# Systematic Comparison of Projection Methods for Monocular 3D Human Pose Estimation on Fisheye Images

Stephanie Käs[1], Sven Peter[1], Henrik Thillmann[1],
Anton Burenko[1], David Benjamin Adrian[2], Dennis Mack[2], Timm Linder[2], Bastian Leibe[1]

*Abstract*— Fisheye cameras offer robots the ability to capture human movements across a wider field of view (FOV) than standard pinhole cameras, making them particularly useful for applications in human-robot interaction and automotive contexts. However, accurately detecting human poses in fisheye images is challenging due to the curved distortions inherent to fisheye optics. While various methods for undistorting fisheye images have been proposed, their effectiveness and limitations for poses that cover a wide FOV has not been systematically evaluated in the context of absolute human pose estimation from monocular fisheye images. To address this gap, we evaluate the impact of pinhole, equidistant and double sphere camera models, as well as cylindrical projection methods, on 3D human pose estimation accuracy. We find that in close-up scenarios, pinhole projection is inadequate, and the optimal projection method varies with the FOV covered by the human pose. The usage of advanced fisheye models like the double sphere model significantly enhances 3D human pose estimation accuracy.

We propose a heuristic for selecting the appropriate projection model based on the detection bounding box to enhance prediction quality.

Additionally, we introduce and evaluate on our novel dataset *FISHnCHIPS*, which features 3D human skeleton annotations in fisheye images, including images from unconventional angles, such as extreme close-ups, ground-mounted cameras, and wide-FOV poses, available at:
https://www.vision.rwth-aachen.de/fishnchips

## I. INTRODUCTION

Human pose estimation (HPE) is crucial for automotive systems [1], surveillance [2], human-robot interaction [3], action recognition [4], and sports analysis [5] [6]. Fisheye cameras, with their wide field of view (FOV), capture extensive body movement and reduce the need for multiple cameras, lowering costs in robotics and surveillance. However, fisheye lenses introduce distortions, especially towards the image boundaries, which are not present with classic pinhole (PH) cameras.

Different solutions for handling fisheye images in HPE have been explored but lack systematic comparison [7]. Some of these methods reproject fisheye images to less distorted images [8], so that HPE models trained on regular PH images can be applied [9]. This paper presents the first comprehensive evaluation of reprojection models for monocular 3D HPE with fisheye images. Our comparison includes reprojecting fisheye crops to PH format, using cylindrical barrel reprojections, applying fisheye camera models without

[1] Chair for Computer Vision, RWTH Aachen University, Germany. Mail: {lastname}@vision.rwth-aachen.de
[2] Robert Bosch GmbH, Corporate Research & Bosch Center for AI, Renningen and Hildesheim, Germany. Mail: {firstname.lastname}@de.bosch.com
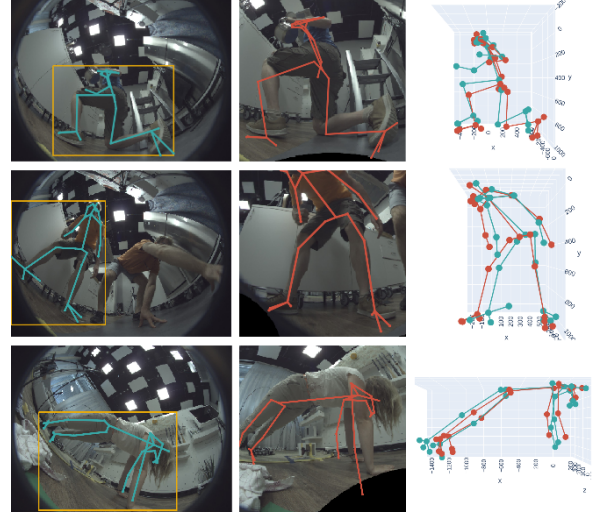
Fig. 1: Sample images from our *FISHnCHIPS* dataset, crops and predicted poses under heavy fisheye distortion.
Color legend: ■ Ground truth ■ Predictions

reprojection as well as employing bounding box heuristics for optimal choice of projection.

TABLE I: Statistics of our novel *FISHnCHIPS* dataset.

|  | Living Room 1 | Living Room 2 | Kitchen | Total |
|---|---|---|---|---|
| Subjects | 4 | 5 | 7 | 7 |
| Subjects/Scene | 1–2 | 1–3 | 1–4 | 1–4 |
| Sequences | 13 | 22 | 21 | 56 |
| Images |  |  |  |  |
| Fisheye | 50,621 | 84,863 | 66,756 | 202,240 |
| Pinhole | 51,407 | 42,532 | 44,359 | 138,298 |
| Camera Perspectives |  |  |  |  |
| Fisheye | 4 | 6 | 6 | 16 |
| Pinhole | 4 | 3 | 4 | 11 |
| Fisheye Perspectives |  |  |  |  |
| Horizontal → | ✓ | ✓ | ✓ | ✓ |
| Upwards ↑ | ✓ | ✓ |  | ✓ |
| Downwards ↓ |  |  | ✓ | ✓ |
| Tilted Downwards ↘ | ✓ | ✓ | ✓ | ✓ |
| Tilted Upwards ↖ | ✓ |  | ✓ | ✓ |

We extend the **MeTRAbs** [10] 3D HPE method to support fisheye images by integrating both equidistant (ES) and generic double sphere (DS) camera models, marking the first integration of these models into an HPE framework. Our results show these models outperform traditional PH reprojection (without requiring retraining on fisheye data) especially for subjects covering large image areas or located close to the camera.

Additionally, we introduce a **new dataset called *FISHnCHIPS*** (Fisheye Imagery in Challenging Human Poses), designed to test HPE methods under unconventional camera angles and body poses. Unlike existing fisheye HPE datasets

(Table II), which are focused on head-worn, single-subject VR/AR applications or surveillance use-cases, our dataset targets robot-like perspectives with complex multi-person scenarios and close-up interactions.

Our main contributions are: (1) a systematic comparison of five projection models, (2) an extension of a state-of-the-art pose estimator with fisheye camera models, (3) a heuristic to dynamically choose the best projection model, and (4) a novel HPE evaluation dataset with challenging fisheye camera angles and unusual poses.

