

Egocentric Whole-Body Motion Capture with FisheyeViT and Diffusion-Based Motion Refinement

Jian Wang^{1,4} Zhe Cao² Diogo Luvizon^{1,4} Lingjie Liu³
 Kripasindhu Sarkar² Danhang Tang² Thabo Beeler² Christian Theobalt^{1,4}
¹MPI Informatics & Saarland Informatics Campus ²Google ³University of Pennsylvania
⁴Saarbrücken Research Center for Visual Computing, Interaction and Artificial Intelligence

Project Page: <https://jianwang-mpi.github.io/egowholemocap>

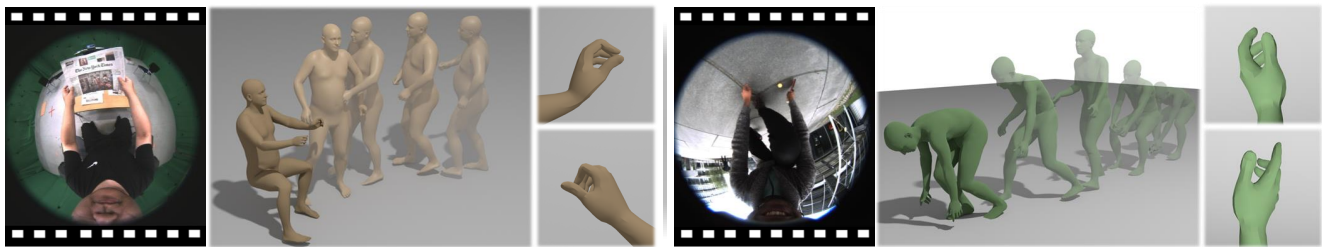


Figure 1. From an image sequence captured by a single head-mounted fisheye camera, our method can predict accurate and temporally coherent whole-body motion, including human body and hand poses. The SMPL-X parameters are obtained using inverse kinematics.

Abstract

In this work, we explore egocentric whole-body motion capture using a single fisheye camera, which simultaneously estimates human body and hand motion. This task presents significant challenges due to three factors: the lack of high-quality datasets, fisheye camera distortion, and human body self-occlusion. To address these challenges, we propose a novel approach that leverages FisheyeViT to extract fisheye image features, which are subsequently converted into pixel-aligned 3D heatmap representations for 3D human body pose prediction. For hand tracking, we incorporate dedicated hand detection and hand pose estimation networks for regressing 3D hand poses. Finally, we develop a diffusion-based whole-body motion prior model to refine the estimated whole-body motion while accounting for joint uncertainties. To train these networks, we collect a large synthetic dataset, EgoWholeBody, comprising 840,000 high-quality egocentric images captured across a diverse range of whole-body motion sequences. Quantitative and qualitative evaluations demonstrate the effectiveness of our method in producing high-quality whole-body motion estimates from a single egocentric camera.

1. Introduction

Egocentric 3D human motion estimation using head-mounted devices [47, 54] has garnered significant traction

in recent years, driven by its diverse applications in VR/AR. Immersed in a virtual world, we can traverse virtual environments, interact with virtual objects, and even simulate real-world interactions. To fully capture the intricacies of human motion during such interaction, understanding both body and hand movements is essential. While existing egocentric motion capture methods [30, 47, 50–52, 54] focus solely on body motion, neglecting the hands, this work proposes the task of egocentric *whole-body* motion capture, *i.e.* simultaneous estimation of the body motion and hand motion from a single head-mounted fisheye camera (shown in Fig. 1). This task is extremely challenging due to three factors: First, the fisheye image introduces significant distortion, making it difficult for existing networks, which are designed for non-distorted images, to extract features. Second, the egocentric camera perspective frequently leads to the occlusion of body parts, such as the feet and hands, further complicating the task of whole-body motion capture. Lastly, large-scale training data with ground truth annotations for both body and hand poses is absent in existing datasets [4, 29, 47, 51, 54].

In this work, we propose a novel egocentric whole-body motion capture method to address the aforementioned challenges. To effectively address fisheye distortion, we propose *FisheyeViT* for extracting image features, along with a joint regressor employing *pixel-aligned 3D heatmap* for predicting 3D body poses. Instead of attempting to undi-

tort the entire fisheye image, which is impractical due to the fisheye lens’s large field of view (FOV), we opt to partition the image into smaller patches aligned with a specific FOV range. This approach enables individual patch-level undistortion and seamlessly aligns with the vision transformer architecture that is employed for extracting the complete image feature map. We further propose an egocentric 3D pose regressor utilizing 3D heatmap representations. Unlike the existing approach [52] that projects image features into 3D space through fisheye reprojection functions and regresses 3D heatmaps with V2V networks [33]—leading to intricate network learning and high computational complexity—our proposed egocentric pose regressor adopts a simpler approach. It employs deconvolutional layers to obtain pixel-aligned 3D heatmaps. Notably, the voxels in the 3D heatmap directly correspond to pixels in 2D features, subsequently linking to image patches in Fisheye-ViT. This streamlined approach significantly simplifies network training. Joint locations from the pixel-aligned 3D heatmap are finally transformed with the fisheye camera model to obtain the 3D human body poses. Due to the large size difference between body and hands, we train a hand detection network and a hand pose estimation network to accurately regress 3D hand poses.

To overcome the challenges posed by self-occlusion and improve the accuracy of pose estimation, we propose a novel method for refining the whole-body motion predictions by incorporating temporal context and a motion prior. Our method learns a whole-body motion prior with the diffusion model [18] from a collection of diverse human motion sequences, capturing intrinsic correlations between hand and body movements. Following this, we extract the joint uncertainties from the pixel-aligned 3D heatmap and utilize them to guide the refinement of the whole-body motion. The joint uncertainties act as indicators of the trustworthiness of the pose regressor’s predictions. By conditioning on joints with low uncertainty, our whole-body motion diffusion model selectively refines joints with high uncertainty. This strategy substantially improves the quality of whole-body pose estimations and effectively mitigates the effects of self-occlusion.

In response to the absence of the egocentric whole-body motion capture datasets, we present *EgoWholeBody*, a new large-scale high-quality synthetic dataset. This dataset encompasses a wide range of whole-body motions, comprising over 870k frames, which significantly surpasses the size of previous egocentric training datasets. *EgoWholeBody* could serve as a valuable resource for advancing research in egocentric whole-body motion capture.

A thorough evaluation across a range of datasets, including SceneEgo [52], GlobalEgoMocap [50] and Mo²Cap² [54], has demonstrated the remarkable improvements of our method in estimating egocentric whole-body

motion compared to previous approaches. This substantiates the effectiveness of our approach in addressing the special challenges encountered in egocentric views, including the fisheye distortion and self-occlusion.

In summary, our key contributions are the following:

- The first egocentric whole-body motion capture method that predicts accurate and temporarily coherent egocentric body and hand motion;
- FisheyeViT for alleviating fisheye camera distortion and pose regressor using pixel-aligned 3D heatmaps for accurate egocentric body pose estimation from a single image;
- Uncertainty-aware refinement method based on motion diffusion models for correcting initial pose estimations and predicting plausible motions even under occlusion;
- *EgoWholeBody*, a new high-quality synthetic dataset for egocentric whole-body motion capture.

