# RAPTR: Radar-based 3D Pose Estimation using Transformer

**Sorachi Kato**[1,2]*, **Ryoma Yataka**[1,3], **Pu (Perry) Wang**[1]†, **Pedro Miraldo**[1],
**Takuya Fujihashi**[2], **Petros Boufounos**[1]
[1]Mitsubishi Electric Research Laboratories (MERL), USA
[2]The University of Osaka, Japan
[3]Information Technology R&D Center (ITC), Mitsubishi Electric Corporation, Japan

## Abstract

Radar-based indoor 3D human pose estimation typically relied on fine-grained 3D keypoint labels, which are costly to obtain especially in complex indoor settings involving clutter, occlusions, or multiple people. In this paper, we propose **RAPTR** (RAdar Pose esTimation using tRansformer) under weak supervision, using only 3D BBox and 2D keypoint labels which are considerably easier and more scalable to collect. Our RAPTR is characterized by a two-stage pose decoder architecture with a pseudo-3D deformable attention to enhance (pose/joint) queries with multi-view radar features: a pose decoder estimates initial 3D poses with a 3D template loss designed to utilize the 3D BBox labels and mitigate depth ambiguities; and a joint decoder refines the initial poses with 2D keypoint labels and a 3D gravity loss. Evaluated on two indoor radar datasets, RAPTR outperforms existing methods, reducing joint position error by 34.3% on HIBER and 76.9% on MMVR. Our implementation is available at https://github.com/merlresearch/radar-pose-transformer.

## 1 Introduction

Accurate human perception is essential for indoor applications, including elderly monitoring, smart building management, and robotic navigation. Although vision sensors offer high spatial resolution, they raise privacy concerns and perform poorly under low light, occlusions, and hazardous conditions (fire or smoke). In contrast, radar provides penetration capability, robustness to adverse conditions, and low deployment cost, ideal for privacy-preserving indoor sensing [44, 18, 23, 30, 26].

By processing 4D radar tensors, RF-Pose 3D [48] demonstrated through-the-wall 3D pose estimation with a convolutional neural network (CNN), while HRRadarPose [11] employed an hourglass neural network HRNet [34]. mRI [1] is a multi-modal 3D human pose estimation dataset that integrates mmWave radar, RGB-D cameras, and inertial sensors to facilitate research in human pose estimation and action detection. QRFPose [33] is a novel approach that adopts a DETR [3]-style query mechanism for end-to-end 3D regression using multi-view radar perceptions. Existing pipelines often rely on expensive fine-grained 3D keypoint labels [35], typically collected using non-portable 3D motion capture systems such as VICON, or using LiDAR, which can still suffer from occlusions and incomplete observations.

Collecting cheaper, lower-cost labels, such as fine-grained 2D keypoints in the image plane and/or coarse-grained 3D bounding boxes (BBoxes), is considerably easier and more scalable particularly in complex indoor settings (e.g., cluttered, occlusion, multi-person), compared with acquiring dense

---

*The work was initiated during his internship at MERL.
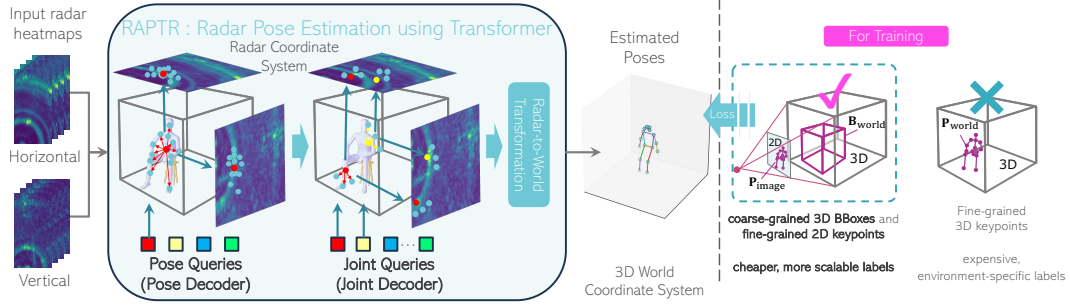†Project Lead.

Figure 1: RAPTR takes multi-view radar heatmaps as inputs and performs a novel Pseudo-3D deformable attention between (pose and joint) queries and multi-view radar features in a two-stage decoder to estimate 3D human poses in a 3D coordinate system. Rather than relying on expensive, environment-specific fine-grained 3D keypoint labels, RAPTR makes use of cheaper, more scalable labels such as coarse-grained 3D BBoxes and fine-grained 2D keypoints to train the model.

3D keypoints labels. Examples include RF-Pose [47], HuPR [15], and, more recently, MMVR [24] datasets. To the best of our knowledge, the use of 2D keypoints and 3D BBoxes, as a substitute for costly 3D keypoints, for radar-based 3D human pose estimation has not been systematically investigated in the literature before.

To address this gap, we propose **RAPTR** (RAdar Pose esTimation using tRansformer) in Fig. 1, a radar-based pipeline designed to take multi-view radar heatmaps as inputs and estimate 3D human poses under weak supervision using only 3D BBox and 2D keypoint labels. RAPTR builds on the two-stage (pose and joint) decoder architecture of the state-of-the-art RGB-based 2D pose estimation PETR framework [28] and introduces a structural loss function that is designed to utilize weak supervision labels to mitigate the depth ambiguity. RAPTR also lifts the 2D deformable attention in PETR to a pseudo-3D deformable attention, wherein reference points (dots in Fig. 1) and offsets (arrows in Fig. 1) are proposed in the 3D radar coordinate system and projected onto multiple radar views (dots on the radar heatmaps in Fig. 1) to eliminate redundant per-view offset estimation and offer better scalability as the number of radar views increases. Our model outperforms a list of radar-based 3D pose baselines over two indoor radar datasets: HIBER [35] and MMVR [24]. The main contributions of this work are:

- To the best of our knowledge, RAPTR is the first radar-based 3D human pose estimation framework to explicitly utilize low-cost weak supervision in the form of 3D BBoxes and 2D keypoints, rather than relying on fine-grained 3D keypoint labels.

- We introduce a structured loss function that tightly couples the two-stage decoder architecture to enable 3D pose estimation under weak supervision. Specifically, we design a 3D Template Loss, which utilizes the 3D BBox labels at the pose decoder, and a combined 3D Gravity and 2D Keypoint Loss at the pose decoder, allowing RAPTR to effectively learn geometrically consistent 3D poses from weak supervision.

- We further introduce a pseudo-3D deformable attention mechanism to bridge the 3D spatial domain and 2D radar views, enabling scalable view association while preserving pose estimation performance.

## 2 Related Work

**Human Pose Estimation with RGB Image:** Human pose estimation from images involves localizing body joints for multiple subjects and associating them for each subject. Existing architectures fall into two main paradigms: top-down and bottom-up. The top-down methods first detect each person using detectors such as Faster R-CNN [25] or Mask R-CNN [10], then applying a single-person pose estimator to each cropped region. These approaches achieve state-of-the-art accuracy with models like Stacked-Hourglass [22], HRNet [31], and DarkPose [46]. In contrast, bottom-up methods such as OpenPose [2], HigherHRNet [6], and SAHR [19] bypass the detection step by predicting all joint candidates across the entire image and grouping them into individuals. PETR [28] introduces an

end-to-end pose estimation framework using a query-based, two-stage transformer decoder architecture. Beyond 2D, recent methods addresses 3D pose from RGB or RGB-D inputs, either by directly regressing 3D joints [20] or by lifting 2D predictions into the 3D space through geometric reasoning or weak supervision [4, 5, 21, 32].

**Human Pose Estimation with radar or radio frequency signals:**  Recent studies have shown that information extracted from commercial radars is sufficiently informative to perform fine-grained human pose estimation, both for 2D and 3D. Despite the coarse-grained nature of the radar point clouds (PCs), deep neural pipelines have achieved a multitude of performance gains [29, 27, 38, 1, 41, 36, 12, 14, 39, 42, 8, 7, 43]. On the other hand, methods using raw radar measurements and radar heatmaps have been widely explored [47, 49, 15, 35, 11, 24, 33, 45, 37]. RF-Pose [47] pioneered multi-view 3D CNNs for through-wall 2D estimation. HuPR [15] refines such heatmaps via a graph convolutional network (GCN). HRRadarPose [11] adopts an HRNet-style [34] single-stage head for 3D output. QRFPose [33], based on a DETR-style Transformer [3] for end-to-end query-based 3D pose estimation, is the closest baseline to ours. It differs by applying per-view 2D deformable attention and using a single decoder for all keypoints, followed by grouping. In contrast, our method employs pseudo-3D deformable attention and a two-stage decoder.

## 3 Preliminary

**Multi-View Radar Heatmaps:**  As shown in Fig. 2, two synchronized radar arrays (horizontal and vertical) collect reflected pulses that form a 3D data cube per array (ADC samples × pulses × elements). A 3D FFT converts each cube into a range–Doppler–angle spectrum, whose angle dimension is azimuth for the horizontal array and elevation for the vertical array. After Doppler-axis integration to boost SNR (signal-to-noise ratio), we obtain two polar 2D heatmaps (range-azimuth and range-elevation). These are mapped to the Cartesian space: $\mathbf{Y}_{\text{hor}}(t) \in \mathbb{R}^{W \times D}$ for horizontal-depth and $\mathbf{Y}_{\text{ver}}(t) \in \mathbb{R}^{H \times D}$ for vertical-depth at frame $t$. The temporal context is captured by stacking $T$ consecutive frames, giving $\mathbf{Y}_{\text{hor}} \in \mathbb{R}^{T \times W \times D}$ and $\mathbf{Y}_{\text{ver}} \in \mathbb{R}^{T \times H \times D}$.
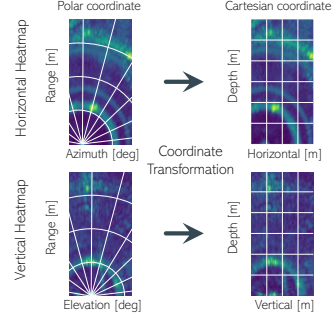


Figure 2: Multi-view radar heatmaps.