## II. RELATED WORK

### A. 3D Human Pose Estimation on Pinhole Images

HPE is the process of locating the positions of key body joints of a person in an image or video. While some 3D HPE methods incorporate additional modalities like depth information [11] or Inertial Measurement Unit data [12], we focus on top-down 3D HPE using monocular RGB images. Such methods [13], [14], including MeTRAbs [10] and its extensions [15], [16], first detect persons in the image and then predict skeletons within detected bounding boxes.

3D HPE estimates joint coordinates either as absolute poses, relative to a global coordinate system, or as root-relative poses, defined with respect to a predefined root joint (e.g., pelvis) [17], [18], [10]. For 3D keypoint estimation, approaches include direct 3D regression [19], 2D heatmaps with 3D uplifting [20], and 3D heatmaps [21], [13].

### B. Fisheye-based Human Pose Estimation

Despite the prevalence of fisheye optics in autonomous driving [22], a recent survey by [7] underlines that research on 3D HPE using fisheye images is limited. Related work mainly targets egocentric pose estimation for AR/VR [23], [24], [25], [26], [27], surveillance with downward-facing cameras [8], [28], or person detection [29], [30]. For human-robot interaction, horizontal below-eye-level perspectives [9], [31] and upward-facing views of robots are particularly relevant, as robots are often smaller than humans. However, to the best of our knowledge, upward-facing views have not been addressed in any study.

*1) Without Image Reprojection:* In AR/VR scenarios with head-mounted cameras [23], [24], training HPE methods on raw fisheye images can be effective, as backbones can learn to manage distortion with a fixed camera setup. Some works use dedicated fisheye models like the omnidirectional model [32], as in [25]. [7] are the only ones to systematically evaluate several HPE and action recognition algorithms on a fisheye dataset. They find that their own approach [8], the only baseline that explicitly incorporates a fisheye-specific (polynomial) camera model, yields best HPE results on their *F-M3DHPE* dataset. It requires two separate backbones for relative and absolute pose recovery and is thus computationally more complex than our proposed method.

*2) Projection to Pinhole:* Reprojecting a fisheye image to PH format causes significant information loss [33]. However, smaller sections can be projected if the FOV is below $120°$–$140°$. [9] applied this to a small dataset with fixed

TABLE II: Comparison of ego- and exocentric HPE/HAR datasets regarding fisheye camera orientation.

| Dataset | Frames | Subj. | Perspective (fisheye only) | 🏠/☀ | Public | Real/ Synth. |
|---|---|---|---|---|---|---|
| ODIN [28] | 332K | 15 | exo: ↓ | 🏠 | ✓ | real |
| CEPDOF [29] | 25.5K | - | exo: ↓ | 🏠 | ✓ | real |
| 3DhUman [8] | 217 | 3 | exo: ↓ | 🏠 | x | real |
| OmniLab [35] | 4.8K | 5 | exo: ↓ | 🏠 | ✓ | real |
| Mo²Cap²-train [23] | 530K | 700+ | ego | 🏠/☀ | ✓ | synth |
| Mo²Cap²-eval [23] | 5.6K | - | ego | 🏠/☀ | ✓ | real |
| xR-EgoPose [24] | 383K | 46 | ego | 🏠/☀ | ✓ | synth |
| EgoCap [25] | 60K | 8 | ego | 🏠/☀ | ✓ | real |
| UnrealEgo [26] | 900K | 17 | ego | 🏠/☀ | ✓ | synth. |
| ECHP [27] | 92K | 11 | ego | 🏠/☀ | x | real |
| EgoExo4D [36] | 9.6M | 740 | ego | 🏠/☀ | ✓ | real |
| First2Third-Pose [37] | - | 14 | ego, exo: →↘ | 🏠/☀ | ✓ | real |
| Nymeria [38] | 260M | 264 | ego, exo: →↘ | 🏠/☀ | ✓ | real |
| EgoHumans [39] | 125k | - | ego, exo: →↘ | 🏠/☀ | ✓ | real |
| F-M3DHPE [7] | 2.8K | 11 | exo: →↘ | 🏠/☀ | ✉ | real |
| F-HAR [7] | - | 13 | exo: →↘ | 🏠/☀ | ✉ | real |
| *FISHnCHIPS-F* (ours) (fisheye subset, see Table I) | 202K | 7 | exo: ↓↑↖↘→ & extreme close-ups | 🏠 | ✓ | real |

Notation: indoor 🏠, outdoor ☀, on request ✉, exocentric fisheye camera orientation ↓↑↖↘→

single horizontal camera setup but did not address cases with subjects close to the camera, where PH projection becomes infeasible.

*3) Projection to Cylinder Surface:* Plaut et al. [33], [34] found that cylindrical projections better preserve translation invariance for CNNs compared to spherical projections. [31] used cylindrical projections with chest-mounted, horizontal fisheye cameras for 3D HPE.

### C. Public Datasets for Fisheye-Based HPE

Most HPE and human action recognition (HAR) datasets use PH images [40], [41], [42]. Existing fisheye datasets (Table II), mainly focus on egocentric or downward-facing surveillance views [23], [24], [25], [26], [27], [8], [29], [30], [28]. Among them, *Nymeria* [38] and *EgoExo4D* [36] use Meta's *Project Aria Glasses*, which allow a fisheye-based human-perspective observation of the scene and other humans. The datasets presented in [9] and [31] focus on third-person perspective and include horizontal views, but are not publicly accessible.

Unlike previous works, we introduce the *FISHnCHIPS* dataset, comprising 202K images from diverse camera perspectives, including horizontal, floor-mounted cameras (tilted at $45°$ or $90°$), wall-mounted cameras angled downwards, various lens types, multi-person scenarios with occlusion, and a range of activities. It includes accurate 3D pseudo-ground truth obtained through multi-view triangulation from time-synchronized cameras. We assess various fisheye reprojection methods using the state-of-the-art MeTRAbs [10] 3D HPE approach on our dataset.

## III. METRABS FISHEYE EXTENSION

The PH camera model is unsuitable for fisheye images due to its inability to account for spherical distortions that increase from center to edge. Fisheye lenses project a spherical view onto a flat image plane, causing straight lines to bend. Therefore, an HPE framework using fisheye input must incorporate camera models that represent these distortions.