2. Related Work

Egocentric 3D Human Body Pose Estimation. Recently, there has been growing interest in estimating egocentric 3D poses from body-worn cameras. Some methods [21, 25, 31, 35, 59, 60] use front-facing cameras and infer the human body motion from the camera view. However, since the user’s body is often unobserved by the camera, these methods fail when the human body is not roaming around. Millerdurai *et al.* [32] leverage event cameras for estimating egocentric body pose. Other methods [4, 5, 7, 23, 39, 65] use head-mounted down-facing stereo cameras to estimate body poses. However, stereo camera setups introduce extra burdens of weight and energy consumption.

Xu *et al.* [54] and Tome *et al.* [47] introduce the single head-mounted down-facing fisheye camera setup for the egocentric 3D human pose estimation task. Zhang *et al.* [64] regressed fisheye camera parameters and 3D human pose simultaneously. To address the self-occlusion issue, Park *et al.* [36] leveraged the temporal information with the spatio-temporal self-attention network, and Liu *et al.* [30] introduced diffusion model to generate 3D human pose conditioned on egocentric image features. Wang *et al.* [50] and Liu *et al.* [28] combined the SLAM and egocentric pose estimation methods to estimate human body poses in the world coordinate. Wang *et al.* [51] and Liu *et al.* [29] leverage the synchronized egocentric camera and external cameras to collect large-scale egocentric pose estimation datasets with pseudo-ground truth. Considering the human-scene interaction, Wang *et al.* [52] estimated the scene geometry from the egocentric camera and constrained the 3D human pose with it.

These methods only focus on estimating human body poses while omitting the hand motion, and they still suffer from fisheye camera distortion since they directly put the highly distorted fisheye images into the neural network. Our proposed method can capture whole-body motion and

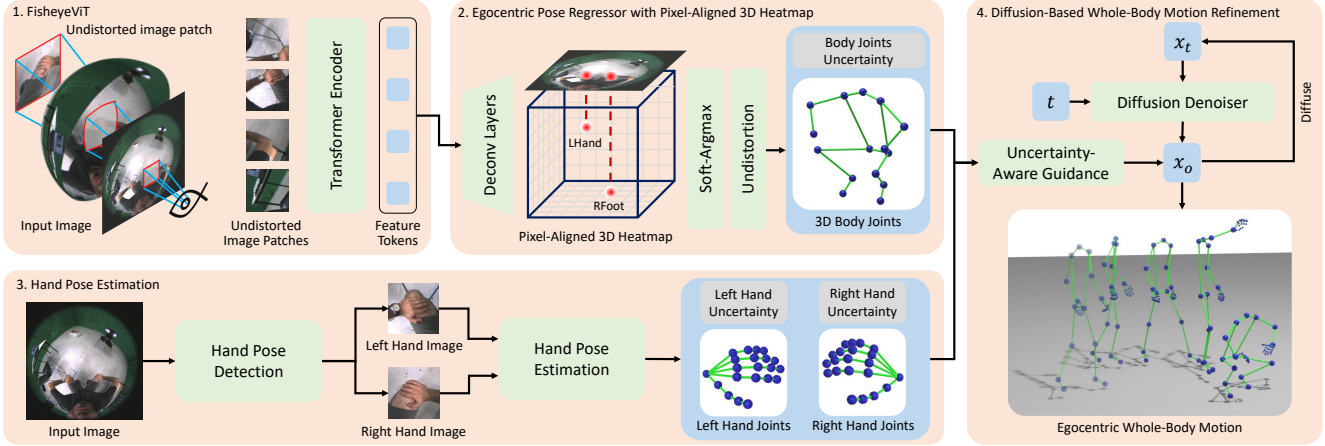


Figure 2. Overview of our whole-body motion capture pipeline. We first use FisheyeViT to undistort the input image and generate image feature tokens (3.1.1). Next, we use a 1D convolutional network to convert the image features to a pixel-aligned 3D heatmap and use soft-argmax and fisheye camera undistortion function to obtain the 3D body joints positions and uncertainty (3.1.2). We further detect the hand location and regress the 3D hand poses from the input image (3.1.3). Finally, the estimated hand motion and human body motion are combined and the uncertainty-aware diffusion model is applied to refine the estimated whole-body motion (3.2).

resolve the fisheye camera distortion issue with the FisheyeViT and pixel-aligned 3D heatmap.

Whole-Body 3D Pose Estimation. Whole-body 3D pose estimation aims to estimate the 3D human body, face, and hands parameters from input images. This task is crucial for many applications, e.g., modeling human activities and human-scene interactions. Some methods [37, 53] fit the 2D body joints estimated from images with optimization algorithms, while these methods suffer from high computation overhead and can fall into local optima. Some other learning-based methods [6, 9, 15, 27, 40, 45, 66] use the neural network to regress the SMPL-X [37] parameters from input images. For example, ExPose [9] introduced body-driven attention to extract face and hand crops and used a refinement module to regress whole-body pose. OSX [27] proposed a one-stage pipeline for whole-body mesh recovery without separate networks for each part. SMPLer-X [6] propose a foundation model for whole-body pose estimation trained with the large model and big data.

Though much progress has been made on whole-body pose estimation from an external view, the task from an egocentric view is still unexplored. In this paper, we introduce the first whole-body 3D pose estimation method from a single egocentric image and also incorporate temporal information with diffusion-based motion refinement.

Diffusion Models for Pose Estimation. Recently, some methods [8, 10, 16, 17, 19, 42] have effectively applied Denoising Diffusion Probabilistic Models (DDPM) [18] to human pose estimation tasks. Building on the success of motion diffusion models in human pose estimation, many methods have extended this approach to egocentric pose estimation, where the human body is only partially visible from RGB cameras or VR sensors. Zhang *et al.*'s work [63]

uses a diffusion model to generate realistic human poses considering scene geometry. AGROL [14] generates body motion based on head and hand 6D pose estimates from a VR headset. EgoEgo [25] estimates head poses from a head-mounted front-facing camera and uses them to generate body poses. EgoHMR [30] extracts image features and uses them as a condition for the diffusion denoising process.

However, the aforementioned pose estimation methods train the *conditioned* diffusion model with image features or IMU signals. This cannot be generalized since the trained network only accepts one specific condition format and is inclined to learn domain-specific distributions of condition features. ZeDO [22] tackles this issue with a zero-shot diffusion-based optimization approach that doesn't require training with 2D-3D or image-3D pairs. Our method leverages the uncertainty value given by the single-frame pose estimation network and refines the initial motion estimation with the uncertainty of each joint. Moreover, different from previous methods that only focus on human body motion, we train a whole-body motion diffusion model to construct the relationship between hand and body motion.

3. Method

In this section, we propose a new method for predicting accurate egocentric whole-body poses from egocentric image sequences. An overview of our approach is shown in Fig. 2.

