes over more conventional joint estimations.
- **Egocentric Position Embedding:** A novel equirectangular projection method with learned position embedding to estimate the 3D spherical coordinates of each fisheye image pixel to address image distortion.
- **Dataset augmentation:** A prompt-based data collection method that captures naturalistic and real-world human movements with real-world occlusions, significantly improving model performance in practical applications.

2. Related Work

2.1. Egocentric Human Body Estimation

2.1.1. Egocentric Pose Estimation

Research in egocentric human body estimation has gained considerable momentum, building on early pose estimation works that used stereo (i.e., dual) camera systems, such as EgoCap [22], EgoGlass [32], and UnrealEgo [1]. However, such stereo-egocentric vision systems were constrained by the alignment and calibration of stereo cameras, the computational overhead required to process stereo imagery, and the physical burden of additional hardware. As such, egocentric body estimation shifted towards monocular (i.e., single-camera) systems [16, 17, 24, 28], which simplify hardware complexity and reduce user burden while addressing the inherent challenges of egocentric vision systems. The estimation of monocular egocentric pose was further improved by the development of new training data sets [24, 28], such as SelfPose [24] and Mo2Cap2 [28], thus eliminating the need for stereo camera setup, reducing complexity and improving usability.

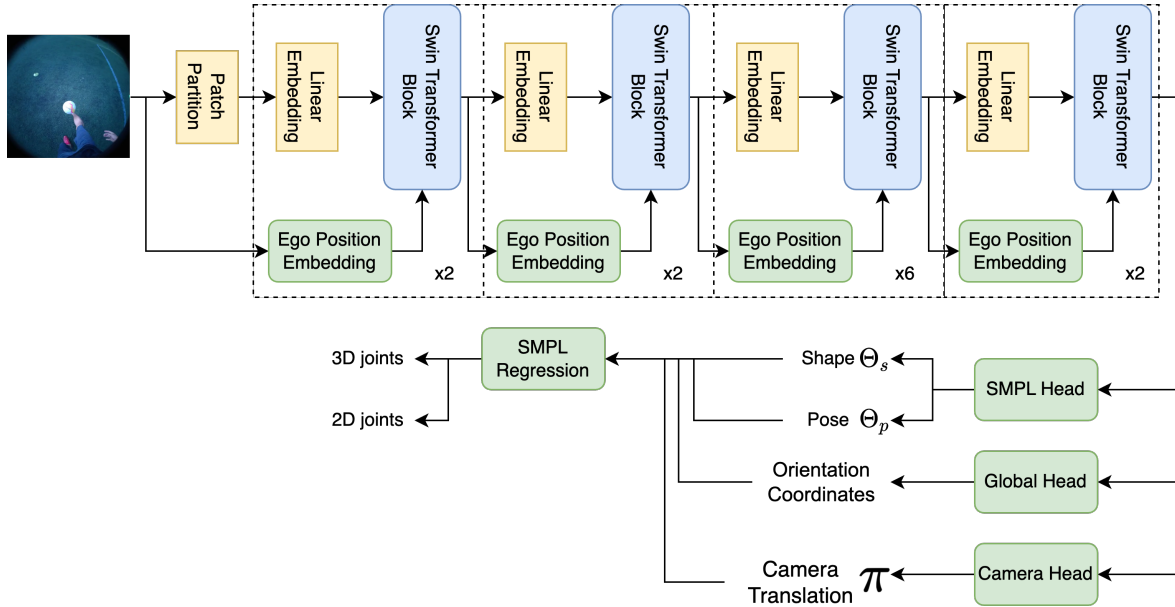


Figure 2. The architecture of the Fish2Mesh transformer model. The SMPL parameters Θ_s , Θ_p are calculated to recover the human mesh, where Θ_s and Θ_p refer to the shape parameters and pose parameters respectively.

2.1.2. Egocentric Human Mesh Recovery

While egocentric human body estimation research has focused mainly on estimating joints through pose estimation [16, 25, 28], EgoBody [31] and EgoMoCap [15] first sought to introduce the challenge of addressing human mesh recovery using original egocentric image datasets. However, EgoBody and EgoMoCap reconstructed the human mesh of a confederate interlocutor that was the subject of the egocentric image data from the first-person view of the wearer. In other words, EgoBody and EgoMoCap achieved third-person human mesh recovery of egocentric videos, which was a departure from the first-person body estimation achieved via egocentric pose estimators like EgoCap [22] and EgoGlass [32].

2.1.3. State of the Art: EgoHMR

In perhaps the most significant contribution to egocentric vision to date, EgoHMR [17] proposed a diffusion model to recover a 3D human mesh from the camera user using an egocentric video dataset with SMPL labels [19]. Despite its strengths, the inherently variable nature of EgoHMR’s diffusion model makes it unsuitable for applications requiring high precision, such as real-time XR systems [3, 10, 27] and assistive human-robot interaction [12, 20, 30]. More critically, the EgoHMR model fails to address fisheye image distortions, which can degrade mesh recovery quality. In spite of EgoHMR’s significant advances, the above limitations underscore the need for further human mesh recovery methods that can meet the rigorous demands of egocentric vision with greater resilience to fisheye images.