**Problem Formulation:**  The 3D pose estimation task takes $T$ consecutive radar frames, $\mathbf{Y}_{\text{hor}}$ and $\mathbf{Y}_{\text{ver}}$, as input and estimates poses $\hat{\mathbf{P}}_{\text{world}}$ in the 3D world coordinate system,

$$\hat{\mathbf{P}}_{\text{world}} = \mathcal{T}_{\text{r2w}}(\hat{\mathbf{P}}_{\text{radar}}) = \mathcal{T}_{\text{r2w}}(f(\mathbf{Y}_{\text{hor}}, \mathbf{Y}_{\text{ver}})), \qquad (1)$$

where $f$ represents the 3D pose estimation pipeline in the 3D radar coordinate system, and $\mathcal{T}_{\text{r2w}}$ is a known radar-to-world coordinate transformation that converts the estimated 3D poses into the 3D world coordinate system. Rather than relying on costly, non-scalable fine-grained 3D keypoint labels $\mathbf{P}_{\text{world}}$, we consider cheaper, more scalable labels such as coarse-grained 3D BBoxes $\mathbf{B}_{\text{world}}$ and fine-grained 2D keypoints $\mathbf{P}_{\text{image}}$ for supervision, as shown in Fig. 1.

## 4 RAPTR: Radar-based 3D Pose Estimation using Transformer

We present the RAPTR architecture in Fig. 3, following a left-to-right order, and highlight radar-specific modifications. Refer to Appendix A for detailed architecture and computational complexity.

### 4.1 Architecture

**Backbone**: Given $\mathbf{Y}_{\text{hor}} \in \mathbb{R}^{T \times W \times D}$ and $\mathbf{Y}_{\text{ver}} \in \mathbb{R}^{T \times H \times D}$, a shared backbone network (e.g., ResNet [9]) generates separate multi-scale horizontal-view and vertical-view radar feature maps: $\mathbf{Z}_{\text{hor}} = \{\mathbf{Z}_{\text{hor},i}\}_{i=1}^{S} = \texttt{backbone}(\mathbf{Y}_{\text{hor}})$ and $\mathbf{Z}_{\text{ver}} = \{\mathbf{Z}_{\text{ver},i}\}_{i=1}^{S} = \texttt{backbone}(\mathbf{Y}_{\text{ver}})$, where the $i$-th scale feature maps $\mathbf{Z}_{\text{hor},i} \in \mathbb{R}^{W_i \times D_i \times d}$ and $\mathbf{Z}_{\text{ver},i} \in \mathbb{R}^{H_i \times D_i \times d}$ have a spatial dimension of $W_i \times D_i$ or $H_i \times D_i$ and a feature dimension of $d$, and $S$ is the number of scales.
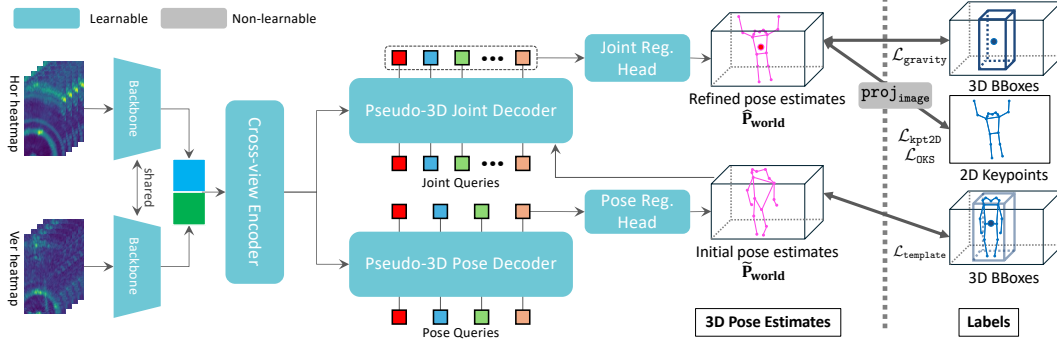
Figure 3: The RAPTR architecture consists of: 1) **Cross-view Encoder** that extracts multi-scale radar features; 2) **Pseudo-3D Pose Decoder** that enhances pose queries via a pseudo-3D deformable attention and predicts initial 3D poses; and 3) **Pseudo-3D Joint Decoder** that further refines joint queries and outputs final 3D poses. In terms of **loss function**, RAPTR leverages 3D BBox and 2D keypoint labels through coarse-grained 3D loss (gravity and template) and 2D keypoint loss.

**Cross-View Encoder** is a Transformer encoder with $L_{\text{enc}}$ layers that fuses the horizontal- and vertical-view radar features. Each layer runs a shared cross-attention twice: first with $\mathbf{Z}_{\text{hor}}$ as key/value and $\mathbf{Z}_{\text{ver}}$ as query, then vice versa. This bidirectional exchange embeds complementary cues, while residual connections keep view-specific details, producing refined features $\mathbf{F}_{\text{enc}}^{(i)}$, $i = 1, \cdots, L_{\text{enc}}$,

$$\mathbf{F}_a^{(i)} = \mathbf{F}_a^{(i-1)} + \texttt{CrossAttn}(\mathbf{F}_a^{(i-1)}, \mathbf{F}_b^{(i-1)}), \quad (a,b) \in \{(\texttt{hor}, \texttt{ver}), (\texttt{ver}, \texttt{hor})\}, \quad (2)$$

where $\mathbf{F}_{\text{hor}}^{(0)} = \mathbf{Z}_{\text{hor}}$, $\mathbf{F}_{\text{ver}}^{(0)} = \mathbf{Z}_{\text{ver}}$, and $\texttt{CrossAttn}(\cdot, \cdot)$ denotes the deformable cross-attention [50] following [28] with fixed positional embeddings added beforehand for efficiency. After $L_{\text{enc}}$ iterations, the encoded features $\mathbf{F}_{\text{hor}}$ and $\mathbf{F}_{\text{ver}}$ are obtained at the output of the cross-view encoder.

**Pseudo-3D Pose Decoder** associates $N$ pose queries $\mathbf{Q}_{\text{pose}} \in \mathbb{R}^{N \times d}$ (embedding dimension $d$) with encoded radar features $(\mathbf{F}_{\text{hor}}, \mathbf{F}_{\text{ver}})$, where each query corresponds to a reference pose refined through pseudo-3D deformable attention over $L_{\text{pose}}$ layers. We define the $l$-th decoder layer as a function $\mathcal{D}_{\text{pose}}^{(l)}$ that updates both the pose queries and reference poses in the 3D radar space:

$$(\mathbf{Q}_{\text{pose}}^{(l)}, \tilde{\mathbf{P}}_{\text{radar}}^{(l)}) = \mathcal{D}_{\text{pose}}^{(l)}(\mathbf{Q}_{\text{pose}}^{(l-1)}, \mathbf{F}_{\text{hor}}, \mathbf{F}_{\text{ver}}, \tilde{\mathbf{P}}_{\text{radar}}^{(l-1)}), \quad (3)$$

where $\tilde{\mathbf{P}}_{\text{radar}}^{(0)}$ is initialized by passing $\mathbf{Q}_{\text{pose}}$ through an MLP. Reference poses are iteratively refined by applying predicted coordinate offsets $\Delta\tilde{\mathbf{P}}_{\text{radar}}^{(l)}$ in the normalized scale:

$$\tilde{\mathbf{P}}_{\text{radar}}^{(l)} \in \mathbb{R}^{N \times 3K} = \sigma(\sigma^{-1}(\tilde{\mathbf{P}}_{\text{radar}}^{(l-1)}) + \Delta\tilde{\mathbf{P}}_{\text{radar}}^{(l-1)}), \quad l = 1, \ldots, L_{\text{pose}}, \quad (4)$$

where $\sigma$ and $\sigma^{-1}$ denote the Sigmoid function and its inverse. The predicted offsets $\Delta\tilde{\mathbf{P}}_{\text{radar}}^{(l)} = H_{\text{pose}}(\mathbf{Q}_{\text{pose}}^{(l)})$ are obtained by passing pose queries at each layer to a shared regression head $H_{\text{pose}}$.

We convert the initial pose estimates $\tilde{\mathbf{P}}_{\text{radar}} = \tilde{\mathbf{P}}_{\text{radar}}^{(L_{\text{pose}})}$ from the radar coordinate system to the world coordinate system via $\tilde{\mathbf{P}}_{\text{world}} = \mathcal{T}_{\text{r2w}}(\tilde{\mathbf{P}}_{\text{radar}})$, along with the corresponding confidence scores $\tilde{\mathbf{c}}$. We defer the pseudo-3D deformable attention to Section 4.2.