3.1. Single Image Based Egocentric Pose Estimation

3.1.1 FisheyeViT

In this section, we introduce FisheyeViT, which is specially designed to alleviate the fisheye distortion issue. Instead of undistorting the entire fisheye image, we extract undistorted

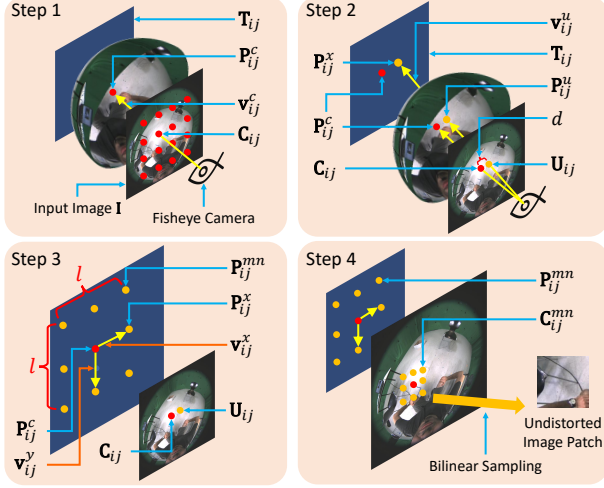


Figure 3. The detailed illustration of FisheyeViT (Sec. 3.1.1).

image patches from the fisheye image and then fit these patches as tokens into the transformer network [13]. To get the undistorted patches, we first warp the fisheye image to a unit semi-sphere, then get the patches with the gnomonic projection (see Fig. 2). The FisheyeViT can be split into five steps, the first four of which are illustrated in Fig. 3.

Step 1. Given an input image \mathbf{I} with size $H \times W$, we first evenly sample $N \times N$ patch center points: $\{\mathbf{C}_{ij} = (u_i, v_j) = (\frac{H}{N}(i + \frac{1}{2}), \frac{W}{N}(j + \frac{1}{2})) \mid i, j \in 0, \dots, N - 1\}$. Then, the patch center points \mathbf{C}_{ij} are projected onto a unit sphere with the fisheye reprojection function: $\mathbf{P}_{ij}^c = (x_{ij}^c, y_{ij}^c, z_{ij}^c) = \mathcal{P}^{-1}(u_i, v_j, 1)$. The fisheye camera model is described in Sec. 8 of the supplementary material. Given a point \mathbf{P}_{ij}^c on the unit sphere, the tangent plane \mathbf{T}_{ij} that passes through the point is defined by the normal vector $\mathbf{v}_{ij}^c = (x_{ij}^c, y_{ij}^c, z_{ij}^c)$. In the following steps, we implement the gnomonic projection by sampling grid points in the plane and projecting them back onto the fisheye image.

Step 2. In this step, we determine the orientation of the grid points in the tangent plane, ensuring that the grid points from different tangent planes \mathbf{T}_{ij} have the same orientation when projected back onto the fisheye image. To achieve this, we select a 2D point $\mathbf{U}_{ij} = (u_i + d, v_j)$ in the fisheye image space that is d pixels to the right of the patch center point and project it to the unit sphere using the fisheye reprojection function: $\mathbf{P}_{ij}^u = (x_{ij}^u, y_{ij}^u, z_{ij}^u) = \mathcal{P}^{-1}(u_i + d, v_j, 1)$. We then calculate the intersection point \mathbf{P}_{ij}^x between the vector $\mathbf{v}_{ij}^u = (x_{ij}^u, y_{ij}^u, z_{ij}^u)$ that is passing the origin and the tangent plane \mathbf{T}_{ij} :

$$\mathbf{P}_{ij}^x = \frac{\langle \mathbf{P}_{ij}^c, \mathbf{v}_{ij}^c \rangle}{\langle \mathbf{v}_{ij}^u, \mathbf{v}_{ij}^c \rangle} \mathbf{v}_{ij}^u = \frac{1}{\langle \mathbf{v}_{ij}^u, \mathbf{v}_{ij}^c \rangle} \mathbf{v}_{ij}^u, \quad (1)$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product.

Step 3. Based on the center point \mathbf{P}_{ij}^c and intersection point \mathbf{P}_{ij}^x on the tangent plane \mathbf{T}_{ij} , we build a coordinate

system with the x axis: $\mathbf{v}_{ij}^x = \text{Norm}(\mathbf{P}_{ij}^x - \mathbf{P}_{ij}^c)$, the z axis: $\mathbf{v}_{ij}^z = \text{Norm}(\mathbf{v}_{ij}^c)$ and the y axis: $\mathbf{v}_{ij}^y = \mathbf{v}_{ij}^z \times \mathbf{v}_{ij}^x$, where Norm denotes the normalize operation. We grid-sample $M \times M$ points in a $l \times l$ square on the x - y plane:

$$\{\mathbf{P}_{ij}^{mn} = \mathbf{P}_{ij}^c + (l \frac{m}{M} \mathbf{v}_{ij}^x, l \frac{n}{M} \mathbf{v}_{ij}^y)\} \quad (2)$$

where $m, n \in -\frac{1}{2}(M - 1), \dots, -\frac{3}{2}, -\frac{1}{2}, \frac{1}{2}, \frac{3}{2}, \dots, \frac{1}{2}(M - 1)$.

Step 4. The points \mathbf{P}_{ij}^{mn} are projected back to the fisheye image with the fisheye projection function: $\mathbf{C}_{ij}^{mn} = \mathcal{P}(\mathbf{P}_{ij}^{mn})$. We then apply bilinear sampling to obtain the colors at points \mathbf{C}_{ij}^{mn} of the input image \mathbf{I} , yielding the undistorted image patch $\mathbf{I}_{ij}^{\text{undis}}$. Please also see the supplementary video for a visual demonstration of undistorted image patches and their movement on the fisheye image.

Step 5. The image patches $\{\mathbf{I}_{ij}^{\text{undis}}\}$ are sent to a ViT transformer network [13] to obtain the feature tokens $\{\mathbf{F}_{ij}\}$. The feature token is further reshaped in $i \times j$ matrix and obtain the image feature \mathbf{F} . In the FisheyeViT, we empirically chose $N = 16$; $M = 16$; $d = 8$; $l = 0.2m$ given the image size $H = W = 256$.

Note that \mathbf{C}_{ij}^{mn} is independent of the image \mathbf{I} . This means that, given a fixed fisheye camera model, we can precompute \mathbf{C}_{ij}^{mn} for all combinations of m, n and i, j in advance. This significantly speeds up both the training and evaluation processes. Furthermore, the number and dimensions of image patches $\{\mathbf{I}_{ij}^{\text{undis}}\}$ match exactly with those in the traditional ViT network. This compatibility allows us to finetune existing ViT networks on our egocentric datasets. Our sampling strategy ensures that each image patch $\mathbf{I}_{ij}^{\text{undis}}$ corresponds to the same FOV range in the fisheye camera. In our ablation study in Sec. 5.3, we show that FisheyeViT enhances the performance of the pose estimation network when applied to egocentric fisheye images.