2.2. Egocentric Vision Transformer Models

One of the significant challenges in developing deep learning models for egocentric vision is that the camera’s close proximity to the body introduces image distortions that complicate tasks like feature extractions and pose estimation [16, 17, 24, 28].

Panoformer [23] attempts to address fisheye distortion challenges by estimating depth from panoramic images using tangent patches from the spherical domain, effectively correcting distortions. However, the local attention mechanisms used in Panoformer have a restricted receptive field, which limits their ability to capture long-range dependencies and global context effectively. Alternatively, EGformer [29] introduces equirectangular-aware, multi-head self-attention to enlarge the receptive field along the horizontal and vertical axes, but this comes at a higher computational cost and a greater number of network parameters, reducing efficiency.

More recently, FisheyeViT [26] employs a patch-based strategy, extracting many undistorted patches from fisheye images and fitting them into a transformer network. However, the FisheyeViT method relies on small patch sizes, leading to a larger number of patches and thus significant preprocessing time.

Our proposed model introduces a novel **position embedding tailored for the Swin Transformer** that allows the model to capture both horizontal and vertical information more effectively in the equirectangular domain.

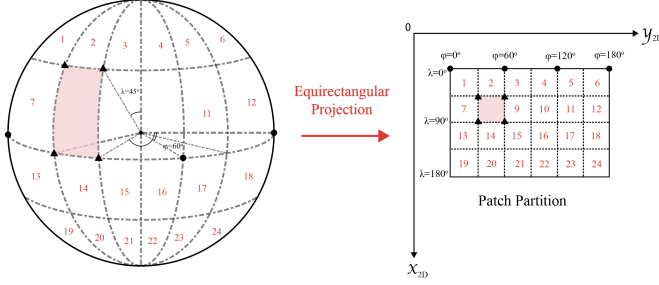


Figure 3. Equirectangular projection from spherical image.

3. Methodology

3.1. Problem Statement

The primary objective of this paper is to achieve accurate HMR from egocentric RGB images captured by a head-mounted fisheye camera. We denote the captured image at each frame as I , where each pixel is defined by its 2D coordinates x_{2D} and y_{2D} , forming the input data of our model. As illustrated in Fig. 2, our proposed *Fish2Mesh* framework consists of three components: (1) **Egocentric Position Embedding** to mitigate fisheye distortion, (2) **Swin Transformer blocks** [18] for hierarchical feature extraction, and (3) **Task-specific heads** for predicting SMPL parameters as output. The outputs of *Fish2Mesh* include the SMPL parameters Θ_s for body shape, Θ_p for body pose, camera transformation Π , global orientation L , 3D key joints L_{3D} and 2D key joints L_{2D} . The positions of 3D body joints can be directly derived from the SMPL regression model.

3.2. Egocentric Position Embedding

Unlike conventional images, fisheye images are projected onto a spherical surface, which introduces image distortions as shown in Fig. 3. Previous approaches [5, 16, 17, 25, 28] often unintentionally amplify these distortions through projection techniques, leading to a further loss in spatial accuracy. To mitigate this, we propose a parameterized Egocentric Position Embedding that applies **equirectangular geometry bias** directly to the raw image elements. This embedding method reduces distortion by adapting to the fisheye images' inherent spherical nature while avoiding computational complexity from traditional correction methods. This equirectangular geometry bias also adapts to the unique spatial constraints of egocentric data, where the camera's proximity to the wearer introduces further challenges that standard positional embeddings struggle to address. Our Vision Transformer-inspired approach [4] also introduces positional information to transformer blocks to effectively retain 3D spatial context.

3.2.1. Equirectangular Projection

As seen in Fig 3, we project the spherical 180-degree fisheye image onto a 2D plane to construct an equirectangular panoramic image using the following formula:

$$x_{2D} = R(\lambda - \lambda_0) \cos \varphi_1 \quad (1)$$

$$y_{2D} = R(\varphi - \varphi_0) \quad (2)$$

$$\lambda = \frac{x_{2D}}{R \cos \varphi_1} + \lambda_0 \quad (3)$$

$$\varphi = \frac{y_{2D}}{R} + \varphi_0 \quad (4)$$

where R is the radius of the spherical projection and depends on the physical lens properties of the camera, and x_{2D} and y_{2D} are the coordinates of input images.