**Pseudo-3D Joint Decoder** associates $K$ joint queries $\mathbf{Q}_{\text{joint}} \in \mathbb{R}^{K \times d}$ with encoded radar features $(\mathbf{F}_{\text{hor}}, \mathbf{F}_{\text{ver}})$, where each query corresponds to a single joint refined by pseudo-3D deformable attention over $L_{\text{joint}}$ layers. Here, $K$ joint queries correspond to the same subject. We define the $l$-th decoder layer as a function $\mathcal{D}_{\text{joint}}^{(l)}$ that updates both the joint queries and corresponding joints:

$$(\mathbf{Q}_{\text{joint}}^{(l)}, \tilde{\mathbf{p}}_{i,\text{radar}}^{(l)}) = \mathcal{D}_{\text{joint}}^{(l)}(\mathbf{Q}_{\text{joint}}^{(l-1)}, \mathbf{F}_{\text{hor}}, \mathbf{F}_{\text{ver}}, \tilde{\mathbf{p}}_{i,\text{radar}}^{(l-1)}), \quad (5)$$

where $\tilde{\mathbf{p}}_{i,\text{radar}}^{(l)} \in \mathbb{R}^{K \times 3}, i = 1, \cdots, N$ is one specific pose in the $N$ poses, and $\tilde{\mathbf{p}}_{i,\text{radar}}^{(0)}$ is $i$-th pose prediction from the pose decoder. Joints in the reference pose are iteratively refined by applying predicted coordinate offsets $\Delta\tilde{\mathbf{p}}_{i,\text{radar}}^{(l)}$:

$$\tilde{\mathbf{p}}_{i,\text{radar}}^{(l)} \in \mathbb{R}^{K \times 3} = \sigma(\sigma^{-1}(\tilde{\mathbf{p}}_{i,\text{radar}}^{(l-1)}) + \Delta\tilde{\mathbf{p}}_{i,\text{radar}}^{(l-1)}), \quad l = 1, \cdots, L_{\text{joint}}, \quad (6)$$

4

where the predicted offsets are given as $\Delta\tilde{\mathbf{p}}_{i,\mathtt{radar}}^{(l)} = H_{\mathtt{joint}}(\mathbf{Q}_{\mathtt{joint}}^{(l)})$ with a shared regression head.

We collect refined reference poses from the joint decoder as $\hat{\mathbf{P}}_{\mathtt{radar}} = \{\tilde{\mathbf{p}}_{i,\mathtt{radar}}^{(L_{\mathtt{joint}})}\}_{i=1}^{N}$ and convert them into the 3D world coordinate system as $\hat{\mathbf{P}}_{\mathtt{world}} = \mathcal{T}_{\mathtt{r2w}}(\hat{\mathbf{P}}_{\mathtt{radar}})$.

## 4.2 Pseudo-3D Deformable Attention

Our two-stage decoder incorporates a pseudo-3D deformable attention module, where "pseudo" highlights that reference points and sampling offsets are defined in 3D space, while feature sampling occurs on the 2D radar views, as illustrated in Fig. 4.

Consider a 3D reference point $(x, y, z)$ in the 3D radar space with a corresponding query $\mathbf{q} \in \mathbb{R}^d$ (from either pose queries in the pose decoder or joint queries in the joint decoder). We first feed $\mathbf{q}$ into a linear projection layer to predict a set of 3D sampling offsets $\{\Delta x_i, \Delta y_i, \Delta z_i\}_{i=1}^{N_{\mathtt{offset}}}$. Given the 3D reference point and sampling offsets, we can locate the 3D sampling coordinates and project them onto the two



Figure 4: The pseudo-3D deformable attention operates on a 3D reference point and 3D sampling offsets that are projected to different radar views for pseudo-3D attention between multi-view radar features and the query.

radar views, extracting deformable multi-view radar features:

$$\mathbf{f}_{\mathtt{hor}}^{(i)} = \mathbf{F}_{\mathtt{hor}}(x + \Delta x_i, z + \Delta z_i), \quad \mathbf{f}_{\mathtt{ver}}^{(i)} = \mathbf{F}_{\mathtt{ver}}(y + \Delta y_i, z + \Delta z_i) \quad i = 1, \cdots, N_{\mathtt{offset}}. \quad (7)$$

We group deformable multi-view radar features as $\mathbf{F}_{\mathtt{attn}} = \{\mathbf{f}_{\mathtt{hor}}^{(1)}, \mathbf{f}_{\mathtt{ver}}^{(1)}, \cdots, \mathbf{f}_{\mathtt{hor}}^{(N_{\mathtt{offset}})}, \mathbf{f}_{\mathtt{ver}}^{(N_{\mathtt{offset}})}\}$.

Meanwhile, multi-view attention weights $f_{\mathtt{attn}} \in \mathbb{R}^{N_{\mathtt{offset}} \times 2}$ (where 2 corresponds to the two radar views) are proposed by linearly projecting the query and applying a softmax normalization. These weights capture the relative importance of radar features across the two views in a unified manner.

Given the deformable multi-view radar features $\mathbf{F}_{\mathtt{attn}}$ and the multi-view attention weights $f_{\mathtt{attn}}$, deformable multi-view attention features can be calculated as

$$\bar{\mathbf{F}}_{\mathtt{attn}} = \sum_{i=1}^{N_{\mathtt{offset}}} (A_{i,0}\mathbf{W}\mathbf{F}_{\mathtt{attn}}^{(2i-1)} + A_{i,1}\mathbf{W}\mathbf{F}_{\mathtt{attn}}^{(2i)}), \quad (8)$$

where $A_{i,0}$ and $A_{i,1}$ are the attention weights in $f_{\mathtt{attn}}$ for the $i$-th deformable radar feature in the horizontal and, respectively, vertical radar views, and $\mathbf{W} \in \mathbb{R}^{C \times C}$ is a learnable weight matrix. We denote the overall pseudo-3D deformable attention as $\bar{\mathbf{F}}_{\mathtt{attn}} = \mathtt{DeformableAttn}(\mathbf{F}_{\mathtt{ver}}, \mathbf{F}_{\mathtt{hor}}, (x, y, z), \mathbf{q})$.

Appendix B provides implementation details of $\mathtt{DeformableAttn}(\cdot, \cdot, \cdot, \cdot)$ and a computational complexity comparison with the decoupled 2D deformable attention used in QRFPose [33], demonstrating better scalability of the proposed pseudo-3D attention as the number of radar views increases. Appendix C describes an optional view mask module (top right of Fig. 4) that adds flexibility in selecting multi-view radar features per query. For example, an all-zero mask can be applied to exclude features from a specific radar view.

## 4.3 Structural Loss Function

As illustrated in Fig. 3, RAPTR utilizes weak supervision labels: coarse-grained BBox labels $\mathbf{B}_{\mathtt{world}}$ in the 3D world coordinate system and fine-grained 2D keypoint labels $\mathbf{P}_{\mathtt{image}}$ in the image plane. The loss function is calculated between these labels $\{\mathbf{B}_{\mathtt{world}}, \mathbf{P}_{\mathtt{image}}\}$ and the initial and refined 3D pose estimates $\{\tilde{\mathbf{P}}_{\mathtt{world}}, \hat{\mathbf{P}}_{\mathtt{world}}\}$ with details included in Appendix D.

**3D Template (T3D) Loss at Pose Decoder** utilizes coarse-grained 3D BBox labels $\mathbf{B}_{\mathtt{world}}$. For each $\mathbf{B}_{\mathtt{world}}$, we construct a 3D keypoint template by computing the centroid of the corresponding 3D BBox, which serves as the 3D gravity center label $\mathbf{g}_{\mathtt{world}} \in \mathbb{R}^{1 \times 3}$.
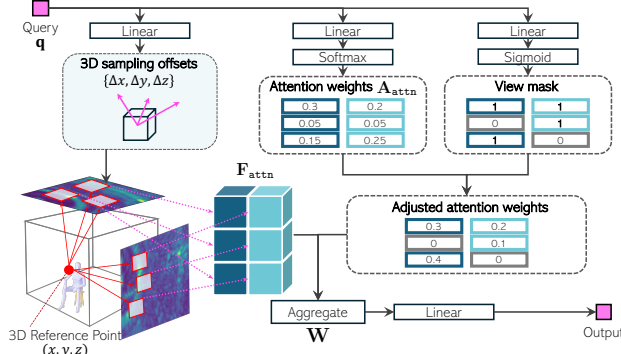
Then, given a keypoint template defined at the coordinate origin $\mathbf{K}_{\mathrm{world}} \in \mathbb{R}^{K \times 3}$, the corresponding template pose $\mathbf{T}_{\mathrm{world}}$ is computed as $\mathbf{T}_{\mathrm{world}} = \mathbf{K}_{\mathrm{world}} + \mathbf{1}^{\top} \mathbf{g}_{\mathrm{world}}$. As illustrated in the lower right of Fig. 3, the T3D loss $\mathcal{L}_{\mathrm{template}}$ is defined as the Euclidean distance between the template poses $\mathbf{T}_{\mathrm{world}}$ and the initial 3D pose estimates $\tilde{\mathbf{P}}_{\mathrm{world}}$ at the pose decoder.