3.1.2 Pose Regressor with Pixel-Aligned 3D Heatmap

After collecting image features with FisheyeViT, we utilize a 3D heatmap-based network to estimate the body poses. The existing 3D heatmap-based pose regressors [34, 44] are designed for the weak-perspective cameras and predict the 3D heatmap in xyz space. Directly applying these regressors will result in misalignment between 3D heatmap features in xyz space and 2D image features in the fisheye image space. Therefore, we introduce a novel egocentric pose regressor that relies on the pixel-aligned 3D heatmap, tailored to address the needs of fisheye cameras. The idea is to regress the 3D heatmap in uvw space rather than traditional xyz space, where uv corresponds to the fisheye image uv space. Specifically, given a feature map $\mathbf{F} \in \mathbb{R}^{C \times N \times N}$, where C is the channel number, N is feature map height and width, we firstly use two deconvolutional layers to convert the feature map \mathbf{F} into shape

$(D_h \times J, H_h, W_h)$, and further reshape it to pixel-aligned 3D heatmap $\mathbf{H} \in \mathbb{R}^{J \times D_h \times H_h \times W_h}$, where J is the joint number and D_h, H_h, W_h is the 3D heatmap depth, height and width. The illustration of pixel-aligned 3D heatmap is shown in Fig. 2. Next, we obtain the max-value positions $\tilde{\mathbf{J}}_b = \{(u_i, v_i, d_i) \mid i \in 0, 1, 2, \dots, J\}$ from \mathbf{H} by the differentiable soft-argmax operation [44]. Here, we note that u_i and v_i correspond to the uv-coordinate of the 3D body joint projected in the fisheye image space, and d_i denotes the distance of the joint to the fisheye camera. Finally, the 3D body joints $\hat{\mathbf{J}}_b = \{(x_i, y_i, z_i) \mid i \in 0, 1, 2, \dots, J\}$ are recovered with the fisheye reprojection function: $(x_i, y_i, z_i) = \mathcal{P}^{-1}(u_i, v_i, d_i)$. The predicted body pose $\hat{\mathbf{J}}_b$ is finally compared with the ground truth body pose \mathbf{J}_b with the MSE loss. By first regressing 3D body poses in uvd space and then reprojecting it, we ensure that the 3D heatmap is pixel-aligned with the end-to-end training.

With the pixel-aligned heatmap, our proposed 3D pose regressor solves problems in all three types of previous egocentric joint regressors. First, Mo²Cap² [54] employs separate networks to predict 2D joint positions and joint distances. However, this method can yield unrealistic joint estimations because small errors in 2D joints can result in large errors in 3D joints due to the projection effect. Second, xR -egopose [47] and EgoHMR [30] directly regress the 3D joint positions. However, this method is agnostic to the fisheye camera parameters, making it suitable only for a specific camera configuration (*e.g.*, camera parameters, head-mounted position, *etc.*). Third, SceneEgo [52] projects 2D features into 3D voxel space and uses a V2V network to regress 3D poses. Because of these, the SceneEgo method suffers from low accuracy and large computation overhead. Different from previous methods, our pose regressor with pixel-aligned 3D heatmap is versatile and efficient since it directly estimates 3D joints while also incorporating an explicitly parametrized fisheye camera model. Moreover, it can preserve the uncertainty of the estimated joints, which will be used in our uncertainty-aware motion refinement method (Sec. 3.2.2). Detailed comparison with other pose prediction heads is shown in Table 3.

3.1.3 Egocentric Hand Pose Estimation

In this section, we first train a network to detect hand pose locations and then train a 3D hand pose estimation network to regress 3D hand poses. Then, we describe how to integrate the estimated hand and body poses.

Hand Detection. Given an input image \mathbf{I} , we finetune the HRNet [49] network to regress the 2D hand poses of left hand \mathbf{J}_{lh}^{2d} and right hand \mathbf{J}_{rh}^{2d} . From the hand poses, we obtain the center point of left hand \mathbf{C}_{lh} and right hand \mathbf{C}_{rh} , along with the bounding box sizes, d_{lh} and d_{rh} . We use our approach described in Sec. 3.1.1 to compute undistorted

image patches of left \mathbf{I}_{lh} and right hands \mathbf{I}_{rh} .

Hand Pose Estimation. Given the cropped image \mathbf{I}_{lh} or \mathbf{I}_{rh} , we regress the 3D hand poses $\hat{\mathbf{J}}_{lh}^{loc}$ and $\hat{\mathbf{J}}_{rh}^{loc}$ with the Hand4Whole [34] network, which is fine-tuned on our Ego-FullBody dataset.

Integration of Body and Hand Poses. It is not straightforward to integrate the hand poses with the body pose in the egocentric camera view primarily due to the fisheye camera’s perspective effects. Take the left hand as an example. Following Step 3 in Sec. 3.1.1, we establish a local coordinate system on the tangent plane of the left-hand image with XYZ axes as follows: $x : \mathbf{v}_{lh}^x; y : \mathbf{v}_{lh}^y; z : \mathbf{v}_{lh}^z$. We define a rotation matrix, denoted as \mathbf{R} , that represents the transformation between the root coordinate system and the local coordinate system on the tangent plane. The estimated hand pose is first rotated with the rotation matrix $\hat{\mathbf{J}}_{lh} = \mathbf{R}\hat{\mathbf{J}}_{lh}^{loc}$ and then translated to align the wrist location of the human body. This same process is also applied to the right hand to get the right hand pose $\hat{\mathbf{J}}_{rh}$. The whole-body joints $\hat{\mathbf{J}}$ are obtained by combining $\hat{\mathbf{J}}_b, \hat{\mathbf{J}}_{lh}$, and $\hat{\mathbf{J}}_{rh}$. The uncertainty of whole-body joints $\hat{\mathbf{U}}$ is also obtained from the maximal value of the 3D heatmap in pose estimation modules.

3.2. Diffusion-Based Motion Refinement

We notice that the single-frame estimations in Sec. 3.1 suffer from inaccuracies and temporal instabilities. In this section, we propose a diffusion-based motion refinement method to tackle this problem. We first learn the whole-body motion prior with the motion diffusion model in Sec. 3.2.1 and then introduce an uncertainty-aware zero-shot motion refinement method in Sec. 3.2.2.

3.2.1 Whole-Body Motion Diffusion Model

We follow DDPM [18] as our diffusion approach to capture the whole-body motion prior $q(\mathbf{x})$. DDPM learns a distribution of whole-body motion \mathbf{x} through a forward diffusion process and an inverse denoising process. The forward diffusion process is a Markov process of adding Gaussian noise over $t \in \{0, 1, \dots, T - 1\}$ steps:

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\sqrt{\alpha_t} \mathbf{x}_{t-1}, (1 - \alpha_t) I) \quad (3)$$

where \mathbf{x}_t denotes the whole-body motion sequence at step t , the variance $(1 - \alpha_t) \in (0, 1]$ denotes a constant hyperparameter increases with t .