φ_1 represents a reference latitude on the spherical projection. It is used to adjust the scaling factor in the transformation, ensuring that the projection accurately maps the spherical coordinates onto the 2D plane. This parameter helps maintain the correct aspect ratio and spatial alignment when converting from spherical to equirectangular representation. λ_0 is the reference longitude of the spherical coordinate system that sets the origin of the longitude measurement, serving as the reference point for horizontal positioning in the equirectangular projection. This ensures that the panoramic image is correctly centered when unwrapped onto the 2D plane. φ_0 is the reference latitude, which determines the vertical center of the projection. It specifies the latitude where the unwrapping process begins, helping to align the image correctly in the 2D coordinate system.

The variables $\lambda \in (0, \pi)$ and $\varphi \in (0, \pi)$ denote the longitude and latitude, respectively. However, since our position embedding is implemented as a learnable table that requires discrete indices, these continuous values cannot be used directly. We restrict φ to the range $(0, \pi)$ because our fisheye image captures a 180 degree field of view.

To recover the 3D spatial context, we convert the equirectangular images back to spherical coordinates:

$$\begin{cases} x_{3D} = R \cdot \sin(\varphi) \cdot \cos(\lambda) \\ y_{3D} = R \cdot \sin(\varphi) \cdot \sin(\lambda) \\ z_{3D} = R \cdot \cos(\varphi) \end{cases} \quad (5)$$

This back-projection is essential for reconstructing the 3D spatial properties of the original scene, ensuring the model retains an accurate understanding of depth and orientation.

The final transformation extends this logic:

$$\begin{cases} x_{3D} = R \cdot \sin\left(\frac{y_{2D}}{R} + \varphi_0\right) \cdot \cos\left(\frac{x_{2D}}{R \cos \varphi_1} + \lambda_0\right) \\ y_{3D} = R \cdot \sin\left(\frac{y_{2D}}{R} + \varphi_0\right) \cdot \sin\left(\frac{x_{2D}}{R \cos \varphi_1} + \lambda_0\right) \\ z_{3D} = R \cos\left(\frac{y_{2D}}{R} + \varphi_0\right) \end{cases} \quad (6)$$

The resulting x_{3D} , y_{3D} , and z_{3D} coordinates provide a rich, structured representation of the spherical distortion.

By discretizing these transformed coordinates, we create a set of indices for our learnable 3D position embeddings $POS[x_{3D}, y_{3D}, z_{3D}]$. Unlike standard position embeddings that use fixed sinusoidal functions or direct spherical coordinates, our approach embeds discretized 3D coordinates into the Swin Transformer blocks, enabling the model to capture the spherical characteristics of the fisheye input and learn complex distortions while preserving essential 3D spatial information.

3.3. Swin Transformer Blocks

Our Fish2Mesh model employs the Swin Transformer [18] architecture as its backbone, utilizing its **patch merging layers** as the network deepens, prior to the four main Swin blocks shown in Fig 2. Due to its hierarchical and adaptable structure, the Swin Transformer supports effective integration with our embedding scheme, improving spatial understanding in egocentric views. We apply our Egocentric Position Embedding $POS[x_{3D}, y_{3D}, z_{3D}]$ of spherical information to the equirectangular images, and feed them into the patch merging layers. In the initial patch merging mechanism, we hierarchically consolidate the feature representations by merging adjacent patches. Specifically, we fuse features from a 2×2 patch grid into a $4C$ vector, where C is the channel depth, then linearly map it via $W : \mathbb{R}^{4C} \rightarrow \mathbb{R}^{2C} \times (H/2, W/2)$, which halves the space and doubles the channels [18]. An embedding layer casts this into Swin tokens, priming EPE for fisheye lens’s curves.

The two successive Swin Transformer blocks consist of a window-based and a shifted window-based multi-head self-attention (MSA) module. The window-based MSA module operates within fixed, non-overlapping windows, while the shifted window-based MSA module in the subsequent block introduces connections between neighboring windows from the previous layer, enabling the model to capture both local and long-range dependencies. Each MSA module is followed by a two-layer MLP with GELU nonlinearity, and a LayerNorm (LN) is applied before each MSA and MLP. Residual connections are added after each module to enhance gradient flow and stabilize training.

We remove the relative position bias in the transformer blocks because our features already contain positional information, making the addition of relative position bias redundant. The final attention equation is shown in Formula 7, where $Q, K, V \in \mathbb{R}^{M^2 \times d}$ represent the query, key, and value matrices, respectively, d is the dimension of the query and key, and M^2 is the number of patches within a window.

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d}} \right) V \quad (7)$$

3.4. Task-Specific Heads and Losses

After extracting key features from the Swin Transformer blocks, we employ three task heads for SMPL parameters

prediction, camera parameters, and global orientation coordinates. The training of our Fish2Mesh model relies on a well structured loss function that integrates SMPL shape Θ_s losses, SMPL pose Θ_p losses, and orientation losses, along with auxiliary 3D and 2D losses shown in Formula 8.