We aim to assess the effectiveness of different fisheye reprojection methods on 3D HPE. To achieve this, we select the non-temporal, heatmap-based monocular HPE method MeTRAbs [10] as baseline, reimplement it in PyTorch and
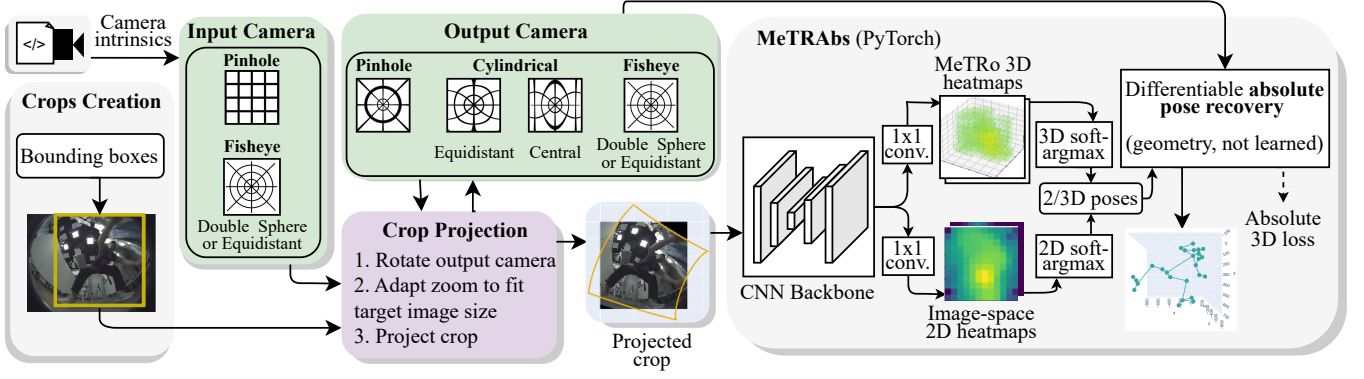
Fig. 2: Our fisheye extension of MeTRAbs. The input camera is the camera that created the original image. Its intrinsics are required input for our method. The output camera can be chosen from 5 different models. We first create crops for each person in an image, then project the crop from the input camera to the output camera (see Section III-.3). The projected crop is then fed into MeTRAbs. Using the 2D prediction, relative 3D prediction and the camera parameters of the output camera, the absolute 3D pose is recovered by solving a system of linear equations. More details in Section III-.3.

extend it to enable handling of both PH and fisheye input images (see Fig. 2). We choose this method as it enables real-time performance even with constrained computational resources due to a lightweight backbone. Our implementation incorporates different camera models and reprojection methods to effectively manage fisheye distortions. Notably, the application of these (re-)projection methods only requires camera parameters of the fisheye optics without any further need for retraining. Our approach offers seamless flexibility for various fisheye models and is adaptable to any similar top-down HPE system.

*1) Fisheye Projection Methods:* We implement two fisheye camera models: The **Equidistant Model (EF)** [43] assumes a direct proportionality between the angle from the optical axis and the radial distance from the image center. The **Double Sphere Model (DS)** [44] improves accuracy by projecting a 3D point onto two concentric spheres with shifted centers before mapping onto the image plane via a translated PH model. DS is better suited to model real fisheye lenses than the EF model. Unlike other fisheye models (e. g. polynomial) which may require an iterative approach for unprojection that is not easily usable in end-to-end training, the closed-form analytic inverse of DS facilitates obtaining normalized image coordinates for absolute pose recovery.

In contrast, a **cylindrical projection** flattens the curved surface onto a cylinder, mitigating radial distortion and preserving proportions more consistently over a $180°$ field [34]. Our framework supports both "equidistant" (EC) and "central" cylindrical (CC) projections. In CC projection [33], rays from a central point intersect the cylindrical surface aligned with the Y-axis. This projection spans $360°$ azimuthally but struggles near the cylinder axis, similar to PH camera limitations around $180°$ [34]. The EC projection maps the vertical image coordinate to the polar angle instead of the cylinder axis, allowing it to represent the entire sphere.

*2) Original MeTRAbs:* **MeTRAbs (Metric-Scale Truncation-Robust Heatmaps for Absolute 3D Human Pose Estimation)** [10], estimates an absolute 3D human pose given an image crop of a person and known input

camera parameters. It feeds the person crop through a CNN backbone and outputs 3D as well as 2D image coordinates for the person's joints. Predicted 2D keypoints are then transformed to normalized image coordinates via the unprojection equation of the crop's camera model. From those, the absolute 3D pose is obtained by solving a strong perspective model via linear least squares [10].

*3) Our Pipeline & MeTRAbs Fisheye Extension:* Fig. 2 shows how we process PH/fisheye images. Given an image of (multiple) persons, we apply an off-the-shelf person detector [45]. Backward warping is then applied to transform the resulting crop using one of the camera projections we intend to compare (see III-.1). The projection applies the intrinsics of the input camera (the physical camera that created the image) and output camera (a virtual output camera with intrinsics chosen by the user). The latter is virtually pointed at the bounding box center, creating an output image which looks as if it was actually taken with the output camera.

The output cameras's zoom factor is chosen so that the centers of the sides of the bounding box lie within the resulting image. Afterwards, we feed the transformed crop into MeTRAbs and obtain the absolute 3D pose, using the output camera's intrinsics for the absolute pose recovery.

*4) Heuristics for Projection Choice:* Our experiments (VII-B) show that for relative pose estimation, different projection types can be beneficial depending on the FOV covered by a person. We thus introduce two metrics to estimate the FOV a person occupies (MPJA) and automatically select the most appropriate projection method for each image based on the spatial expansion of the bounding box (MBBA).

(a) Analysis Tool: The ***Maximum Pairwise Joint Angle (MPJA)*** calculates the maximum angular difference for all joint pairs in a human's pose using the ground truth skeleton (Fig. 3a). MPJA helps quantifying fisheye distortions of a person by distinguishing between poses that occupy a small FOV (suitable for PH projection) from those spanning a larger FOV. It allows to reveal limitations of the PH projection model as seen in our experiments (Figure 4).