The inverse process uses a denoising network $D(\cdot)$ to remove the added Gaussian noise at each time step t . Here we use the transformer-based framework in EDGE [48] as the motion-denoising network $D(\cdot)$. We follow Ramesh *et al.*’s work [38] to make the network predict the original signal itself, *i.e.* $\hat{\mathbf{x}}_0 = D(\mathbf{x}_t, t)$ and train it with the simple objective [18]:

$$\mathcal{L}_{\text{simple}} = E_{\mathbf{x}_0 \sim q(\mathbf{x}_0), t \sim [1, T]} [\|\mathbf{x}_0 - D(\mathbf{x}_t, t)\|_2^2] \quad (4)$$

3.2.2 Uncertainty-Aware Motion Refinement

Given the learned whole-body motion prior, we leverage the uncertainty value for each pose to guide the diffusion denoising process with the classifier-guided diffusion sampling [12]. Given an initial sequence of whole-body pose estimation $\mathbf{x}_e = \{\hat{\mathbf{J}}_i\}$ and the uncertainty value for each pose $\mathbf{u} = \{\hat{\mathbf{U}}_i\}$, where i denotes the i th pose in the sequence, we keep the joints with low uncertainty but use the diffusion model to generate joints with high uncertainty conditioned on the low-uncertainty joints. Specifically, in the t th sampling step of the diffusion process, the denoising network predicts $\hat{\mathbf{x}}_0 = D(\mathbf{x}_t, t)$, which is noised back to \mathbf{x}_{t-1} by sampling from the Gaussian distribution:

$$\mathbf{x}_{t-1} \sim \mathcal{N}(\hat{\mathbf{x}}_0 + \mathbf{w}(\mathbf{x}_e - \hat{\mathbf{x}}_0), \Sigma_t) \quad (5)$$

where Σ_t is a scheduled Gaussian distribution in DDPM [18] and \mathbf{w} controls the weight of a specific joint between the predicted motion $\hat{\mathbf{x}}_0$ and the estimated motion \mathbf{x}_e . Generally, we expect $\mathbf{w} \rightarrow \vec{0}$ when $t \rightarrow 0$ such that the temporal stability is guaranteed through the generation of the denoising process, and $\mathbf{w} \rightarrow \vec{1}$ when $t \rightarrow T$ such that the denoising process is initialized by the estimated motion \mathbf{x}_e . We also expect that $w_{ij} = \mathbf{w}[i][j]$, which is the weight of j th joints in the i th pose, is smaller when the uncertainty value $u_{ij} = \mathbf{u}[i][j]$ of the j th joints in the i th pose is large. Based on this requirement, we design \mathbf{w} as:

$$\mathbf{w} = 1 / \left(1 + e^{-k(t-T\mathbf{u})} \right) \quad (6)$$

where T is the overall diffusion steps, k is a hyperparameter which is empirically set to 0.1. From the experimental results in Sec. 5, we demonstrate the effectiveness of uncertain-aware motion refinement and our uncertainty-guided diffusion sampling strategy.

4. EgoWholeBody Dataset

In this section, we introduce EgoWholeBody, a large-scale high-quality synthetic dataset built for the task of egocentric whole-body motion capture. The EgoWholeBody dataset is organized into two sections. The first part, containing over 700k frames, is rendered with 14 different rigged Renderpeople [3] models driven by 2367 Mixamo [2] motion sequences. The second part focuses on hand motions and contains 170k frames with the SMPL-X model. This data is constructed from 24 different shapes and textures, driven by 262 motion sequences selected from the GRAB [46] and TCDHandMocap dataset [20]. We also created synthetic test sequences, which include 133k images rendered with 3 Renderpeople models and Mixamo motions.

During the rendering process, we first attach a virtual fisheye camera to the forehead of human body models and render the images, semantic labels, and depth map with

Blender [1]. Our dataset is larger and more diverse than previous egocentric training datasets—see Sec. 10 in the supplementary material for a detailed comparison.

5. Experiments

5.1. Datasets and Evaluation Metrics

Training Datasets. To train our body pose estimation module (Sec. 3.1.1 and Sec. 3.1.2), we use our EgoWholeBody dataset and the EgoPW dataset [50]. Additionally, the EgoWholeBody dataset is used to train the hand pose estimation module in Sec. 3.1.3. For training the whole-body diffusion model (Sec. 3.2), we utilize a combined motion capture dataset that includes EgoBody [62], Mixamo [2], TCDHandMocap dataset [20] and GRAB dataset [46].

Evaluation Datasets. In our experiment, we evaluate our methods on four datasets: the GlobalEgoMocap test datasets [50], the Mo²Cap² test dataset [54], the SceneEgo test dataset [52] and out EgoWholeBody test dataset. The details of the datasets are shown in Sec. 12 of supplementary materials. Note that evaluating whole-body poses requires accurate annotations for human hands, which is absent in real-world datasets. To resolve the issue, we request the multi-view videos of the SceneEgo test dataset [52] from the authors and use a multi-view motion capture system to obtain the hand motion. The hand pose annotations will be made publicly available.

Evaluation Metrics. To evaluate the precision of human body poses on the SceneEgo test dataset [52], we use MPJPE and PA-MPJPE. For the GlobalEgoMocap test dataset [50] and Mo²Cap² test dataset [54], where egocentric camera poses are unavailable, we evaluate PA-MPJPE and BA-MPJPE. For hand pose accuracy, we align the predicted and ground truth hand poses at the root position, followed by computing MPJPE and PA-MPJPE. Detailed explanations of these metrics are in Sec. 11 of the supplementary materials. All reported metrics are in millimeters.

5.2. Comparisons on Whole-Body Pose Estimation

For a fair comparison with existing methods focusing solely on body or hand pose, we split our evaluation into two parts, reporting results of body poses in Table 1 and hand pose in Table 2. We first compare the accuracy of the human body poses with state-of-the-art methods, including EgoPW [51] and SceneEgo [52], on EgoWholeBody and SceneEgo [52] test datasets. The comparison with more previous methods [47, 50, 54] and on more evaluation datasets [50, 54] are shown in Sec. 7 of the supplementary materials. Since our motion refinement method incorporates random Gaussian noise, we generate five samples and calculate the average MPJPE values. The standard deviation is low ($< 0.01\text{mm}$) and is discussed in Sec. 13 of supplementary materials. Results are presented in Table 1, where our single-frame re-

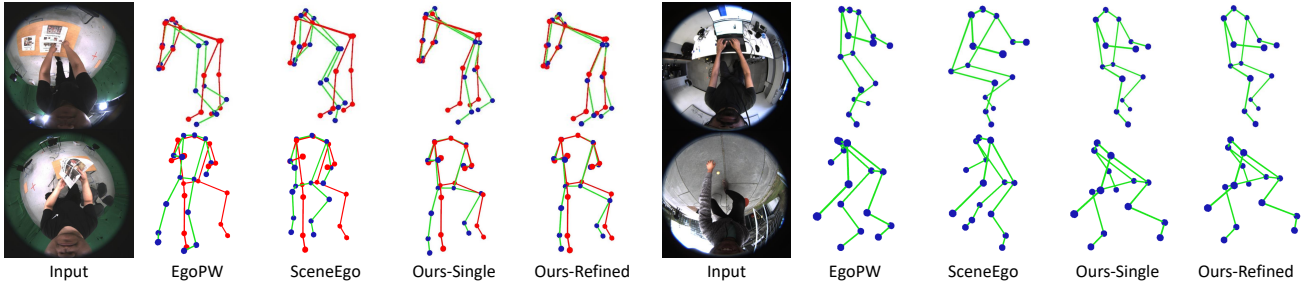


Figure 4. Qualitative comparison on human body pose estimations between our methods and the state-of-the-art egocentric pose estimation methods on in-the-studio (left column) and in-the-wild scenes (right column). The red skeleton is the ground truth while the green skeleton is the predicted pose. Our methods predict more accurate body poses compared with EgoPW [51] and SceneEgo [52].