The total loss function is defined as:

$$\mathcal{L} = a * (\mathcal{L}_{\text{SMPL}} + \mathcal{L}_{\text{orient}}) + b * \mathcal{L}_{3D} + c * \mathcal{L}_{2D} \quad (8)$$

$$\mathcal{L}_{\text{SMPL}} = \|GT_{\Theta_s} - \Theta_s\|_2^2 + \|GT_{\Theta_p} - \Theta_p\|_2^2 \quad (9)$$

$$\mathcal{L}_{\text{orient}} = \|GT_O - O\|_1 \quad (10)$$

$$\mathcal{L}_{3D} = \|GT_{3D} - X\|_1 \quad (11)$$

$$\mathcal{L}_{2D} = \|GT_{2D} - \pi(X)\|_1 \quad (12)$$

where \mathcal{L} is the total loss function used for model training and a, b , and c are weights assigned to the different loss components. $\mathcal{L}_{\text{SMPL}}$ represents the SMPL parameter loss, which includes shape loss Θ_s and pose loss Θ_p , and GT_{Θ_s} and GT_{Θ_p} are ground truth SMPL shape and pose parameters, respectively. The orientation loss $\mathcal{L}_{\text{orient}}$ ensures that the predicted orientation O matches the ground truth GT_O using the L1-based loss, \mathcal{L}_{3D} is the auxiliary 3D loss with GT_{3D} as the ground truth 3D coordinates, and \mathcal{L}_{2D} is the auxiliary 2D loss with GT_{2D} as the ground truth 2D coordinates. X denotes the predicted 3D coordinates, and $\pi(X)$ is X ’s projection onto the 2D plane.

This multi-loss approach ensures that the model optimizes 3D reconstructions and the alignment of the 2D projections with ground truth. Balancing between 3D and 2D consistency allows the model to converge more reliably, enhancing both accuracy and generalization. By carefully weighing the different loss components, our robust training process strikes a balance between optimizing the SMPL parameters, refining the global orientation, and improving overall mesh recovery performance. We train our Fish2Mesh model from scratch using loss function (8) in an end-to-end manner.

4. Experimental Results

4.1. Datasets

To comprehensively evaluate the performance of our Fish2Mesh model, we utilized three distinct datasets: Ego4D [6], ECHP [16], and an **enhanced ECHP dataset** using prompt-based collection that captures natural movements and real-world occlusions. For further details about experimental setup and enhancements using prompts, please refer to the supplementary material. Each of these datasets offers unique characteristics and contributions that enhance the diversity and robustness of our model.

4.1.1. Ego4D Dataset

The Ego4D dataset is an expansive dataset that captures a wide range of egocentric interactions. Comprising of 3,670

Table 1. Evaluation results across 3 datasets, including MPJPE, MPVPE, PA-MPJPE, and PA-MPVPE (all in mm).

Model	MPJPE	MPVPE	PA-MPJPE	PA-MPVPE	Dataset
4DHuman [15]	390.037	521.349	90.037	129.849	ECHP
Fisheye ViT [26]	4594.004	/	94.184	/	ECHP
EgoHMR [17]	84.332	99.983	64.112	79.031	ECHP
Fish2Mesh (Ours)	79.699	98.111	57.671	75.322	ECHP
4DHuman [15]	320.005	402.222	120.305	132.832	Ego4D
Fisheye ViT [26]	491.975	/	91.975	/	Ego4D
EgoHMR [17]	224.423	311.129	114.423	128.999	Ego4D
Fish2Mesh (Ours)	71.934	84.116	41.931	54.756	Ego4D
4DHuman [15]	298.233	320.944	98.613	120.304	Our
Fisheye ViT [26]	493.117	/	93.547	/	Our
EgoHMR [17]	227.552	294.484	127.55	144.184	Our
Fish2Mesh (Ours)	57.352	71.233	37.242	51.58	Our

hours of video footage from 923 participants across nine countries and 74 locations, this dataset serves as a cornerstone for research in first-person visual perception. Its richness in capturing various daily activities and diverse environments ensures that our model is trained and tested on a broad spectrum of real-world interactions, enhancing its generalization to unseen scenarios.

4.1.2. ECHP Dataset

The ECHP dataset is tailored specifically for egocentric 3D human pose estimation. It uses a GoPro camera with a fish-eye lens to capture daily human actions in varied environments. The ECHP dataset features a multi-camera setup and includes ground truth annotations, allowing for precise evaluation of pose estimation models. This dataset is invaluable for assessing the accuracy of our model in real-world scenarios, where precise 3D reconstructions are critical.

4.1.3. Expanding The ECHP Dataset

To address the scarcity of high-quality egocentric camera data, we adopt the EgoCentric-Human-Pose (ECHP) Dataset [16] for data expansion, leveraging consistent equipment and processing techniques. To improve label accuracy, we utilize the 4DHuman [15] model as a self-supervisor rather than PARE [11]. Incorporating the 4DHuman model significantly enhances data diversity, as the actions in the existing ECHP dataset are overly simplistic or artificial (e.g., walking actions are simulated by participants merely walking in place). As a result of participants’ confinement to a limited space, they cannot exhibit natural head turns or shakes that are naturalistic to everyday human activity. Ultimately, the ECHP data lacks the richness and realism seen in real-life scenarios, where head movements are more varied, leading to common issues like self-occlusion, as illustrated in Fig. 4.