(b) Prediction Tool: Calculating MPJA requires ground

truth poses, which are unavailable during inference. To maximize HPE accuracy, we rely solely on the person's bounding box coordinates to estimate the FOV. After projecting the 2D bounding box into 3D, we calculate the angles relative to the camera center and select the maximum angle (***Maximum Bounding Box Angle, MBBA***), as shown in Fig. 3b. The optimal projection technique is selected based on MBBA during inference. We define a threshold $\alpha_t$ for MBBA and suggest to apply PH projection below the threshold and DS above it, resulting in a **hybrid projection method** (**H**). Our experiments (Table IV) demonstrate that MBBA closely approximates MPJA, yielding nearly identical inference results.



(a) Example for MPJA ($\alpha$).

(b) Example for MBBA ($\alpha$).



(c) Dist. [mm] vs. MPJA [°].
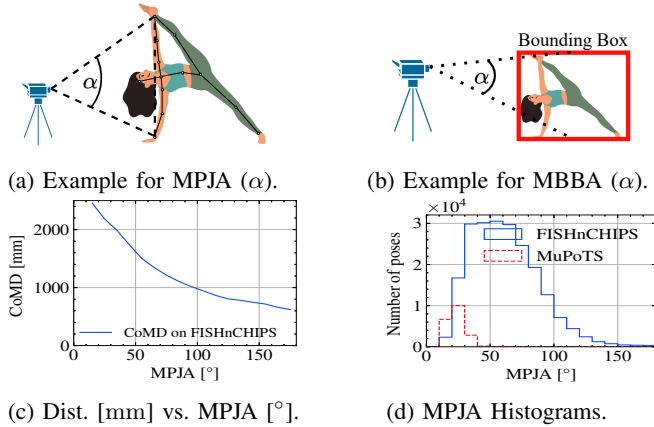
(d) MPJA Histograms.

Fig. 3: Top row: Our novel spatial expansion metrics. MPJA (a) takes the maximum pairwise angle of joints. MBBA (b) computes the maximum opening angle of the (reprojected) bounding box. Bottom row: (c) Plot depicts correlation between a skeleton's CoM-camera-distance (CoMD) and MPJA. (d) Difference in MPJA counts for MuPoTS-3D-val [42] and our *FISHnCHIPS* dataset. *(Yoga pose credit: [46])*

## IV. NOVEL FISHEYE EVALUATION DATASET

For evaluation, we introduce *FISHnCHIPS*, a novel multi-view fisheye dataset, tailored to household robotics scenarios (living room, kitchen), unlike existing datasets focused on AR/VR and surveillance applications. The dataset counts more than 200,000 fisheye images and 138,000 PH (Kinect) images, with fisheye images captured from four different camera orientations (upwards, downwards, angled, and horizontal). It also includes near-floor camera placements tilted at $45°$ and $90°$ upwards, which are particularly relevant for low-height domestic service robots. The dataset features diverse indoor environments, clothing styles, textures, and lighting conditions. Seven subjects perform activities such as embracing, yoga, unloading boxes, climbing ladders and robot navigation gestures. The distribution of MPJA values is shown in Fig. 3d as histogram and compared to the MuPoTS dataset [42]; 5500 images capture poses with MPJA$>120°$.

*1) Camera Setup:* Data was recorded using 10 synchronized cameras at 15 Hz and later downsampled to 5 Hz: 6 machine vision color cameras with fisheye lenses ($2.4$ MP and $5.1$ MP) and 4 Azure Kinect RGB-D cameras ($3.1$ MP)

with PH characteristics. We used S-mount and C-mount fisheye lenses with varying costs (€4 to €800) and distortion levels. Calibration was performed using the DS model for fisheye cameras and OpenCV's PH model for Kinects, with extrinsic calibration via commercial software (cf. [47]).

*2) Ground Truth:* We obtain precise 3D GT through multi-view triangulation across all 10 camera views. Accurate 2D bounding box detection is achieved using Co-DINO-DETR [48] with Swin-L [49] backbone, followed by RTMPose-L trained on 7 public datasets [50] for 2D pose estimation. In multi-person scenarios, we associate 2D skeletons of the same human across all views based on distances to projected, pre-recorded Kinect 3D skeletons. Triangulation is conducted using respective camera models with a non-linear least-squares solver [51], a bone symmetry constraint, and several human anatomy-based skeleton plausibility checks.

Statistics of our dataset are provided in Table I, with examples in Fig. 1. *FISHnCHIPS* features 7 subjects in 3 setups, with a total of 56 video sequences. Overall, it contains approximately 340,000 images from 16 fisheye and 11 PH camera perspectives. For more details, please refer to our supplementary video. An anonymized version of the dataset with blurred faces will be released on our project website.

## V. IMPLEMENTATION & TRAINING

For our HPE framework, we re-implemented the newest MeTRAbs version [15] which uses an autoencoder to allow training on datasets with different skeleton formats. We use the consistency-finetuning training variant presented in [15] on an EfficientNetV2-S [52] backbone, using the same model for base training and autoencoder-based finetuning. We use *Sárándi et al.*'s collection of 28 PH datasets [15] as training data and compare our predictions to the original MeTRAbs paper on *MuPoTS-3D* [42] with respect to the Percentage of Correct Keypoints (PCK), as this metric was the most common among recent methods. As shown in Table III, our re-implementation of MeTRAbs is still within the current state of the art considering that we are using a lightweight backbone and small crop resolution of $256×256$ px.

We do not require any retraining of the core MeTRAbs to use fisheye input images. However, the joint annotation schema for our novel dataset is different from any in the training data. We thus need to apply a mapping between training skeleton formats and our skeleton format. To this end, we adopt the autoencoder proposed in [15] and train a linear transformation on a separate held-out sub-dataset of *FISHnCHIPS* (21k images) recorded with different human subjects in an another room, but with a similar camera setup as described in Section IV.

## VI. METRICS & EXPERIMENTS

Building on our MeTRAbs fisheye extension, we explore best practices for HPE on fisheye images by evaluating our novel spatial expansion metrics, systematically comparing (re-)projection methods, and testing our hybrid reprojection approach.