**Combined 3D Gravity (G3D) Loss and 2D Keypoint (K2D) Loss at Joint Decoder** utilizes both the coarse-grained 3D BBox labels $\mathbf{B}_{\mathrm{world}}$ and the fine-grained 2D keypoint labels $\mathbf{P}_{\mathrm{image}}$ in the image plane, as illustrated in the upper right of Fig. 3

For the G3D loss, the refined 3D pose estimate $\hat{\mathbf{P}}_{\mathrm{world}}$ is collapsed into its centroid as $\hat{\mathbf{g}}_{\mathrm{world}} \in \mathbb{R}^{1 \times 3}$ by averaging the keypoint coordinates along each spatial axis. The resulting G3D loss $\mathcal{L}_{\mathrm{gravity}}$ is then defined as the Euclidean distance between the predicted and ground-truth 3D gravity centers, $\hat{\mathbf{g}}_{\mathrm{world}}$ and $\mathbf{g}_{\mathrm{world}}$.

For the K2D loss, the refined 3D pose estimate $\hat{\mathbf{P}}_{\mathrm{radar}}$ in the radar coordinate system are first transformed into the 3D camera coordinate system via a calibrated coordinate transformation: $\hat{\mathbf{P}}_{\mathrm{camera}} = \mathbf{R}\hat{\mathbf{P}}_{\mathrm{radar}} + \mathbf{1}^{\top}\mathbf{t}$ where $\mathbf{R}$ and $\mathbf{t}$ denote the calibrated 3D rotation matrix and the translation vector, respectively. The resulting 3D camera-space pose estimates are then projected onto the 2D image plane via a known 3D-to-2D projection: $\hat{\mathbf{P}}_{\mathrm{image}} = \mathtt{proj}_{\mathrm{image}}(\hat{\mathbf{P}}_{\mathrm{camera}})$. Finally, the fine-grained 2D loss combines the image-plane Euclidean error $\mathcal{L}_{\mathrm{kpt2D}}$ and the object keypoint similarity (OKS) loss $\mathcal{L}_{\mathrm{OKS}}$ [28] between $\mathbf{P}_{\mathrm{image}}$ and $\hat{\mathbf{P}}_{\mathrm{image}}$.

**Structural Loss Function**: Following the set-based loss in [3], we employ bipartite matching to associate predictions $\{\hat{\mathbf{g}}_{\mathrm{world}}, \tilde{\mathbf{P}}_{\mathrm{world}}, \hat{\mathbf{P}}_{\mathrm{world}}\}$ with their ground-truth labels $\{\mathbf{g}_{\mathrm{world}}, \mathbf{T}_{\mathrm{world}}, \mathbf{P}_{\mathrm{world}}\}$. Based on these associations, we define the structural loss function as

$$\mathcal{L} = \frac{1}{N'} \sum_{i=1}^{N'} (\lambda_1 \mathcal{L}_{\mathrm{template}} + \lambda_2 \mathcal{L}_{\mathrm{gravity}} + \lambda_3 \mathcal{L}_{\mathrm{kpt2D}} + \lambda_4 \mathcal{L}_{\mathrm{OKS}}) + \lambda_5 \mathcal{L}_{\mathrm{cls}}, \qquad (9)$$

where $N'$ is the number of matched pairs, $\lambda_i$ is the corresponding weighting factor for each loss term, and $\mathcal{L}_{\mathrm{cls}}$ is the classification loss of the focal loss [17] with the confidence scores of the matched estimates.

# 5 Evaluation

## 5.1 Settings

**Datasets:** We assess the performance of RAPTR and baseline models on the HIBER dataset[3] [35] and the MMVR dataset[4] [24], both of which are publicly available multi-view mmWave radar datasets designed for indoor human perception tasks. The HIBER dataset includes two-view radar heatmaps from 10 different viewpoints, the corresponding 3D keypoint labels, and the 3D BBox labels. We use data protocols "MULTI" and "WALK", and use views 2 through 10 for training, validation, and testing. The MMVR dataset includes two-view radar heatmaps in various indoor scenarios, the corresponding 2D keypoint labels, and the 3D BBoxes. We use a data split "P1S1", a single-person case in an open space. A detailed description of the datasets is provided in Appendix E.

**Parameter Settings for RAPTR:** We use $T = 4$ consecutive frames as input to our RAPTR network. For the point decoder, the number of pose queries $N$ is 10. For the joint decoder, the number of joint queries $K$ depends on the dataset to be evaluated: $K = 14$ for HIBER and $K = 17$ for MMVR. The parameters relating to model training are summarized in Appendix F.

**Baselines:** We consider the following competitive radar/RF-based 3D pose estimation baselines: **Person-in-WiFi 3D** [40], **HRRadarPose** [11], and **QRFPose** [33]. We evaluate Person-in-WiFi 3D and HRRadarPose using their open-source implementations. As QRFPose has no public code, we reimplement it from scratches and verify similar performance to the original report [33] using 3D keypoint labels. For fair comparison, we adopt a loss function, similar to RAPTR, combining 2D keypoint loss and 3D gravity loss. Baseline implementation details are provided in Appendix G.

---

[3] https://github.com/Intelligent-Perception-Lab/HIBER
[4] https://zenodo.org/records/12611978

Table 1: 3D pose estimation performance on HIBER (MPJPE: cm).

| Env | Method | Head | Neck | Shoulder | Elbow | Wrist | Hip | Knee | Ankle | Overall | (h) | (v) | (d) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| WALK | Person-in-WiFi 3D | 54.28 | 57.01 | 54.18 | 54.81 | 59.98 | 53.98 | 60.32 | 68.84 | 58.25 | 25.60 | 23.94 | 36.20 |
| | QRFPose | 42.23 | 34.21 | 37.37 | 38.05 | 41.25 | 31.24 | 34.39 | 46.87 | 38.20 | 14.78 | 13.40 | 26.76 |
| | HRRadarPose | 30.23 | 25.44 | 33.70 | 34.15 | 42.33 | 27.71 | 31.55 | 40.46 | 33.96 | 15.14 | 13.13 | 19.85 |
| | RAPTR (ours) | 21.75 | 17.41 | 20.72 | 23.23 | 26.55 | 18.97 | 21.06 | 26.10 | **22.32** | **8.41** | **4.85** | **17.73** |
| MULTI | Person-in-WiFi 3D | 88.48 | 85.14 | 89.44 | 84.33 | 84.29 | 88.69 | 81.70 | 81.53 | 85.25 | 34.06 | 28.57 | 58.93 |
| | QRFPose | 49.49 | 44.48 | 45.54 | 46.77 | 49.06 | 40.99 | 41.87 | 51.57 | 46.11 | 18.20 | 14.13 | 34.39 |
| | HRRadarPose | 30.24 | 24.24 | 30.14 | 35.17 | 44.34 | 28.76 | 31.38 | 35.31 | 33.19 | 16.77 | 10.75 | 21.84 |
| | RAPTR (ours) | 18.39 | 13.13 | 16.44 | 20.12 | 24.62 | 15.01 | 17.76 | 23.22 | **18.99** | **7.80** | **4.38** | **14.54** |

**Metrics:** We employ **Mean Per Joint Pose Error (MPJPE)** with the unit of centimeters in the world coordinate. In addition, we evaluate this MPJPE for each body joint and along each 3D axis, horizontal (h), vertical (v), and depth (d), independently. For MMVR, since 3D keypoint labels are not available, we construct a 3D bounding box (BBox) that encloses the estimated 3D keypoints and then use **the distance between the center points** of this BBox and the 3D BBox labels, as well as **the absolute error in the lengths of the edges** along each axis of the box, as metrics to approximate the 3D pose estimation performance. Detailed evaluation metrics are described in Appendix H.

## 5.2 Main results

**HIBER:** Table 1 shows the performance of 3D pose estimation for HIBER, using 2D and coarse 3D labels for baselines and our RAPTR. The qualitative results are provided in Fig. 5a.

For WALK, RAPTR achieves a significantly lower overall MPJPE of 22.32 cm and outperforms all other baselines in the metric. More specifically, RAPTR reduces the overall error by 61.7%, 41.6%, and 34.3% compared to Person-in-WiFi 3D, QRFPose, and HRRadarPose, respectively. The per-joint breakdown demonstrates that RAPTR maintains its performance on relatively challenging joints, such as the wrist and ankle, where other baselines exhibit significant degradation. For example, HRRadarPose reports a wrist error of 42.33 cm, whereas RAPTR reports an error of 26.55 cm. Moreover, RAPTR maintains the error gap between the best- and worst-estimated joints within 10 cm, showing a consistent level of accuracy throughout the body. In terms of directional components, RAPTR shows much lower errors in the horizontal and vertical dimensions than baselines, indicating that RAPTR estimates well-proportioned 3D poses across all axes.
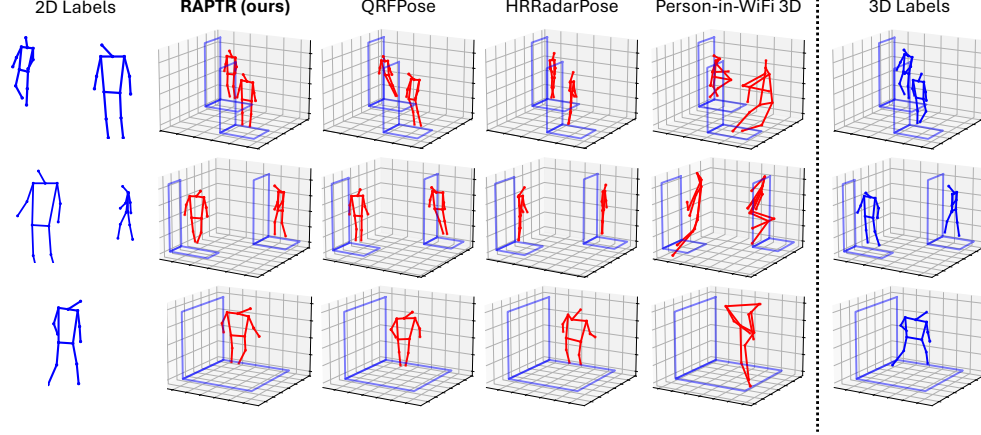
For MULTI, the more challenging multi-person scenario, RAPTR continues to outperform with an overall MPJPE of 18.99 cm and shows a substantial margin compared to the second-best HRRadar-Pose at 33.19 cm. RAPTR reduces the overall error by 77.7%, 58.8%, and 42.7% compared to Person-in-WiFi 3D, QRFPose, and HRRadarPose, respectively. Although the overall accuracy of Person-in-WiFi 3D and QRFPose, noticeably degrades on the MULTI split compared to WALK, likely due to the increased complexity of handling multiple objects, RAPTR maintains a nearly consistent level of performance.