To counter these limitations, we introduce a **prompt-based data collection system**. Instead of directing participants to perform rigid actions, we guide them using open-ended prompts such as “Stirring a Big Pot.” By engaging participants in more interactive and imaginative scenarios rather than constrained actions like walking, we aim to produce more natural and dynamic movements. The benefit of this approach is two fold: it helps enhance the realism of the recorded data and increases the variety of head movements, thus leading to improved model training and more accurate human pose estimations. The details of the ECHP extension are shown in the supplement material.

Additionally, to address challenges like self-occlusions, we use prompts to create situations where these issues are more likely to occur. For example, participants may be asked to perform high-frequency movements that induce self-occlusion or result in parts of the body being outside the frame. This enables us to collect data specifically designed to train and evaluate the model on how well it can overcome these problems.

4.2. Evaluation Metrics

For evaluating the performance of our Fish2Mesh Transformer, we employed PA-MPJPE (Procrustes-Aligned Mean Per Joint Position Error) and PA-MPVPE (Procrustes-Aligned Mean Per Vertex Position Error) as our primary evaluation metrics. The main reason for this choice is the evaluation of models across different datasets. Previous state-of-the-art (SOTA) models often trained on one dataset and were then evaluated on another, which may not provide a fair comparison. In real-world applications, however, it is not feasible to retrain models on different datasets. Therefore, PA-MPJPE and PA-MPVPE offer a more reliable metric for evaluating the performance of models in varying camera coordinates, as they align

the results across datasets and allow for a more accurate comparison of model performance. Nevertheless, we still provide MPJPE and MPVPE for reference to give a comprehensive understanding of the model’s performance on the traditional metrics.

4.2.1. Challenges with Traditional Metrics

Traditional metrics like MPJPE and MPVPE struggle with inherent flaws [8]: 3D poses estimated in camera coordinates vary with positioning, body orientation misaligns with the world frame, and inversion (e.g., left-right symmetry) inflates errors despite near-correct poses. These issues, exacerbated by fisheye lens distortions, compromise reliability and motivate our adoption of PA-aligned metrics to fairly assess EPE’s pose accuracy.

To address these challenges, we use PA-MPJPE and PA-MPVPE with Procrustes alignment (PA), which adjusts predicted poses for scale, rotation, and translation [17]. This minimizes external distortions like body tilt—focusing metrics on pose accuracy, enhancing evaluation of our EPE’s 2D-to-3D mapping.

4.3. Quantitative Results

The results for PA-MPJPE and PA-MPVPE across the three datasets are summarized in Table 1. Our model outperforms all state-of-the-art (SOTA) models, achieving the best results. Since FisheyeViT [26] is a pose estimation method, PA-MPVPE cannot be calculated for it, so we focus on comparing PA-MPJPE values with other SOTA models. To ensure a fair comparison, we use the official checkpoints released by the original authors for both EgoHMR and FisheyeViT, which were trained on egocentric datasets.

Notably, while EgoHMR [17] shows overfitting on its own dataset (ECHP), its performance is still 11.2% higher on PA-MPJPE and 4.9% higher on PA-MPVPE (lower values indicate better performance). Furthermore, Fish2Mesh outperforms other SOTA models on our extended dataset, demonstrating its superior ability to handle challenges such as self-occlusion and body parts outside the frame. This is because our dataset includes high-frequency data where these issues frequently occur, thus specifically addressing these two challenges. We present visual results in the next section to further illustrate the model’s performance.

Fish2Mesh produces more realistic and accurate human meshes, closely matching the ground truth. Even in challenging cases, such as occlusions caused by body parts or limbs moving out of frame, our model continues to deliver superior performance compared to SOTA models. Fig. 4 depicts one example from four datasets to compare the performance of all models (blue meshes) against the ground truth (red meshes). Notably, our Fish2Mesh demonstrates superior estimation of the lower body parts compared to other SOTA models (e.g., sitting in a chair, squatting).

Table 2. Results of ablation study across 3 datasets. Both PA-MPJPE and PA-MPVPE are measured in millimeters (mm).