(a) PCK$_{150}\uparrow$             (b) A-PCK$_{150}\uparrow$

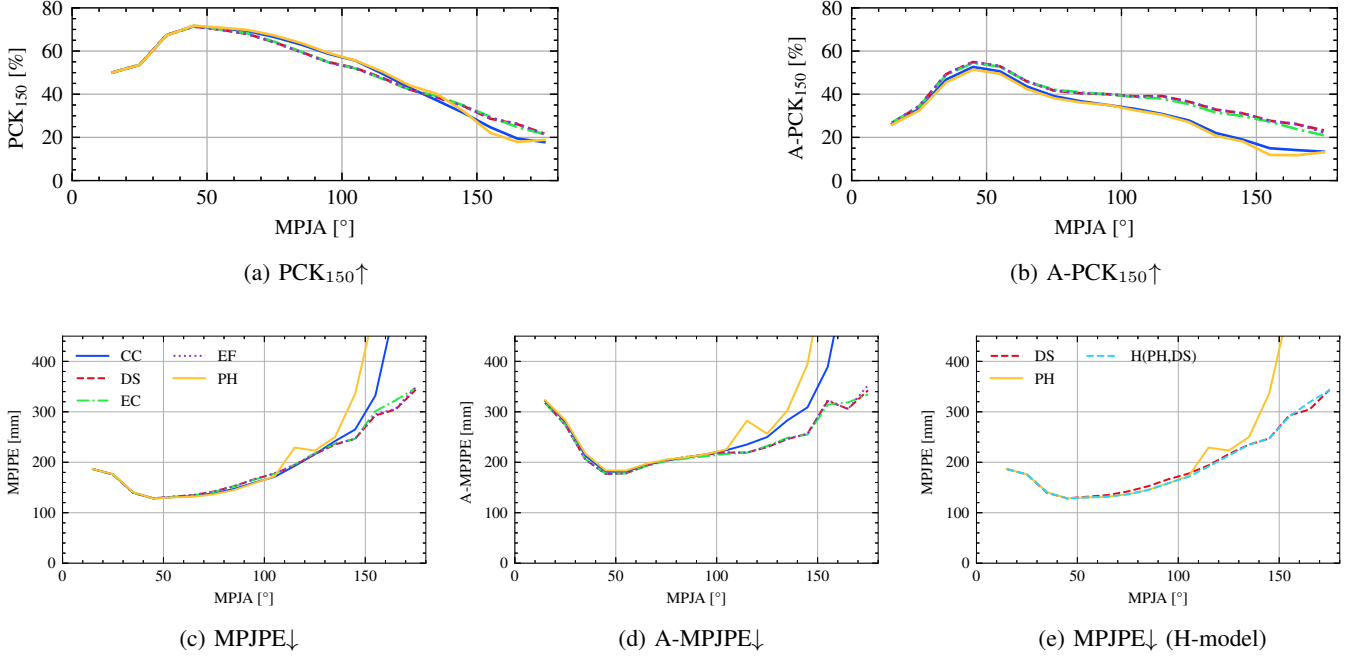(c) MPJPE$\downarrow$      (d) A-MPJPE$\downarrow$      (e) MPJPE$\downarrow$ (H-model)

Fig. 4: MPJPE/PCK and A-MPJPE/A-PCK results depending on MPJA, evaluated on our fisheye dataset *FISHnCHIPS* for different projection methods. Measured by MPJPE and PCK, Pinhole projection (PH) performs best for poses within pinhole-like FOVs ($< 120°$), whereas fisheye camera projections (DS, EF) and EC-projection are more suitable for large FOVs ($> 120°$). DS and EF perform extremely similarly. Fig. 4e depicts that the H-model combines the benefits from PH and DS: it performs identical to PH between $50°$ and $110°$ and comparable to DS for $110°$–$150°$.

TABLE III: Comparison of our MeTRAbs baseline to other recent methods on MuPoTs-3D [42] (pinhole images, no temporal information used).

| Method | Backbone | Training Data | PCK$_{150}$ [%] Abs.↑ | PCK$_{150}$ [%] Rel. ↑ |
|---|---|---|---|---|
| **Dual Network** [53] | HRNet-w32 [54] | MuCo [42] + COCO [55] | 48.1 | 89.6 |
| **PIRN** [56] | RootNet [57] | MuPoTs [42] cross-validation | 44.1 | 85.8 |
| **VirtualPose** [14] | ResNet-18/-152 [58] | MuCo [42] | 44.0 | - |
| **GR-M3D** [59] | Hourglass [60] | MuCo [42] | 41.2 | 84.6 |
| **MeTRAbs** | ResNet-50 [58] | MuCo [42] | 40.2 | 81.1 |
| **MeTRAbs** + AE | EffNetV2-L [52] | 28 Datasets | - | **95.4** |
| **MeTRAbs** + AE | EffNetV2-S [52] | 28 Datasets | - | 94.9 |
| **MeTRAbs** (our re-implementation) | EffNetV2-S [52] | 28 Datasets | **59.5** | 89.4 |
| AE = Autoencoder, | as in [15] | | | |

To investigate the reprojection method's effects on absolute and relative human pose reconstruction quality, we evaluate our results with respect to the HPE metrics **Percentage of Correct Keypoints (PCK)** and **Mean Per Joint Position Error (MPJPE)** as well as their counterparts for absolute poses (A-PCK, A-MPJPE), see [61] for details. To ensure detector-independent results, all experiments are conducted using GT bounding boxes. In real-world usage, they would be provided by a person detector.

### A. Experiments on Spatial Expansion Metrics

Fisheye images distort individuals close to the camera or with poses covering a wide FOV. To validate MPJA as a heuristic for wide FOV coverage, assumed to correlate with higher distortion, we visually assess the relationship between a person's FOV coverage and proximity to the camera. This is done by plotting MPJA against the **Center of Mass Distance to Camera (CoMD)**, which we define as the mean Euclidean distance from each joint to the camera center.

### B. Experiments on Reprojection Methods

We used our MeTRAbs extension to systematically evaluate various reprojection methods on the *FISHnCHIPS* dataset by reprojecting cropped fisheye images to:

- **PH**: Pinhole camera model.
- **EF**: Equidistant fisheye model.
- **DS**: Double sphere model.
- **CC & EC**: Cylindrical projections, distinguishing between central and equidistant cylindrical models.
- **H**: Our MBBA-threshold-based hybrid method, dynamically selecting PH or DS, based on the estimated FOV covered by the person.