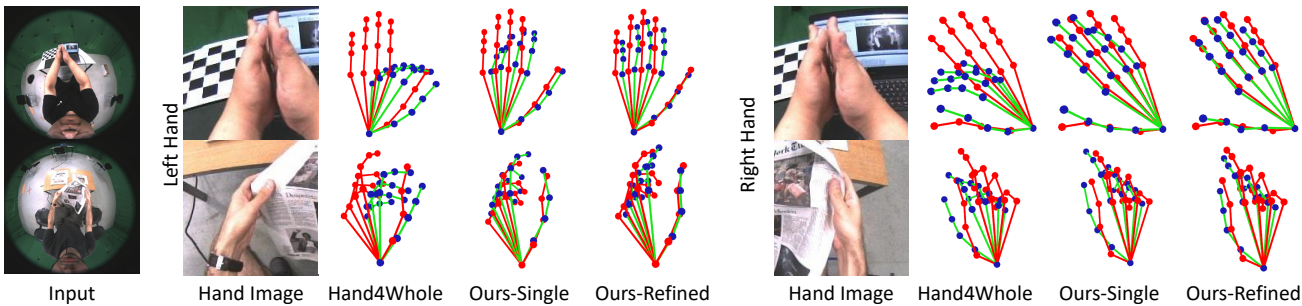


Figure 5. Qualitative comparison on human hand pose estimations between our methods and the state-of-the-art third-view pose estimation methods. Our single-view and refined hand poses are more accurate than the poses from Hand4Whole [34] method. The red skeleton is the ground truth while the green skeleton is the predicted pose.

results are labeled as “Ours-Single” and our refinement results are labeled as “Ours-Refined”. Our single-frame body pose estimation method outperforms all previous methods by a large margin. Our diffusion-based motion refinement method can further improve the accuracy of body poses estimated by the single-frame methods.

Note that previous methods [47, 50–52, 54] use training datasets different from each other. For a fair comparison, we re-train previous methods with our training datasets in Sec. 5.1 and show the results with “*” in Table 1. This retraining led to significant improvements across all previous methods, demonstrating our dataset’s broad applicability. However, these methods still underperformed compared to ours, highlighting our approach’s superiority.

To evaluate the accuracy of our hand pose estimation method, we first crop the hand images with the hand detection method in Sec. 3.1.3. Then we show the results of our single-frame hand pose estimation (labeled as “Ours-Single”) and whole-body motion refinement methods (labeled as “Ours-Refined”) in Table 2. Our single-frame hand pose estimation method outperforms the state-of-the-art method Hands4Whole [34], demonstrating the effectiveness of training the network on our EgoWholeBody dataset. Our whole-body motion refinement method can also enhance the accuracy of hand motion.

For a qualitative comparison, we compare the body and hand poses of our method with existing methods on the

Method	MPJPE	PA-MPJPE
SceneEgo test dataset [52]		
EgoPW [51]	189.6	105.3
SceneEgo [52]	118.5	92.75
EgoPW* [51]	90.96	64.33
SceneEgo* [52]	89.06	70.10
Ours-Single	<u>64.19</u>	<u>50.06</u>
Ours-Refined	57.59	46.55
EgoWholeBody test dataset		
EgoPW* [51]	84.21	63.02
SceneEgo* [52]	87.57	69.46
Ours-Single	<u>66.28</u>	<u>43.14</u>
Ours-Refined	60.32	40.35

Table 1. Egocentric human body pose accuracy of our method on SceneEgo test datasets and EgoWholeBody test dataset. Our method outperforms all previous state-of-the-art methods. * denotes the method trained with the datasets in Sec. 5.1.

SceneEgo dataset and the in-the-wild EgoPW [51] evaluation sequences. The results are shown in Fig. 4 and Fig. 5, showing that our method can predict high-quality whole-body poses from an egocentric camera. Please refer to our supplementary video for more qualitative evaluation results.

5.3. Ablation Study

EgoWholeBody Dataset. Compared to existing egocentric datasets, our EgoWholeBody dataset contains diverse body

Method	MPJPE	PA-MPJPE
SceneEgo test dataset [52]		
Hand4Whole [34]	49.66	13.85
Ours-Single	<u>23.63</u>	<u>9.59</u>
Ours-Refined	19.37	9.05
EgoWholeBody test dataset		
Hand4Whole [34]	52.85	35.04
Ours-Single	<u>33.10</u>	<u>19.68</u>
Ours-Refined	28.29	14.51

Table 2. Egocentric hand pose accuracy of our method. Our method outperforms the Hand4Whole [34] on both datasets.

Method	MPJPE	PA-MPJPE
Body Pose Results		
w/o EgoWholeBody	75.10	58.62
w/o FisheyeViT	67.36	53.44
w/ Mo ² Cap ² [54] head	87.47	65.10
w/ xR -egopose [47] head	116.5	95.78
w/ SceneEgo [52] head	77.73	62.69
Ours-Single	64.19	50.06
w/ GlobalEgoMocap [†]	69.83	56.73
w/o uncert. guidance [†]	62.16	48.40
Only body diffusion	58.95	47.03
Ours-Refined [†]	57.59	46.55
Hand Pose Results		
Only hand diffusion	21.69	9.24
Ours-Refined	19.37	9.05

Table 3. Ablation Study on SceneEgo test dataset [52]. [†] denotes the temporal-based method.

and hand motions, larger quantity of images, and higher image quality. We show this by training our body pose estimation network without our dataset, using the Mo²Cap² [54] and EgoPW [51] training dataset. The results, labeled as “w/o EgoWholeBody” in Table 3, show that performance without the EgoWholeBody dataset is inferior to our proposed method. This highlights that training with our EgoWholeBody dataset enhances the performance of the pose estimation method. We also compare this result with existing methods on the SceneEgo test set (Table 1). Trained without EgoWholeBody, our approach still outperforms previous methods, showing the effectiveness of our method.

FisheyeViT and Pose Regressor with Pixel-Aligned 3D Heatmap. To assess the individual contributions of FisheyeViT and the pixel-aligned 3D heatmap in our single-frame pose estimation pipeline, we perform experiments to measure their impact on the overall performance. First, we substitute the FisheyeViT module in our single-frame pose estimation method to ViT [13]. The result is shown in “w/o FisheyeViT” in Table 3 and it is worse than our full method. This demonstrates the effectiveness of FisheyeViT in addressing fisheye distortion and feature extraction.

Next, we analyze the performance of the single-frame pose estimation network when substituting our pose regres-

sor based on pixel-aligned 3D heatmap with the pose estimation heads of previous works [47, 52, 54]. The results of the three experiments, labeled as “w/ Mo²Cap² head”, “w/ xR -egopose head” and “w/ SceneEgo head”, show a performance drop compared to our full method. This emphasizes the crucial role of the pixel-aligned 3D heatmap in accurately estimating egocentric 3D body joint positions.

Diffusion-based Motion Refinement. We assess the effectiveness of our diffusion-based motion refinement with the following experiments: First, we compare the performance of our diffusion-based motion refinement with GlobalEgoMocap [50] by applying the GlobalEgoMocap optimizer on the single-frame body pose estimation results. The result, labeled as “w GlobalEgoMocap” in Table 3, indicates that our refinement method outperforms GlobalEgoMocap.