Model	PA-MPJPE	PA-MPVPE	Dataset
Without EPE	92.448	106.005	ECHP
Without our dataset	71.446	90.049	ECHP
Fish2Mesh (Ours)	57.671	75.322	ECHP
Without EPE	90.111	109.995	Ego4D
Without our dataset	65.559	87.529	Ego4D
Fish2Mesh (Ours)	41.931	54.756	Ego4D
Without EPE	87.731	99.981	Our
Without our dataset	60.446	80.782	Our
Fish2Mesh (Ours)	37.242	51.58	Our

4.4. Ablation Studies

To verify the contributions of our EPE and **proposed dataset**, we conducted two ablation studies (Table 2). Each study isolates individual components to assess impact on model performance utilizing PA-MPJPE (Procrustes-Aligned Mean Per Joint Position Error) and PA-MPVPE (Procrustes-Aligned Mean Per Vertex Position Error).

4.4.1. Removing Egocentric Position Embedding

As EPE is essential for encoding spherical positional information, our first ablation study removed EPE from the Swin Transformer input of our model and replaced the position embedding to control for its impact on fisheye image distortions. With just a standard position embedding, the model relied solely on raw image features without fisheye-specific encoding nor structured positional guidance, yielding reduced spatial understanding of the input data (Table 2). This replacement resulted in a significant drop in performance, with PA-MPJPE increasing from **57.67 mm** to **92.45 mm** on the ECHP dataset. This substantial increase in error highlights that EPE enables the model to better interpret pixel information in fisheye images, ensuring accurate depth and positional estimation in such distorted contexts.

4.4.2. Excluding Proposed Dataset

Next, we trained Fish2Mesh without our proposed dataset, relying solely on the existing Ego4D and ECHP datasets. The proposed dataset contains realistic human activities, varied camera movements, and natural head motions, adding diversity and complexity that enhances the model’s generalization to real-world scenarios. When trained without our dataset, the model exhibited a measurable but smaller drop in performance than seen in the first study. Specifically, PA-MPJPE increased from **57.67 mm** to **71.44 mm** on the ECHP dataset (Table 2). This decline can be attributed to the diversity of our augmented ECHP dataset enhancing the model’s overall performance.