Referring to the qualitative results provided in Fig. 5a, RAPTR estimates structurally consistent 3D poses that match the 3D labels in both position and orientation, while baselines often suffer from misaligned limbs and implausible joint configurations. While baselines often fail to maintain human-like pose structure in the MULTI setting despite performing well in WALK, RAPTR consistently produces plausible estimates in both scenarios, indicating its robustness to multi-person scenes.
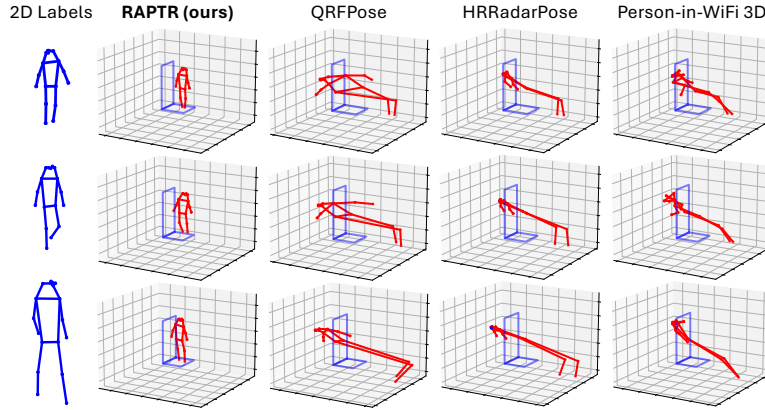
Table 2: Pose estimation performance on MMVR (P1S1).

| Method | Center distance (cm) | Edge length error (cm) | | |
|---|---|---|---|---|
| | | (h) | (v) | (d) |
| Person-in-WiFi 3D | 136.14 | 33.18 | 95.43 | 242.86 |
| QRFPose | 210.75 | 38.12 | 73.69 | 409.38 |
| HRRadarPose | 164.46 | 37.84 | 74.00 | 313.81 |
| RAPTR (ours) | **31.41** | **22.90** | **10.66** | **50.56** |

**MMVR:** Table 2 shows the performance comparison for baselines and RAPTR with MMVR, and Fig. 5b provides qualitative results. Although we cannot directly evaluate the precise 3D pose estimation performance for MMVR due to the absence of 3D pose labels, the results demonstrate that RAPTR effectively preserves reasonable human pose and location accuracy in the 3D space. Specifically, RAPTR shows improvements in center distance by 76.9%, 85.1%, and 80.9% compared to Person-in-WiFi 3D, QRFPose, and HRRadarPose, respectively. As shown in Fig. 5b, other

(a) Visualization of 3D pose estimation by RAPTR and baseline methods on the HIBER dataset.



(b) Visualization of 3D pose estimation by RAPTR and baseline methods on the MMVR dataset.

Figure 5: Qualitative results. Blue lines indicate the keypoint labels, blue rectangles indicate the 3D BBox labels, and red lines indicate the predictions.

baselines exhibit degraded performance due to structural collapse in 3D space, caused by overfitting to 2D alignment when projected onto the image plane. We assume that RAPTR effectively avoids this issue by not directly predicting the keypoints, but instead refining the final output through 2D keypoint supervision applied to each joint of a template pose that is placed in the 3D space.

## 6 Ablation Study

In this section, we present ablation studies of our RAPTR on the HIBER dataset. Unless otherwise stated, all reported evaluation results are reported as the mean $\pm$ standard deviation, computed over three random seeds. Additional ablation results and visualizations are provided in Appendix I.

**Visualization of Pose Refinement Process:** Fig. 6 illustrates the refinement process of a 3D prediction through the two-stage decoder architecture. The pose decoder first establishes coarse 3D structures of the human body under the constraint of the 3D template loss (1st row). Subsequently, the joint decoder fine-tunes the keypoints to better capture the subject orientation and limb configuration (2nd row), while preserving the structure consistency provided by the pose decoder.

**Effect of Loss Terms:** Table 3 provides an ablation study on the effect of different combinations of loss terms in the two-stage decoder. When only the K2D loss is applied at the joint decoder (row 1), the 3D pose estimation suffers from depth ambiguity due to the absence of any 3D constraint, resulting in a substantial increase in MPJPE to $381.18\,\mathrm{cm}$ and $375.73\,\mathrm{cm}$ on the WALK and MULTI splits, respectively. From rows 2 to 4 in Table 3, we remove or modify one loss term at a time from
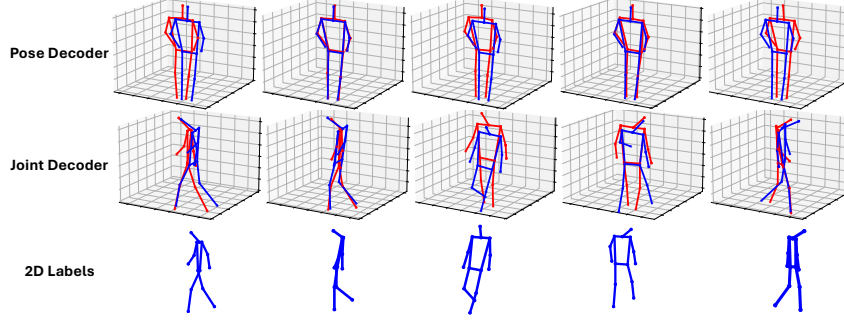
Figure 6: Pose Refinement Process: the pose prediction is first constrained by the 3D template at the pose decoder, and subsequently refined at the joint decoder.

the proposed structural loss. Removing the T3D loss at the pose decoder (row 2), replacing the T3D with the K2D+G3D loss at the pose decoder (row 3), or removing the G3D loss at the joint decoder (row 4) leads to a noticeable degradation in 3D pose estimation.

Table 3: Effect of loss terms for RAPTR (MPJPE: cm).

| Loss | | WALK | MULTI | Notes |
|---|---|---|---|---|
| Pose Dec. | Joint Dec. | | | |
| – | K2D | $381.18 \pm 0.28$ | $375.73 \pm 6.31$ | 2D keypoint loss only at joint decoder |
| – | K2D+G3D | $28.54 \pm 4.57$ | $57.90 \pm 9.81$ | without 3D template loss at pose decoder |
| K2D+G3D | K2D+G3D | $27.49 \pm 3.40$ | $23.43 \pm 3.44$ | with 2D keypoint + 3D gravity loss at both decoders |
| T3D | K2D | $25.96 \pm 4.95$ | $25.83 \pm 3.87$ | without 3D gravity loss at joint decoder |
| T3D | K2D+G3D | $\mathbf{22.32 \pm 0.06}$ | $\mathbf{18.99 \pm 0.16}$ | proposed structural loss |

K2D = 2D Keypoint loss, T3D = 3D Template loss, G3D = 3D Gravity loss

**Effect of Deformable Attention Mechanisms:** Table 4 presents an ablation study on the effect of the deformable attention mechanism for RAPTR. In this study, the pseudo-3D deformable attention is replaced with the decoupled 2D deformable attention used in QRFPose [33], while keeping the cross-view encoder, two-stage decoder architecture, and the proposed structural loss unchanged. The results show that the pseudo-3D deformable attention yields marginal performance improvements, approximately 4% and 2.5% on the WALK and MULTI splits, respectively.

Table 4: Effect of deformable attention mechanisms for RAPTR (MPJPE: cm).

| Attn. | WALK | MULTI | Notes |
|---|---|---|---|
| 2D | $23.25 \pm 1.38$ | $19.47 \pm 0.95$ | RAPTR with decoupled 2D deformable attention |
| 3D | $\mathbf{22.32 \pm 0.06}$ | $\mathbf{18.99 \pm 0.16}$ | RAPTR with pseudo-3D deformable attention |

**Comparison with a 2D-to-3D Pose Uplifting Model:** We further compare RAPTR with a baseline that first estimates 2D keypoints in the image plane and subsequently lifts them to 3D space using a pre-trained 2D-to-3D pose uplifting model [21] trained on vision-based datasets such as Human3.6M [13]. To ensure a fair comparison, this baseline adopts the same network architecture as RAPTR, but both the pose and joint decoders are supervised only by the 2D keypoint loss. Because the 3D poses predicted by the uplifting model are defined in a pelvis-centered coordinate system, we additionally estimate a translation offset to align the estimated poses with their correct position in the world coordinate system. As shown in Table 5, the pose uplifting baseline performs significantly worse than RAPTR, with MPJPEs of $43.43\,\mathrm{cm}$ and $41.76\,\mathrm{cm}$ on the WALK and MULTI splits, respectively.