Second, we remove the uncertainty-aware guidance in the motion refinement. Instead, we use fixed Gaussian denoising steps to refine the motion. The result “w/o uncert. guidance” in Table 3, shows that our uncertainty-aware refinement method performs better. Our approach relies on the uncertainty values for each joint, using low-uncertainty joints to guide the generation of high-uncertainty joints. This helps reduce errors in joint predictions caused by egocentric self-occlusion, leading to improved results.

Third, we replace our whole-body motion diffusion model with the separate human body and left/right-hand diffusion models and show the accuracy of refined body and hand motion in “Only body diffusion” and “Only hand diffusion” in Table 3. From the results, we observe improvements in the accuracy of motion refined by our whole-body diffusion method, proving that learning the whole-body motion prior can help both the refinement of the body and hand motion by learning the correlation between them.

6. Conclusion

In this work, we have introduced an innovative approach to capture egocentric whole-body human motion. Our method comprises a single-frame-based whole-body pose estimation process, which includes FisheyeViT and pixel-aligned 3D heatmap representations. To enhance the initial whole-body pose estimates, we have integrated an uncertainty-aware diffusion-based motion refinement technique. Our experimental results demonstrate that both our single-frame method and the temporal-based method surpass all existing state-of-the-art techniques in terms of both quality and accuracy. Looking ahead, we see the potential for extending the applications of FisheyeViT to other vision tasks involving fisheye cameras. Future work could also involve incorporating facial expressions in whole-body motion capture.

Acknowledgments This project was supported by the Saarbrücken Research Center for Visual Computing, Interaction and AI. Christian Theobalt was supported by ERC Consolidator Grant 4DRely (770784).

References

- [1] Blender. <https://www.blender.org>. 6
- [2] Mixamo. <https://www.mixamo.com>. 6
- [3] Renderpeople. <https://renderpeople.com>. 6, 3
- [4] 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. 1, 2
- [5] Hiroyasu Akada, Jian Wang, Vladislav Golyanik, and Christian Theobalt. 3d human pose perception from egocentric stereo videos. In *37th IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, 2024. 2
- [6] Zhongang Cai, Wanqi Yin, Ailing Zeng, Chen Wei, Qingping Sun, Yanjun Wang, Hui En Pang, Haiyi Mei, Mingyuan Zhang, Lei Zhang, et al. Smpler-x: Scaling up expressive human pose and shape estimation. *arXiv preprint arXiv:2309.17448*, 2023. 3
- [7] Young-Woon Cha, True Price, Zhen Wei, Xinran Lu, Nicholas Rewkowski, Rohan Chabra, Zihe Qin, Hyoungun Kim, Zhaoqi Su, Yebin Liu, et al. Towards fully mobile 3d face, body, and environment capture using only head-worn cameras. *IEEE transactions on visualization and computer graphics*, 24(11):2993–3004, 2018. 2
- [8] Jeongjun Choi, Dongseok Shim, and H Jin Kim. Diffupose: Monocular 3d human pose estimation via denoising diffusion probabilistic model. *arXiv preprint arXiv:2212.02796*, 2022. 3
- [9] Vasileios Choutas, Georgios Pavlakos, Timo Bolkart, Dimitrios Tzionas, and Michael J Black. Monocular expressive body regression through body-driven attention. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part X 16*, pages 20–40. Springer, 2020. 3
- [10] Hai Ci, Mingdong Wu, Wentao Zhu, Xiaoxuan Ma, Hao Dong, Fangwei Zhong, and Yizhou Wang. Gfpose: Learning 3d human pose prior with gradient fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4800–4810, 2023. 3
- [11] Benjamin Coors, Alexandru Paul Condurache, and Andreas Geiger. Spherenet: Learning spherical representations for detection and classification in omnidirectional images. In *Proceedings of the European conference on computer vision (ECCV)*, pages 518–533, 2018. 5
- [12] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 6
- [13] 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, 8
- [14] Yuming Du, Robin Kips, Albert Pumarola, Sebastian Starke, Ali Thabet, and Artsiom Sanakoyeu. Avatars grow legs: Generating smooth human motion from sparse tracking inputs with diffusion model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 481–490, 2023. 3
- [15] Yao Feng, Vasileios Choutas, Timo Bolkart, Dimitrios Tzionas, and Michael J Black. Collaborative regression of expressive bodies using moderation. In *2021 International Conference on 3D Vision (3DV)*, pages 792–804. IEEE, 2021. 3
- [16] Lin Geng Foo, Jia Gong, Hossein Rahmani, and Jun Liu. Distribution-aligned diffusion for human mesh recovery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9221–9232, 2023. 3
- [17] Jia Gong, Lin Geng Foo, Zhipeng Fan, Qihong Ke, Hossein Rahmani, and Jun Liu. Diffpose: Toward more reliable 3d pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13041–13051, 2023. 3
- [18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2, 3, 5, 6
- [19] Karl Holmquist and Bastian Wandt. Diffpose: Multi-hypothesis human pose estimation using diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15977–15987, 2023. 3
- [20] Ludovic Hoyet, Kenneth Ryall, Rachel McDonnell, and Carol O’Sullivan. Sleight of hand: perception of finger motion from reduced marker sets. In *Proceedings of the ACM SIGGRAPH symposium on interactive 3D graphics and games*, pages 79–86, 2012. 6
- [21] Hao Jiang and Kristen Grauman. Seeing invisible poses: Estimating 3d body pose from egocentric video. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3501–3509. IEEE, 2017. 2
- [22] Zhongyu Jiang, Zhuoran Zhou, Lei Li, Wenhao Chai, Cheng-Yen Yang, and Jenq-Neng Hwang. Back to optimization: Diffusion-based zero-shot 3d human pose estimation. *arXiv preprint arXiv:2307.03833*, 2023. 3
- [23] Taeho Kang, Kyungjin Lee, Jinrui Zhang, and Youngki Lee. Ego3dpose: Capturing 3d cues from binocular egocentric views. *arXiv preprint arXiv:2309.11962*, 2023. 2
- [24] David G Kendall. A survey of the statistical theory of shape. *Statistical Science*, 4(2):87–99, 1989. 3
- [25] 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*, pages 17142–17151, 2023. 2, 3
- [26] Yuyan Li, Yuliang Guo, Zhixin Yan, Xinyu Huang, Ye Duan, and Liu Ren. Omnifusion: 360 monocular depth estimation via geometry-aware fusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2801–2810, 2022. 5
- [27] Jing Lin, Ailing Zeng, Haoqian Wang, Lei Zhang, and Yu Li. One-stage 3d whole-body mesh recovery with component aware transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21159–21168, 2023. 3
- [28] Yuxuan Liu, Jianxin Yang, Xiao Gu, Yao Guo, and Guang-Zhong Yang. Ego+ x: An egocentric vision system for global