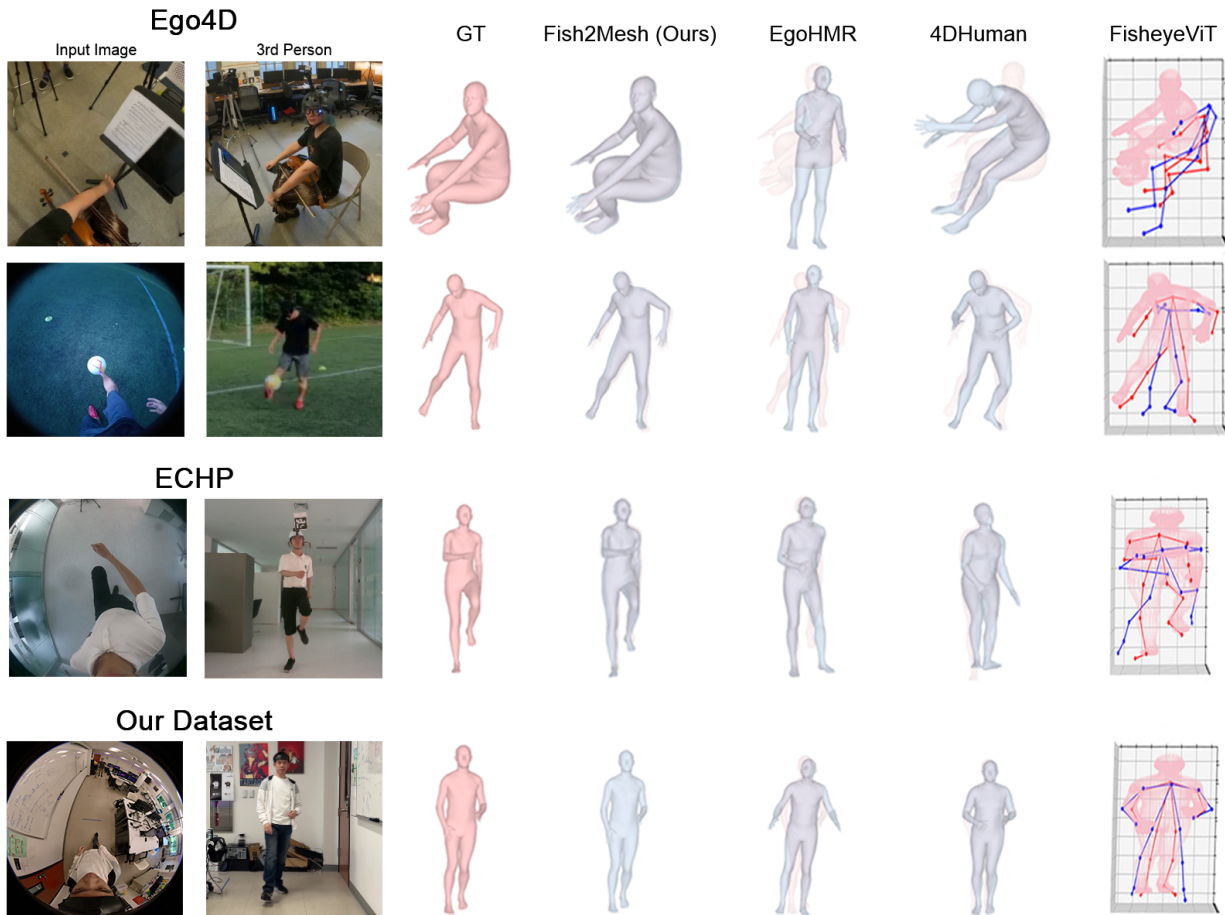


Figure 4. Visual results of four examples from the four datasets, showing ground truth (red) and related models (blue). FisheyeViT is a pose estimation model, so we visualize the skeleton to compare the resulting joints. The third-person view is not used as model input and is provided purely as an environmental reference.

5. Conclusion

In this paper, we introduce Fish2Mesh, a fisheye-aware transformer model that is uniquely tailored for egocentric 3D human mesh reconstruction. Our model addresses three primary challenges in egocentric vision: (1) limited diversity in available datasets for egocentric 3D reconstruction, (2) distortions from fisheye lenses, and (3) self-occlusions due to the egocentric perspective. By leveraging multi-task heads for SMPL parametric regression and camera transformations, Fish2Mesh effectively mitigates these challenges, producing robust, high-fidelity human mesh reconstructions from distorted egocentric input.

Key contributions include a parameterized Egocentric Position Embedding (EPE) that reduces fisheye distortions, as well as a comprehensive training dataset that integrates Ego4D, ECHP, and our newly proposed dataset. This ex-

panded dataset significantly enhances model robustness by capturing diverse scenarios and natural head movements that are common in real-world applications.

The proposed framework demonstrated superior performance over other state-of-the-art methods when evaluated on multiple datasets, including Ego4D, ECHP, and our proposed dataset. In our evaluations, Fish2Mesh achieved the lowest Procrustes-Aligned Mean Per Joint Position Error (PA-MPJPE) and Procrustes-Aligned Mean Per Vertex Position Error (PA-MPVPE) across all tested datasets. These results indicate that Fish2Mesh is not only accurate but also generalizes well across diverse egocentric scenarios, offering a reliable solution for applications requiring high-quality 3D mesh reconstruction from egocentric views. This work has potential applications in the emerging industry of XR glasses and other wearable technologies that leverage egocentric camera perspectives with fisheye lenses.

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 *European Conference on Computer Vision*, pages 1–17. Springer, 2022. 2
- [2] Hiroyasu Akada, Jian Wang, Vladislav Golyanik, and Christian Theobalt. 3d human pose perception from egocentric stereo videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 767–776, 2024. 1
- [3] Tejo Chalasani, Jan Ondrej, and Aljosa Smolic. Egocentric gesture recognition for head-mounted ar devices. In *2018 IEEE international symposium on mixed and augmented reality adjunct (ISMAR-Adjunct)*, pages 109–114. IEEE, 2018. 1, 3
- [4] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 4
- [5] Shubham Goel, Georgios Pavlakos, Jathushan Rajasegaran, Angjoo Kanazawa, and Jitendra Malik. Humans in 4d: Reconstructing and tracking humans with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14783–14794, 2023. 1, 2, 4
- [6] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012, 2022. 1, 5
- [7] Kristen Grauman, Andrew Westbury, Lorenzo Torresani, Kris Kitani, Jitendra Malik, Triantafyllos Afouras, Kumar Ashutosh, Vijay Baiyya, Siddhant Bansal, Bikram Boote, et al. Ego-exo4d: Understanding skilled human activity from first-and third-person perspectives. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19383–19400, 2024. 1, 2
- [8] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3. 6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE transactions on pattern analysis and machine intelligence*, 36(7):1325–1339, 2013. 2, 7
- [9] Hao Jiang and Vamsi Krishna Ithapu. Egocentric pose estimation from human vision span. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10986–10994. IEEE, 2021. 1
- [10] Jason Karakostas, Aikaterini Valakou, Despoina Gavgiotaki, Zinovia Stefanidi, Ioannis Pastaltzidis, Grigorios Tsipouridis, Nikolaos Kilis, Konstantinos C Apostolakis, Stavroula Ntoa, Nikolaos Dimitriou, et al. A real-time wearable ar system for egocentric vision on the edge. *Virtual Reality*, 28(1):44, 2024. 1, 3
- [11] Muhammed Kocabas, Chun-Hao P Huang, Otmar Hilliges, and Michael J Black. Pare: Part attention regressor for 3d human body estimation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11127–11137, 2021. 6
- [12] Mohammed Kutbi, Haoxiang Li, Yizhe Chang, Bo Sun, Xin Li, Changjiang Cai, Nikolaos Agadakos, Gang Hua, and Philippos Mordohai. Egocentric computer vision for hands-free robotic wheelchair navigation. *Journal of Intelligent & Robotic Systems*, 107(1):10, 2023. 1, 3
- [13] Jiaman Li, Karen Liu, and Jiajun Wu. Ego-body pose estimation via ego-head pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 17142–17151, 2023. 1
- [14] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 2
- [15] Miao Liu, Dexin Yang, Yan Zhang, Zhaopeng Cui, James M Rehg, and Siyu Tang. 4d human body capture from egocentric video via 3d scene grounding. In *Proceedings of the 2021 International Conference on 3D Vision (3DV)*, pages 930–939. IEEE, 2021. 3, 6
- [16] Yuxuan Liu, Jianxin Yang, Xiao Gu, Yijun Chen, Yao Guo, and Guang-Zhong Yang. EgoFish3d: Egocentric 3d pose estimation from a fisheye camera via self-supervised learning. *IEEE Transactions on Multimedia*, 25:8880–8891, 2023. 1, 2, 3, 4, 5, 6
- [17] Yuxuan Liu, Jianxin Yang, Xiao Gu, Yao Guo, and Guang-Zhong Yang. EgoHmr: Egocentric human mesh recovery via hierarchical latent diffusion model. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9807–9813. IEEE, 2023. 1, 2, 3, 4, 6, 7
- [18] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. 2, 4, 5
- [19] 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*, pages 851–866. 2023. 3
- [20] Javier Marina-Miranda and V Javier Traver. Head and eye egocentric gesture recognition for human-robot interaction using eyewear cameras. *IEEE Robotics and Automation Letters*, 7(3):7067–7074, 2022. 1, 3
- [21] Thi-Hoa-Cuc Nguyen, Jean-Christophe Nebel, and Francisco Florez-Revuelta. Recognition of activities of daily living with egocentric vision: A review. *Sensors*, 16(1):72, 2016. 1
- [22] Helge Rhodin, Christian Richardt, Dan Casas, Eldar Insafutdinov, Mohammad Shafiei, 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. 2, 3
- [23] Zhijie Shen, Chunyu Lin, Kang Liao, Lang Nie, Zishuo Zheng, and Yao Zhao. Panoformer: Panorama transformer