We evaluate these methods by plotting PCK and MPJPE against MPJA to analyze their suitability for poses spanning different FOVs.

### VII. RESULTS & INTERPRETATION

### A. Experimental Results on Spatial Expansion Metrics

Our assumption was that subjects close to the camera (less than 1 m away) would cover a large FOV in the image. Likewise, small FOV coverage would correspond to very distant or crouched poses. This assumption is supported by our data: a plot of CoMD against the MPJA shows that low MPJA values correspond to high CoMD and vice versa (see Fig. 3c). This insight helps interpret the results of our

comparison of different fisheye reprojection methods. *(Note: All MPJA plots bin the data in $10°$-intervals.)*

### B. Experimental Results on Reprojection Methods

To compare the quality of the projection methods, we assess both the overall (A-)MPJPE and (A-)PCK on our full *FISHnCHIPS* dataset. As seen from Table IV, for absolute pose estimation, using DS as output camera slightly outperforms other methods as it has the highest A-PCK ($45.5\,\%$) and second-lowest A-MPJPE ($205.2\,\mathrm{mm}$). For relative pose estimation, no single best practice emerges.

To evaluate each projection's applicability to handle various amounts of FOV covered by the depicted person, we analyse MPJPE and PCK metrics against MPJA. Fig. 4 reveals that for relative pose estimation, reprojection method performance varies with the MPJA of the pose:

(1) **MPJA $< 50°$**: No projection method shows clear dominance. Small MPJA values indicate a small FOV covered by the pose. Thus, in this MPJA interval, we assume minimal fisheye distortion and therefore less pronounced differences among projections.

(2) $50° <$ **MPJA** $< 150°$: In this MPJA range, PH projection slightly outperforms others in relative pose estimation (Fig. 4c), as the FOV covered by the person remains within PH's effective range, minimizing pinhole reprojection artifacts. As shown in Fig. 4c, poses with MPJA between $40°$ and $70°$ yield optimal results, likely because their reprojections closely resemble the model's pinhole training data. MPJPE increases from $115\,\mathrm{mm}$ at $50°$ MPJA to over $250\,\mathrm{mm}$ at $150°$, reflecting the growing spherical distortion. For absolute pose estimation, however, PH projection does not outperform other methods in this range, as indicated by A-MPJPE and A-PCK metrics.

(3) **MPJA $> 120°$**: DS and EF projections achieve the best absolute and relative pose predictions, followed by EC. At high MPJA, PH projection's MPJPE rises rapidly, likely because the FOV exceeds PH's effective range, causing detail loss due to the small zoom factor during reprojection. In contrast, errors in other projections increase more gradually, as they are designed for larger FOVs.

As seen in Fig. 4a-d, CC projection underperforms compared to EC, DS, and EF, likely due to its limited $180°$

TABLE IV: (A-)MPJPE [mm] and (A-)PCK$_{150}$ [%] results on our *FISHnCHIPS* fisheye dataset. H stands for the proposed H projection, which uses a heuristic based upon MPJA/MBBA thresholds $\alpha_t$ to dynamically switch between PH and DS projection models. **Bold**: best, underlined: 2nd.

| Projection | MPJPE↓ | A-MPJPE↓ | PCK$_{150}$ ↑ | A-PCK$_{150}$ ↑ |
|---|---|---|---|---|
| CC | 146.5 | 209.9 | 64.8 | 42.5 |
| DS | 147.7 | <u>205.2</u> | 63.6 | **45.5** |
| EC | 147.4 | **204.4** | 63.7 | <u>45.4</u> |
| PH | 151.6 | 219.4 | <u>65.1</u> | 41.5 |
| EF | 147.8 | 205.3 | 63.5 | **45.5** |
| H w/MPJA | | | | |
| $\alpha_t = 110°$ | **145.1** | 210.4 | <u>65.1</u> | 42.0 |
| $\alpha_t = 135°$ | 145.9 | 212.1 | **65.2** | 41.7 |
| H w/MBBA | | | | |
| $\alpha_t = 110°$ | **145.1** | 210.2 | <u>65.1</u> | 42.0 |
| $\alpha_t = 135°$ | <u>145.8</u> | 211.8 | **65.2** | 41.7 |

vertical FOV (Section III-.1), restricting its effectiveness across different camera orientations. Generally, cylindrical projections, constrained by their symmetry, are less versatile than spherical models. DS and EF appear to better handle diverse camera orientations, which may explain their slight advantage over EC.

We conclude that, for large FOV poses (MPJA $> 120°$), DS, EC, and EF projections are preferred over PH and CC. For smaller FOVs, PH projection is recommended for relative pose estimation. For absolute pose estimation, using DS is optimal irrespective of the pose's FOV coverage.

### C. Results on Predicting the Best Projection Model

Section VII-B illustrates the benefits of selecting projections based on MPJA. As described in Section III-.4, we developed a hybrid method (H) that uses an MBBA threshold $\alpha_t$ to switch projections: PH for MBBA $< \alpha_t$ and DS for MBBA $> \alpha_t$. In Table IV, we exemplary depict results for $\alpha_t \in \{110°, 135°\}$, as these yielded the best MPJPE and PCK, respectively. As expected, these values align with the PH-DS curve intersections in Fig. 4a and Fig. 4c.

Additionally, we compute the heuristic using MPJA from GT joints. The H-method results in Fig. 4e confirm that the H-curve follows PH below $\alpha_t = 110°$ and DS above it. The heuristic improves MPJPE by $2\,\mathrm{mm}$ over using a single PH/DS projection, demonstrating its effectiveness for relative pose estimation by leveraging both projections' strengths. As expected, A-MPJPE drops by $6\,\mathrm{mm}$, since PH remains the weakest model for absolute pose estimation (Fig. 4d, Fig. 4b), independent of MPJA or MBBA.

### VIII. CONCLUSIONS

In this paper, we have explored the integration of various fisheye camera and reprojection models into a state-of-the-art 3D HPE framework. Our novel *FISHnCHIPS* dataset, with focus on challenging camera perspectives encountered in robotics use-cases, allowed us to assess these projection models on close-up interaction scenarios where the person covers a significant part of the camera's FOV. Here, we found that PH and CC projections deliver suboptimal performance. In contrast, EF and DS projections performed best, with EC projections performing slightly worse. We also introduced a heuristic for dynamically selecting a suitable projection model based on estimated bounding box geometry.