**Limitation:** Given that the process of refining the template to the actual pose in the joint decoder is supervised by the 2D keypoint labels, the accuracy of the 3D pose estimation is highly dependent on the precision of the labels in the image plane. In this context, since the 2D keypoint label lacks the

Table 5: Comparison with a 2D-to-3D pose uplifting model (MPJPE: cm).

| | Loss | | WALK | MULTI |
|---|---|---|---|---|
| | Pose Dec. | Joint Dec. | | |
| Pose Lifting [21] | K2D | K2D | $43.43 \pm 2.66$ | $41.76 \pm 6.85$ |
| RAPTR (ours) | T3D | K2D+G3D | $\mathbf{22.32 \pm 0.06}$ | $\mathbf{18.99 \pm 0.16}$ |

ability to discern whether the person is facing forward or backward to the camera, the estimated 3D poses may have joints that are bent in the opposite direction in depth from the actual pose. In addition, real-world conditions such as occlusion and human-to-human interference can further degrade the pose estimation performance. These effects become more pronounced in crowded or interactive environments.

## 7  Conclusion

We introduced RAPTR, a radar-based 3D human pose estimation system using reliable 2D keypoint labels and 3D BBoxes as the coarse-grained 3D information. We designed the network architecture and the loss function to integrate multi-view radar features and consistently represent human poses in the 3D space, whose effectiveness was demonstrated through experimental results.

**Broader Impacts:**  Indoor radar perception technologies, such as RAPTR, provide diverse indoor applications. These technologies may improve the safety and energy efficiency of indoor systems while preserving privacy. However, it is paramount that the perception results remain secure and private to prevent misuse.

## References

[1] Sizhe An, Yin Li, and Umit Ogras. mRI: Multi-modal 3D human pose estimation dataset using mmWave, RGB-D, and inertial sensors. In *Advances in Neural Information Processing Systems*, pages 27414–27426, 2022.

[2] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2D pose estimation using part affinity fields. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7291–7299, 2017.

[3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, page 213–229, 2020.

[4] Ching-Hang Chen and Deva Ramanan. 3D human pose estimation = 2D pose estimation + matching. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5759–5767, 2017.

[5] Ching-Hang Chen, Ambrish Tyagi, Amit Agrawal, Dylan Drover, Rohith Mv, Stefan Stojanov, and James M Rehg. Unsupervised 3D pose estimation with geometric self-supervision. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5714–5724, 2019.

[6] Bowen Cheng, Bin Xiao, Jingdong Wang, Honghui Shi, Thomas S Huang, and Lei Zhang. HigherHRNet: Scale-aware representation learning for bottom-up human pose estimation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[7] Fangqiang Ding, Zhen Luo, Peijun Zhao, and Chris Xiaoxuan Lu. milliFlow: Scene flow estimation on mmWave radar point cloud for human motion sensing. In *European Conference on Computer Vision (ECCV)*, pages 202–221, 2024.

[8] Junqiao Fan, Jianfei Yang, Yuecong Xu, and Lihua Xie. Diffusion model is a good pose estimator from 3D RF-vision. In *European Conference on Computer Vision*, page 1–18, 2024.

[9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.

[10] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask R-CNN. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(2):386–397, 2020.

[11] Yuan-Hao Ho, Jen-Hao Cheng, Sheng Yao Kuan, Zhongyu Jiang, Wenhao Chai, Hsiang-Wei Huang, Chih-Lung Lin, and Jenq-Neng Hwang. RT-Pose: A 4D radar tensor-based 3D human pose estimation and localization benchmark. In *European Conference on Computer Vision*, page 107–125, 2024.

[12] Shuting Hu, Siyang Cao, Nima Toosizadeh, Jennifer Barton, Melvin G. Hector, and Mindy J. Fain. mmPose-FK: A forward kinematics approach to dynamic skeletal pose estimation using mmWave radars. *IEEE Sensors Journal*, 24(5):6469–6481, 2024.

[13] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6M: Large scale datasets and predictive methods for 3D human sensing in natural environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(7):1325–1339, 2014.

[14] Niraj Prakash Kini, Ruey-Horng Shiue, ryan chandra, Wen-Hsiao Peng, Ching-Wen Ma, and Jenq-Neng Hwang. TransHuPR: Cross-view fusion transformer for human pose estimation using mmWave radar. In *British Machine Vision Conference*, 2024.

[15] Shih-Po Lee, Niraj Prakash Kini, Wen-Hsiao Peng, Ching-Wen Ma, and Jenq-Neng Hwang. HuPR: A benchmark for human pose estimation using millimeter wave radar. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 5715–5724, 2023.

[16] Jiefeng Li, Siyuan Bian, Ailing Zeng, Can Wang, Bo Pang, Wentao Liu, and Cewu Lu. Human pose regression with residual log-likelihood estimation. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.

[17] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollar. Focal loss for dense object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(2):318–327, 2020.

[18] Chris Xiaoxuan Lu, Stefano Rosa, Peijun Zhao, Bing Wang, Changhao Chen, John A. Stankovic, Niki Trigoni, and Andrew Markham. See through smoke: robust indoor mapping with low-cost mmWave radar. In *International Conference on Mobile Systems, Applications, and Services (MobiSys)*, page 14–27, 2020.

[19] Zhengxiong Luo, Zhicheng Wang, Yan Huang, Liang Wang, Tieniu Tan, and Erjin Zhou. Rethinking the heatmap regression for bottom-up human pose estimation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13259–13268, 2021.

[20] Sebastian Lutz, Richard Blythman, Koustav Ghosal, Matthew Moynihan, Ciaran Simms, and Aljosa Smolic. Jointformer: Single-frame lifting transformer with error prediction and refinement for 3D human pose estimation. In *International Conference on Pattern Recognition (ICPR)*, pages 1156–1163, 2022.

[21] Julieta Martinez, Rayat Hossain, Javier Romero, and James J Little. A simple yet effective baseline for 3D human pose estimation. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2659–2668, 2017.

[22] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. *arXiv [cs.CV]*, 2016.

[23] Ashish Pandharipande, Chih-Hong Cheng, Justin Dauwels, Sevgi Z. Gurbuz, Javier Ibanez-Guzman, Guofa Li, Andrea Piazzoni, Pu Wang, and Avik Santra. Sensing and machine learning for automotive perception: A review. *IEEE Sensors Journal*, 23(11):11097–11115, 2023.

[24] M. Mahbubur Rahman, Ryoma Yataka, Sorachi Kato, Pu (Perry) Wang, Peizhao Li, Adriano Cardace, and Petros Boufounos. MMVR: Millimeter-wave multi-view radar dataset and benchmark for indoor perception. In *European Conference on Computer Vision (ECCV)*, pages 306–322, 2024.

[25] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(6):1137–1149, 2017.

[26] Avik Santra, Pu Wang, George Shaker, Bhavani Shankar Mysore, Guido Dolmans, Yan Chen, Negin Shariati, and Ashish Pandharipande. Machine learning-powered radio frequency sensing: A review. *IEEE Sensors Journal*, 25(13):23164–23183, 2025.

[27] Arindam Sengupta, Feng Jin, Renyuan Zhang, and Siyang Cao. mm-Pose: Real-time human skeletal posture estimation using mmWave radars and CNNs. *IEEE Sensors Journal*, 20(17):10032–10044, 2020.

[28] Dahu Shi, Xing Wei, Liangqi Li, Ye Ren, and Wenming Tan. End-to-end multi-person pose estimation with transformers. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11059–11068, 2022.

[29] Akash Deep Singh, Sandeep Singh Sandha, Luis Garcia, and Mani Srivastava. RadHAR: Human activity recognition from point clouds generated through a millimeter-wave radar. In *ACM Workshop on Millimeter-Wave Networks and Sensing Systems*, pages 51–56, 2019.

[30] Mikael Skog, Oleksandr Kotlyar, Vladimír Kubelka, and Martin Magnusson. Human detection from 4D radar data in low-visibility field conditions. *arXiv:2404.05307*, 2024.

[31] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5686–5696, 2019.

[32] Denis Tome, Chris Russell, and Lourdes Agapito. Lifting from the deep: Convolutional 3D pose estimation from a single image. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[33] Hong Wan, Ruiyuan Song, Chunyang Xie, Zhi Lu, Qi Chen, Zhi Wu, Dongheng Zhang, Yang Hu, and Yan Chen. QRFPose: Query-based 3D pose estimation using radio signals. In *International Conference on Wireless Communications and Signal Processing (WCSP)*, pages 1271–1277, 2024.

[34] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, Wenyu Liu, and Bin Xiao. Deep high-resolution representation learning for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(10):3349–3364, 2021.

[35] Zhi Wu, Dongheng Zhang, Chunyang Xie, Cong Yu, Jinbo Chen, Yang Hu, and Yan Chen. RFMask: A simple baseline for human silhouette segmentation with radio signals. *IEEE Transactions on Multimedia*, 25:4730–4741, 2023.

[36] Qian Xie, Qianyi Deng, Ta Ying Cheng, Peijun Zhao, Amir Patel, Niki Trigoni, and Andrew Markham. mmPoint: Dense human point cloud generation from mmWave. In *British Machine Vision Conference (BMVC)*, pages 194–196, 2023.

[37] Qian Xie, Xinyu Hou, Qianyi Deng, Amir Patel, Niki Trigoni, and Andrew Markham. mmDiffusion: mmWave diffusion for sequential 3d human dense point cloud generation. In *2025 International Conference on 3D Vision (3DV)*, pages 781–790, 2025.