- 3d human pose estimation and social interaction characterization. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5271–5277. IEEE, 2022. [2](#)
- [29] 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*, 2023. [1](#), [2](#), [4](#)
- [30] 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](#), [5](#)
- [31] Zhengyi Luo, Ryo Hachiuma, Ye Yuan, and Kris Kitani. Dynamics-regulated kinematic policy for egocentric pose estimation. *Advances in Neural Information Processing Systems*, 34:25019–25032, 2021. [2](#)
- [32] Christen Millerdurai, Hiroyasu Akada, Jian Wang, Diogo Luvizon, Christian Theobalt, and Vladislav Golyanik. Eventego3d: 3d human motion capture from egocentric event streams. In *37th IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, 2024. [2](#)
- [33] Gyeongsik Moon, Ju Yong Chang, and Kyoung Mu Lee. V2v-posenet: Voxel-to-voxel prediction network for accurate 3d hand and human pose estimation from a single depth map. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5079–5088, 2018. [2](#)
- [34] Gyeongsik Moon, Hongsuk Choi, and Kyoung Mu Lee. Accurate 3d hand pose estimation for whole-body 3d human mesh estimation. In *Computer Vision and Pattern Recognition Workshop (CVPRW)*, 2022. [4](#), [5](#), [7](#), [8](#), [2](#)
- [35] Evonne Ng, Donglai Xiang, Hanbyul Joo, and Kristen Grauman. You2me: Inferring body pose in egocentric video via first and second person interactions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9890–9900, 2020. [2](#)
- [36] Jinman Park, Kimathi Kaai, Saad Hossain, Norikatsu Sumi, Sirisha Rambhatla, and Paul Fieguth. Domain-guided spatio-temporal self-attention for egocentric 3d pose estimation. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1837–1849, 2023. [2](#), [1](#)
- [37] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and Michael J Black. Expressive body capture: 3d hands, face, and body from a single image. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10975–10985, 2019. [3](#)
- [38] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022. [5](#)
- [39] 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](#)
- [40] Yu Rong, Takaaki Shiratori, and Hanbyul Joo. Frankmocap: A monocular 3d whole-body pose estimation system via regression and integration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1749–1759, 2021. [3](#)
- [41] Davide Scaramuzza, Agostino Martinelli, and Roland Siegwart. A toolbox for easily calibrating omnidirectional cameras. In *2006 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5695–5701. IEEE, 2006. [1](#)
- [42] Wenkang Shan, Zhenhua Liu, Xinfeng Zhang, Zhao Wang, Kai Han, Shanshe Wang, Siwei Ma, and Wen Gao. Diffusion-based 3d human pose estimation with multi-hypothesis aggregation. *arXiv preprint arXiv:2303.11579*, 2023. [3](#)
- [43] Soshi Shimada, Vladislav Golyanik, Weipeng Xu, and Christian Theobalt. Physcap: Physically plausible monocular 3d motion capture in real time. *ACM Transactions on Graphics (TOG)*, 39(6):1–16, 2020. [5](#)
- [44] Xiao Sun, Bin Xiao, Fangyin Wei, Shuang Liang, and Yichen Wei. Integral human pose regression. In *Proceedings of the European conference on computer vision (ECCV)*, pages 529–545, 2018. [4](#), [5](#)
- [45] Yu Sun, Tianyu Huang, Qian Bao, Wu Liu, Wenpeng Gao, and Yili Fu. Learning monocular mesh recovery of multiple body parts via synthesis. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2669–2673. IEEE, 2022. [3](#)
- [46] Omid Taheri, Nima Ghorbani, Michael J Black, and Dimitrios Tzionas. Grab: A dataset of whole-body human grasping of objects. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16*, pages 581–600. Springer, 2020. [6](#)
- [47] Denis Tomè, Patrick Peluse, Lourdes Agapito, and Hernán Badino. xr-egopose: Egocentric 3d human pose from an HMD camera. In *ICCV*, pages 7727–7737, 2019. [1](#), [2](#), [5](#), [6](#), [7](#), [8](#), [4](#)
- [48] Jonathan Tseng, Rodrigo Castellon, and Karen Liu. Edge: Editable dance generation from music. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 448–458, 2023. [5](#), [2](#)
- [49] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, et al. Deep high-resolution representation learning for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, 43(10):3349–3364, 2020. [5](#)
- [50] Jian Wang, Lingjie Liu, Weipeng Xu, Kripasindhu Sarkar, and Christian Theobalt. Estimating egocentric 3d human pose in global space. *ICCV*, 2021. [1](#), [2](#), [6](#), [7](#), [8](#), [3](#), [4](#), [5](#)
- [51] 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 *CVPR*, pages 13157–13166, 2022. [1](#), [2](#), [6](#), [7](#), [8](#), [4](#)
- [52] 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, 5, 6, 7, 8, 3, 4
- [53] Donglai Xiang, Hanbyul Joo, and Yaser Sheikh. Monocular total capture: Posing face, body, and hands in the wild. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10965–10974, 2019. 3
- [54] Weipeng Xu, Avishek Chatterjee, Michael Zollhöfer, Helge Rhodin, Pascal Fua, Hans-Peter Seidel, and Christian Theobalt. Mo²cap²: Real-time mobile 3d motion capture with a cap-mounted fisheye camera. *IEEE Trans. Vis. Comput. Graph.*, 25(5):2093–2101, 2019. 1, 2, 5, 6, 7, 8, 3, 4
- [55] Yufei Xu, Jing Zhang, Qiming Zhang, and Dacheng Tao. Vit-pose: Simple vision transformer baselines for human pose estimation. *Advances in Neural Information Processing Systems*, 35:38571–38584, 2022. 2
- [56] Dianyi Yang, Jiadong Tang, Yu Gao, Yi Yang, and Mengyin Fu. Sector patch embedding: An embedding module conforming to the distortion pattern of fisheye image. *arXiv preprint arXiv:2303.14645*, 2023. 5
- [57] Shangrong Yang, Chunyu Lin, Kang Liao, and Yao Zhao. Dual diffusion architecture for fisheye image rectification: Synthetic-to-real generalization. *arXiv preprint arXiv:2301.11785*, 2023.
- [58] Fanghua Yu, Xintao Wang, Mingdeng Cao, Gen Li, Ying Shan, and Chao Dong. Osrt: Omnidirectional image super-resolution with distortion-aware transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13283–13292, 2023. 5
- [59] Ye Yuan and Kris Kitani. 3d ego-pose estimation via imitation learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 735–750, 2018. 2
- [60] Ye Yuan and Kris Kitani. Ego-pose estimation and forecasting as real-time pd control. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10082–10092, 2019. 2
- [61] Ye Yuan, Jiaming Song, Umar Iqbal, Arash Vahdat, and Jan Kautz. Physdiff: Physics-guided human motion diffusion model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16010–16021, 2023. 5
- [62] 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. 6
- [63] Siwei Zhang, Qianli Ma, Yan Zhang, Sadegh Aliakbarian, Darren Cosker, and Siyu Tang. Probabilistic human mesh recovery in 3d scenes from egocentric views. *arXiv preprint arXiv:2304.06024*, 2023. 3
- [64] Yahui Zhang, Shaodi You, and Theo Gevers. Automatic calibration of the fisheye camera for egocentric 3d human pose estimation from a single image. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1772–1781, 2021. 2
- [65] 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
- [66] Yuxiao Zhou, Marc Habermann, Ikhsanul Habibie, Ayush Tewari, Christian Theobalt, and Feng Xu. Monocular real-time full body capture with inter-part correlations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4811–4822, 2021. 3
- [67] Christian Zimmermann, Duygu Ceylan, Jimei Yang, Bryan Russell, Max Argus, and Thomas Brox. Freihand: A dataset for markerless capture of hand pose and shape from single rgb images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 813–822, 2019. 2