In summary, we believe that our experiments shed valuable insights into how fisheye cameras can successfully be used for 3D HPE in various human-robot interaction scenarios.

### ACKNOWLEDGMENT

## REFERENCES

[1] H. Ai, Z. Cao, J. Zhu, H. Bai, Y. Chen, and L. Wang, "Deep learning for omnidirectional vision: A survey and new perspectives," arXiv:2205.10468, 2022.

[2] M. Cormier, A. Clepe, A. Specker, and J. Beyerer, "Where are we with human pose estimation in real-world surveillance?" in *WACV Workshops*, 2022.

[3] Y. Cheng, P. Yi, R. Liu, J. Dong, D. Zhou, and Q. Zhang, "Human-robot interaction method combining human pose estimation and motion intention recognition," in *International Conference on Computer Supported Cooperative Work in Design*, 2021.

[4] J. Lee, M. Lee, D. Lee, and S. Lee, "Hierarchically Decomposed Graph Convolutional Networks for Skeleton-Based Action Recognition," in *ICCV*, 2023.

[5] C. K. Ingwersen, C. Mikkelstrup, J. N. Jensen, M. R. Hannemose, and A. B. Dahl, "SportsPose: A dynamic 3D sports pose dataset," in *CVPR Workshops*, 2023.

[6] J. Jj and P. Maheswari, "AI and augmented reality for 3D Indian dance pose reconstruction cultural revival," *Scientific Reports*, vol. 14, 2024.

[7] Y. Zhang, S. You, S. Karaoglu, and T. Gevers, "3D human pose estimation and action recognition using fisheye cameras: A survey and benchmark," *Pattern Recognition*, vol. 162, 2025.

[8] ——, "Multi-person 3D pose estimation from a single image captured by a fisheye camera," *CVIU*, 2022.

[9] K. Minoda and T. Yairi, "3D human pose estimation in weightless environments using a fisheye camera," in *IROS*, 2022.

[10] I. Sárándi, T. Linder, K. O. Arras, and B. Leibe, "MeTRAbs: Metric-scale truncation-robust heatmaps for absolute 3D human pose estimation," *IEEE T-BIOM*, vol. 3, no. 1, 2021.

[11] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake, "Real-time human pose recognition in parts from single depth images," in *CVPR*, 2011.

[12] Y. Ren, C. Zhao, Y. He, P. Cong, H. Liang, J. Yu, L. Xu, and Y. Ma, "Lidar-aid inertial poser: Large-scale human motion capture by sparse inertial and lidar sensors," *IEEE T-VCG*, vol. 29, no. 5, 2023.

[13] T. D. Nguyen and M. Kresovic, "A survey of top-down approaches for human pose estimation," in *arXiv:2202.02656*, 2022.

[14] J. Su, C. Wang, X. Ma, W. Zeng, and Y. Wang, "VirtualPose: Learning generalizable 3D human pose models from virtual data," in *ECCV*, 2022.

[15] I. Sárándi, A. Hermans, and B. Leibe, "Learning 3D human pose estimation from dozens of datasets using a geometry-aware autoencoder to bridge between skeleton formats," in *WACV*, 2023.

[16] A. Matsune, S. Hu, G. Li, S. Wen, X. Zhu, and Z. Tan, "A geometry loss combination for 3D human pose estimation," in *WACV*, 2024.

[17] I. Sárándi, T. Linder, K. O. Arras, and B. Leibe, "Metric-scale truncation-robust heatmaps for 3D human pose estimation," in *IEEE Conf. Autom. Face and Gesture Recog.*, 2020.

[18] J. Zhen, Q. Fang, J. Sun, W. Liu, W. Jiang, H. Bao, and X. Zhou, "SMAP: single-shot multi-person absolute 3D pose estimation," in *ECCV*, 2020.

[19] E. Samkari, M. Arif, M. Alghamdi, and M. A. Al Ghamdi, "Human pose estimation using deep learning: A systematic literature review," *Machine Learning and Knowledge Extraction*, no. 4, 2023.

[20] J. Martinez, R. Hossain, J. Romero, and J. J. Little, "A simple yet effective baseline for 3D human pose estimation," in *ICCV*, 2017.

[21] W. Liu, Q. Bao, Y. Sun, and T. Mei, "Recent advances of monocular 2D and 3D human pose estimation: A deep learning perspective," *ACM CS.*, vol. 55, no. 4, 2022.

[22] V. R. Kumar, C. Eising, C. Witt, and S. Yogamani, "Surround-view fisheye camera perception for automated driving: Overview, survey & challenges," *IEEE T-ITS*, 2023.

[23] W. Xu, A. Chatterjee, M. Zollhoefer, H. Rhodin, P. Fua, H.-P. Seidel, and C. Theobalt, "Mo$^2$Cap$^2$: Real-time mobile 3D motion capture with a cap-mounted fisheye camera," *IEEE T-VCG*, 2019.

[24] D. Tome, P. Peluse, L. Agapito, and H. Badino, "xR-EgoPose: Egocentric 3D human pose from an HMD camera," in *ICCV*, 2019.

[25] H. Rhodin, C. Richardt, D. Casas, E. Insafutdinov, M. Shafiei, H.-P. Seidel, B. Schiele, and C. Theobalt, "EgoCap: egocentric marker-less motion capture with two fisheye cameras," *ACM ToG*, 2016.

[26] H. Akada, J. Wang, S. Shimada, M. Takahashi, C. Theobalt, and V. Golyanik, "UnrealEgo: A new dataset for robust egocentric 3D human motion capture," in *ECCV*, 2022.

[27] Y. Liu, J. Yang, X. Gu, Y. Chen, Y. Guo, and G.-Z. Yang, "EgoFish3D: Egocentric 3D pose estimation from a fisheye camera via self-supervised learning," *IEEE TTM*, 2023.