[38] Hongfei Xue, Yan Ju, Chenglin Miao, Yijiang Wang, Shiyang Wang, Aidong Zhang, and Lu Su. mmMesh: towards 3D real-time dynamic human mesh construction using millimeter-wave. In *International Conference on Mobile Systems, Applications, and Services (MobiSys)*, page 269–282, 2021.

[39] Hongfei Xue, Qiming Cao, Yan Ju, Haochen Hu, Haoyu Wang, Aidong Zhang, and Lu Su. M$^4$esh: mmWave-based 3D human mesh construction for multiple subjects. In *Proceedings of the 20th ACM Conference on Embedded Networked Sensor Systems*, pages 391–406, 2022.

[40] Kangwei Yan, Fei Wang, Bo Qian, Han Ding, Jinsong Han, and Xing Wei. Person-in-WiFi 3D: End-to-end multi-person 3D pose estimation with Wi-Fi. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 969–978, 2024.

[41] Jianfei Yang, He Huang, Yunjiao Zhou, Xinyan Chen, Yuecong Xu, Shenghai Yuan, Han Zou, Chris Xiaoxuan Lu, and Lihua Xie. MM-Fi: Multi-modal non-intrusive 4D human dataset for versatile wireless sensing. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 18756–18768, 2023.

[42] Jiarui Yang, Songpengcheng Xia, Yifan Song, Qi Wu, and Ling Pei. mmBaT: A multi-task framework for mmWave-based human body reconstruction and translation prediction. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 8446–8450, 2024.

[43] Jiarui Yang, Songpengcheng Xia, Zengyuan Lai, Lan Sun, Qi Wu, Wenxian Yu, and Ling Pei. mmDEAR: mmWave point cloud density enhancement for accurate human body reconstruction. In *IEEE International Conference on Robotics and Automation*, pages 11227–11233, 2025.

[44] Shanliang Yao, Runwei Guan, Zitian Peng, Chenhang Xu, Yilu Shi, Weiping Ding, Eng Gee Lim, Yong Yue, Hyungjoon Seo, Ka Lok Man, Jieming Ma, Xiaohui Zhu, and Yutao Yue. Exploring radar data representations in autonomous driving: A comprehensive review. *arXiv:2312.04861*, 2024.

[45] Ryoma Yataka, Adriano Cardace, Pu (Perry) Wang, Petros Boufounos, and Ryuhei Takahashi. RETR: Multi-view radar detection transformer for indoor perception. *Neural Information Processing Systems*, 37: 19839–19869, 2024.

[46] Feng Zhang, Xiatian Zhu, Hanbin Dai, Mao Ye, and Ce Zhu. Distribution-aware coordinate representation for human pose estimation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7091–7100, 2020.

[47] Mingmin Zhao, Tianhong Li, Mohammad Abu Alsheikh, Yonglong Tian, Hang Zhao, Antonio Torralba, and Dina Katabi. Through-wall human pose estimation using radio signals. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7356–7365, 2018.

[48] Mingmin Zhao, Yonglong Tian, Hang Zhao, Mohammad Abu Alsheikh, Tianhong Li, Rumen Hristov, Zachary Kabelac, Dina Katabi, and Antonio Torralba. RF-based 3D skeletons. In *Conference of the ACM Special Interest Group on Data Communication (SIGCOMM)*, page 267–281, 2018.

[49] Peijun Zhao, Chris Xiaoxuan Lu, Bing Wang, Niki Trigoni, and Andrew Markham. CubeLearn: End-to-end learning for human motion recognition from raw mmWave radar signals. *IEEE Internet of Things Journal*, 10(12):10236–10249, 2023.

[50] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable DETR: Deformable transformers for end-to-end object detection. In *International Conference on Learning Representations*, 2021.

# A RAPTR Architecture

**Cross-view Encoder:** Fig. 7a illustrates the layer structure of the cross-view encoder in our RAPTR architecture. The cross-view encoder associates multi-view radar features using the cross-attention and skip-connection mechanism. A shared-weight backbone extracts sets of multi-scale features for both horizontal and vertical radar heatmaps as $\mathbf{Z}_{\text{hor}} = \{\mathbf{Z}_{\text{hor},i}\}_{i=1}^{S}$ and $\mathbf{Z}_{\text{ver}} = \{\mathbf{Z}_{\text{ver},i}\}_{i=1}^{S}$ with $S$ scale levels. Each scale-level feature map is enriched with spatial positional encoding and a learnable level embedding, following [50], and then fed into the cross-view encoder. The cross-view encoder consists of a stack of $L_{\text{enc}}$ multi-head multi-scale deformable attention layers, and we denote the input horizontal/vertical features to the $i$-th layer as $\mathbf{F}_{\text{hor}}^{(i)}, \mathbf{F}_{\text{ver}}^{(i)}$, respectively. In our case, $\mathbf{F}_{\text{hor}}^{(0)} = \mathbf{Z}_{\text{hor}}, \mathbf{F}_{\text{ver}}^{(0)} = \mathbf{Z}_{\text{ver}}$. Each $i$-th layer runs a shared cross-attention bidirectionally: first with $\mathbf{F}_{\text{hor}}^{(i-1)}$ as key/value and $\mathbf{F}_{\text{ver}}^{(i-1)}$ as query, then vice versa. Residual connection, layer normalization, and an FFN follow as in standard Transformer encoders to further refine the features, and additional residual connections are incorporated to preserve view-specific details. As a whole, the encoding process in the $i$-th layer is written as:

$$\bar{\mathbf{F}}_{\text{a}}^{(i-1)} = \mathbf{F}_{\text{a}}^{(i-1)} + \texttt{CrossAttn}(\mathbf{F}_{\text{a}}^{(i-1)}, \mathbf{F}_{\text{b}}^{(i-1)}),$$

$$\mathbf{F}_{\text{a}}^{(i)} = \mathbf{F}_{\text{a}}^{(i-1)} + \texttt{FFN}(\texttt{layernorm}(\bar{\mathbf{F}}_{\text{a}}^{(i-1)})), \quad (a,b) \in \{(\text{hor}, \text{ver}), (\text{hor}, \text{ver})\}. \tag{10}$$

The output of the last layer, or encoder memory, is denoted as $\mathbf{F}_{\text{hor}}, \mathbf{F}_{\text{ver}}$.

**Pose/Joint Decoder:** Fig. 7b illustrates the layer structure of the pose/joint decoder in RAPTR. The pose decoder and the joint decoder share the architecture: they receive pose/joint queries $\mathbf{Q}_{\text{pose}}, \mathbf{Q}_{\text{joint}}$, multi-scale encoder memory from the cross-view encoder $\mathbf{F}_{\text{hor}}, \mathbf{F}_{\text{ver}}$, and reference points $\tilde{\mathbf{P}}_{\text{radar}}$. The decoders then generate refined embeddings through multi-head self-attention and pseudo-3D deformable attention layers, which is the process denoted as $\mathcal{D}_{\text{pose}}$ and $\mathcal{D}_{\text{joint}}$ in Section 4. $N$ pose queries correspond to $N$ pose predictions in the pose decoder, whereas $K$ joint queries correspond to $K$ joints on the same subject in the joint decoder. The queries are first fed into self-attention, followed by residual connection and layer normalization, and then passed into the pseudo-3D deformable attention layer, as defined in Section 4.2. The pseudo-3D deformable attention layer is a cross-attention layer, using encoder memory to produce keys and values, which correlate with the refined queries. In addition, reference points are fed into this layer to determine the sampling locations and aggregate sparse features on the multi-view encoder memory across space and scales in the pseudo-3D deformable attention mechanism. The outputs are then passed through another residual connection, layer normalization, and an FFN. In the pose decoder, a class regression
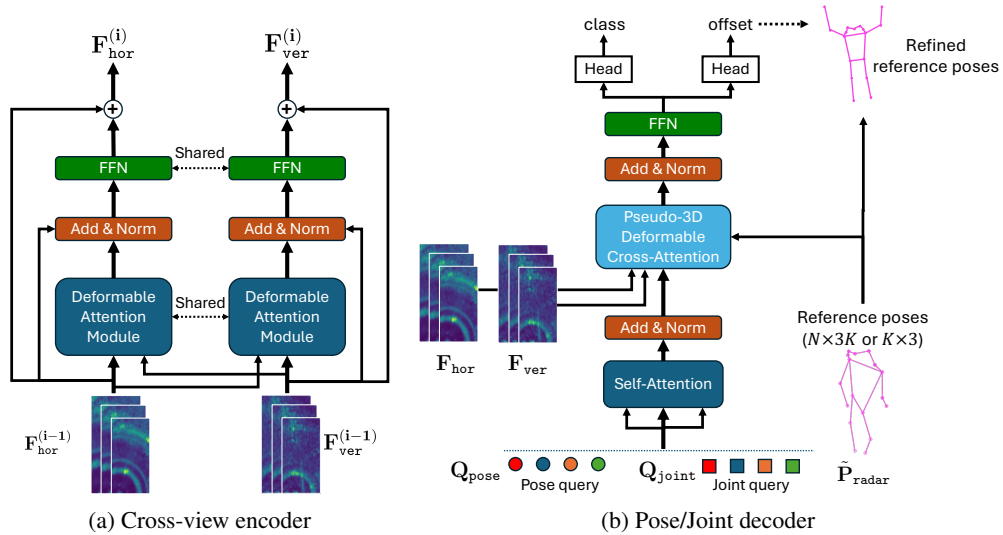


(a) Cross-view encoder       (b) Pose/Joint decoder

Figure 7: Transformer design in RAPTR.

head takes the resulting $N$ queries to calculate confidence scores $\hat{\mathbf{c}}$ for each corresponding person. A pose regression head also takes the resulting queries to calculate pose coordinate offsets $\mathbb{R}^{N \times 3K}$ to refine the reference poses. In the joint decoder, a pose regression head, shared with the pose decoder, takes the resulting $K$ queries to calculate the joint coordinate offsets $\mathbb{R}^{K \times 3}$ to refine the reference points. The pose decoder and the joint decoder consist of $L_{\texttt{pose}}, L_{\texttt{joint}}$ decoder layers, respectively, and their outputs are the initial pose estimates and the refined pose estimates, as shown in Fig. 3.