- for indoor 360 depth estimation. In *European Conference on Computer Vision*, pages 195–211. Springer, 2022. [3](#)
- [24] Denis Tome, Thiemo Alldieck, Patrick Peluse, Gerard Pons-Moll, Lourdes Agapito, Hernan Badino, and Fernando De la Torre. Selfpose: 3d egocentric pose estimation from a headset mounted camera. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(6):6794–6806, 2020. [2](#), [3](#)
- [25] Jian Wang, Diogo Luvizon, Weipeng Xu, Lingjie Liu, Kripasindhu Sarkar, and Christian Theobalt. Scene-aware egocentric 3d human pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13031–13040, 2023. [1](#), [2](#), [3](#), [4](#)
- [26] Jian Wang, Zhe Cao, Diogo Luvizon, Lingjie Liu, Kripasindhu Sarkar, Danhang Tang, Thabo Beeler, and Christian Theobalt. Egocentric whole-body motion capture with fisheye and diffusion-based motion refinement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 777–787, 2024. [1](#), [3](#), [6](#), [7](#)
- [27] Fang Xu, Tianyu Zhou, Hengxu You, and Jing Du. Improving indoor wayfinding with ar-enabled egocentric cues: A comparative study. *Advanced Engineering Informatics*, 59: 102265, 2024. [1](#), [3](#)
- [28] 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 graphics*, 25(5):2093–2101, 2019. [1](#), [2](#), [3](#), [4](#)
- [29] Ilwi Yun, Chanyong Shin, Hyunku Lee, Hyuk-Jae Lee, and Chae Eun Rhee. Egformer: Equirectangular geometry-biased transformer for 360 depth estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6101–6112, 2023. [3](#)
- [30] Jingzhe Zhang, Lishuo Zhuang, Yang Wang, Yameng Zhou, Yan Meng, and Gang Hua. An egocentric vision based assistive co-robot. In *2013 IEEE 13th International Conference on Rehabilitation Robotics (ICORR)*, pages 1–7. IEEE, 2013. [1](#), [3](#)
- [31] Siwei Zhang, Qianli Ma, Yan Zhang, Zhiyin Qian, Taein Kwon, Marc Pollefeys, Federica Bogo, and Siyu Tang. Ego-body: Human body shape and motion of interacting people from head-mounted devices. In *European conference on computer vision*, pages 180–200. Springer, 2022. [3](#)
- [32] Dongxu Zhao, Zhen Wei, Jisan Mahmud, and Jan-Michael Frahm. Egoglass: Egocentric-view human pose estimation from an eyeglass frame. In *2021 International Conference on 3D Vision (3DV)*, pages 32–41. IEEE, 2021. [2](#), [3](#)