[28] S. Ravi, P. Climent-Perez, T. Morales, C. Huesca-Spairani, K. Hashemifard, and F. Florez-Revuelta, "ODIN: An omnidirectional indoor dataset capturing activities of daily living from multiple synchronized modalities," in *CVPR Workshops*, 2023.

[29] Z. Duan, M. Ozan Tezcan, H. Nakamura, P. Ishwar, and J. Konrad, "RAPiD: Rotation-aware people detection in overhead fisheye images," in *CVPR Workshops*, 2020.

[30] M. O. Tezcan, Z. Duan, M. Cokbas, P. Ishwar, and J. Konrad, "WEPDTOF: A dataset and benchmark algorithms for in-the-wild people detection and tracking from overhead fisheye cameras," in *WACV*, 2022.

[31] K. Aso, D.-H. Hwang, and H. Koike, "Portable 3D human pose estimation for human-human interaction using a chest-mounted fisheye camera," in *AHs*, ser. AHs '21. ACM, 2021.

[32] D. Scaramuzza, A. Martinelli, and R. Siegwart, "A toolbox for easily calibrating omnidirectional cameras," in *IROS*, 2006.

[33] E. Plaut, E. Ben-Yaacov, and B. E. Shlomo, "3D object detection from a single fisheye image without a single fisheye training image," in *CVPR Workshops*, 2021.

[34] E. Plaut, "Tutorial on computer vision with fisheye cameras," accessed: 2023-08-15. [Online]. Available: https://plaut.github.io/fisheye_tutorial

[35] J. Yu, D. Nandi, R. Seidel, and G. Hirtz, "NToP: NeRF-powered large-scale dataset generation for 2D and 3D human pose estimation in top-view fisheye images," in *arXiv:2402.18196*, 2024.

[36] K. Grauman, A. Westbury, L. Torresani, K. Kitani, J. Malik, T. Afouras, *et al.*, "Ego-Exo4D: Understanding skilled human activity from first- and third-person perspectives," in *CVPR*, 2024.

[37] A. Dhamanaskar, M. Dimiccoli, E. Corona, A. Pumarola, and F. Moreno-Noguer, "Enhancing egocentric 3D pose estimation with third person views," *PR*, 2023.

[38] L. Ma, Y. Ye, F. Hong, V. Guzov, Y. Jiang, R. Postyeni, L. Pesqueira, A. Gamino, V. Baiyya, H. J. Kim, K. Bailey, D. S. Fosas, C. K. Liu, Z. Liu, J. Engel, R. D. Nardi, and R. Newcombe, "Nymeria: A massive collection of multimodal egocentric daily motion in the wild," arXiv:2406.09905, 2024.

[39] R. Khirodkar, A. Bansal, L. Ma, R. Newcombe, M. Vo, and K. Kitani, "EgoHumans: An egocentric 3D multi-human benchmark," arXiv:2305.16487, 2023.

[40] T. von Marcard, R. Henschel, M. Black, B. Rosenhahn, and G. Pons-Moll, "Recovering accurate 3D human pose in the wild using IMUs and a moving camera," in *ECCV*, 2018.

[41] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu, "Human3.6M: Large scale datasets and predictive methods for 3D human sensing in natural environments," *PAMI*, 2014.

[42] D. Mehta, O. Sotnychenko, F. Mueller, W. Xu, S. Sridhar, G. Pons-Moll, and C. Theobalt, "Single-shot multi-person 3D pose estimation from monocular RGB," in *3DV*, 2018.

[43] C. Hughes, R. McFeely, P. Denny, M. Glavin, and E. Jones, "Equidistant (fθ) fish-eye perspective with application in distortion centre estimation," *Image and Vision Computing*, 2010.

[44] V. Usenko, N. Demmel, and D. Cremers, "The double sphere camera model," in *3DV*, 2018.

[45] C. Lyu, W. Zhang, H. Huang, Y. Zhou, Y. Wang, Y. Liu, S. Zhang, and K. Chen, "RTMDet: An empirical study of designing real-time object detectors," in *arXiv:2212.07784*, 2022.

[46] "Yoga pose," accessed: 2024-09-13. [Online]. Available: https://pixabay.com/illustrations/yoga-yoga-pose-side-plank-5494710/

[47] T. Linder, K. Yilmaz, D. Adrian, and B. Leibe, "Acquisition of high-quality images for camera calibration in robotics applications via speech prompts," in *German Robotics Conference (GRC)*, 2025.

[48] Z. Zong, G. Song, and Y. Liu, "DETRs with collaborative hybrid assignments training," in *ICCV*, 2023.

[49] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin Transformer: Hierarchical vision transformer using shifted windows," in *ICCV*, 2021.

[50] T. Jiang, P. Lu, L. Zhang, N. Ma, R. Han, C. Lyu, Y. Li, and K. Chen, "RTMPose: Real-time multi-person pose estimation based on MMPose," in *arXiv:2303.07399*, 2023.

[51] S. Agarwal, K. Mierle, and T. C. S. Team, "Ceres Solver," 2023. [Online]. Available: https://github.com/ceres-solver/ceres-solver

[52] M. Tan and Q. Le, "EfficientNetV2: Smaller models and faster training," in *ICML*, 2021.

[53] Y. Cheng, B. Wang, and R. T. Tan, "Dual networks based 3D multi-person pose estimation from monocular video," *PAMI*, 2023.

[54] K. Sun, B. Xiao, D. Liu, and J. Wang, "Deep high-resolution representation learning for human pose estimation," in *CVPR*, 2019.

[55] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common objects in context," in *ECCV*, 2014.

[56] N. Ugrinovic, A. Ruiz, A. Agudo, A. Sanfeliu, and F. Moreno-Noguer, "Permutation-invariant relational network for multi-person 3D pose estimation," in *arXiv:2204.04913*, 2022.

[57] G. Moon, J. Y. Chang, and K. M. Lee, "Camera distance-aware top-down approach for 3D multi-person pose estimation from a single RGB image," in *ICCV*, 2019.

[58] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.

[59] Z. Qiu, Q. Yang, J. Wang, and D. Fu, "Dynamic graph reasoning for multi-person 3D pose estimation," in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022.

[60] A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," in *ECCV*, 2016.

[61] I. Sarandi, "Robust and efficient methods in visual 3D human pose estimation," Dissertation, RWTH Aachen University, 2023.