In the implementation, we have a technically involved step regarding the application of bipartite matching to the batched estimations from the pose decoder. We train the model using mini-batches, where the reference poses are first shaped as $B \times N \times 3K$, with $B$ denoting the mini-batch size. While Section 4.3 describes the bipartite matching in general terms, it is in practice applied to the initial pose estimations $\tilde{\mathbf{P}}_{\texttt{world}}$ output by the pose decoder and the corresponding 3D keypoint labels $\mathbf{P}_{\texttt{world}}$ so that the computational cost of the subsequent joint decoder would be reduced. The matching is guided by the regressed confidence scores $\hat{\mathbf{c}}$ and the Euclidean distance between the initial estimations and the labels. Out of all estimations in the mini-batch, we retain only the $N'$ matched ones and reshape them to $N' \times K \times 3$ to serve as the reference poses for the joint decoder, in which $N'$ is considered the new mini-batch size.

**Computational Complexity:** The cross-view encoder takes two sets of multi-scale feature maps $\mathbf{F}_{\texttt{hor}}, \mathbf{F}_{\texttt{ver}}$ for horizontal-depth and vertical-depth radar perceptions. As shown in Fig. 7a, we apply deformable cross-attention in both directions: from $\mathbf{F}_{\texttt{hor}}$ to $\mathbf{F}_{\texttt{ver}}$ and vice versa, treating one set as queries and the other as keys and values in each direction. For each direction, given $S$-level feature scales, each with $N_s$ spatial positions, and $N_{\texttt{offset}}$ sampling points per head, the total computational cost is

$$\mathcal{O}(2(d^2 + N_{\texttt{offset}}d) \sum_{s=1}^{S} N_s) \approx \mathcal{O}((d^2 + N_{\texttt{offset}}d) \sum_{s=1}^{S} N_s), \qquad (11)$$

where $d$ is the feature dimension.

The pose decoder takes $N$ object queries and the encoded memory and performs self-attention and pseudo-3D deformable cross-attention. Therefore, the computational cost is written as

$$\mathcal{O}(N^2 d + N d^2 + N N_{\texttt{offset}} S d). \qquad (12)$$

Here, we omit the constant factor associated with bilinear interpolation and regression of sampling offsets and attention weight matrices in the last term, as it does not affect the asymptotic complexity.

Finally, the joint decoder takes $K$ joint queries for $N'$ poses, selected through a bipartite matching procedure out of $N$, and the encoded memory and performs pseudo-3D deformable attention as well, and thus the cost is

$$\mathcal{O}((N'K)^2 d + (N'K)d^2 + N'K N_{\texttt{offset}} S d). \qquad (13)$$

In conclusion, the total computational cost of our RAPTR is

$$\mathcal{O}((d^2 + N_{\texttt{offset}}d) \sum_{s=1}^{S} N_s + (N^2 + (N'K)^2)d + (N + N'K)d^2 + (N + N'K)N_{\texttt{offset}} S d). \qquad (14)$$

## B  Details of Pseudo-3D Deformable Attention

**Multi-scale Multi-head Extension of Pseudo-3D Deformable Attention:** We can extend the pseudo-3D deformable attention defined by Eq. 8 in Section 4.2 to multi-scale and multi-head operation. First, given $M$ heads and $S$ feature scale levels, Eq. 7 is extended as

$$\mathbf{f}_{ms,\texttt{hor}}^{(i)} = \mathbf{F}_{s,\texttt{hor}}(x + \Delta x_{msi}, z + \Delta z_{msi}), \quad \mathbf{f}_{ms,\texttt{ver}}^{(i)} = \mathbf{F}_{s,\texttt{ver}}(y + \Delta y_{msi}, z + \Delta z_{msi}), \qquad (15)$$

where $\mathbf{F}_{s,\texttt{hor}}, \mathbf{F}_{s,\texttt{ver}}$ are the $s$-th level feature maps for horizontal and vertical view, respectively, and $\Delta\{x, y, z\}_{msi}$ is the $i$-th sampling offset in the $m$-th head on the $s$-th level feature. Subsequently, we collect the sampled features as $\mathbf{F}_{ms,\texttt{attn}} = \{\mathbf{f}_{ms,\texttt{hor}}^{(1)}, \mathbf{f}_{ms,\texttt{ver}}^{(1)}, \cdots, \mathbf{f}_{ms,\texttt{hor}}^{(N_{\texttt{offset}})}, \mathbf{f}_{ms,\texttt{ver}}^{(N_{\texttt{offset}})}\}$, and then extend Eq. 8 as

$$\bar{\mathbf{F}}_{\texttt{attn}} = \sum_{m=1}^{M} \mathbf{W}_m [\sum_{s=1}^{S} \sum_{i=1}^{N_{\texttt{offset}}} (A_{msi,0} \mathbf{W}'_m \mathbf{F}_{ms,\texttt{attn}}^{(2i-1)} + A_{msi,1} \mathbf{W}'_m \mathbf{F}_{ms,\texttt{attn}}^{(2i)})], \qquad (16)$$

15

(a) Pseudo-3D Deformable Attention in RAPTR

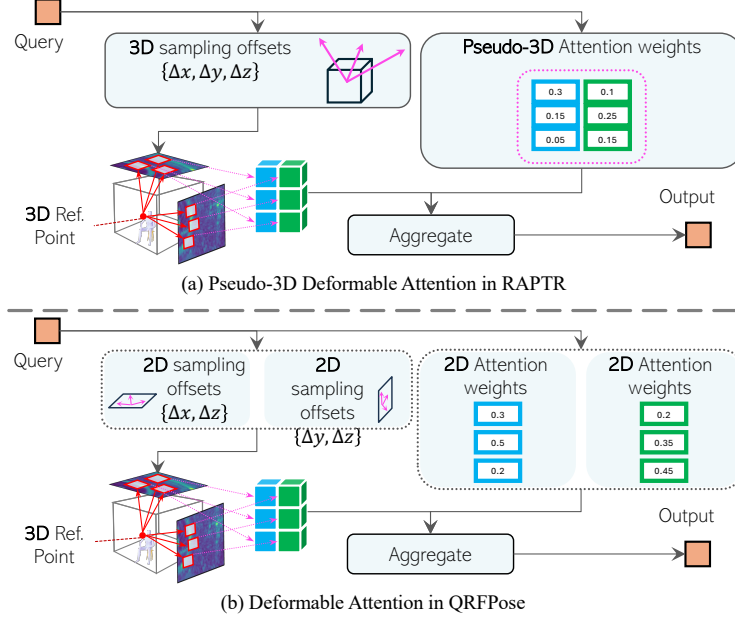(b) Deformable Attention in QRFPose

Figure 8: Comparison between (a) pseudo-3D and (b) decoupled 2D deformable attention mechanisms.

where $\mathbf{W}_m \in \mathbb{R}^{d \times d_v}$ and $\mathbf{W}'_m \in \mathbb{R}^{d_v \times d}$ with $d_v = d/M$ are learnable weight matrices, and $A_{msi,0}$ and $A_{msi,1}$ are the attention weights of the $i$-th sampled deformable radar feature on the $s$-th level feature in the $m$-th attention head, for the horizontal and vertical radar views. Attention weights are normalized per head by $\sum_s^S \sum_i^{N_{\texttt{offset}}} A_{msi} = 1$.

**Implementation of Sampling Location Determination:** We describe the pseudo-3D deformable attention in Section 4.2 as computing sampling locations by adding $N_{\texttt{offset}}$ offsets, which are derived from each query, to a corresponding reference point. From an implementation point of view, our pose decoder adopts a structurally extended design. In our implementation, each query corresponds not to a single keypoint, but to the entire pose of a person, represented as $K \times 3$ coordinates. Accordingly, we replace the notion of a reference point with a reference *pose*, formatted as an $N \times 3K$ tensor for $N$ queries. Sampling locations are determined by calculating the $3K$ offsets per query and adding them element-wise to the corresponding reference pose $K \times 3$. In this way, we virtually have $K$ sampling points for each pose. This effectively enables a single query to take care of $K$ spatial locations, allowing feature aggregation within the context of a unified pose.

On the other hand, the implementation in the joint decoder more aligns with the description in Section 4.2: a joint query in the joint decoder corresponds to a single joint so that the sampling locations are determined by calculating the $N_{\texttt{offset}}$ sampling offsets per query and adding them to the corresponding reference joint. We set $N_{\texttt{offset}} = 4$ as listed in Table 7.

**Computational Complexity Comparison with Decoupled 2D Deformable Attention:** As illustrated in Fig. 8 (a), the pseudo-3D deformable attention adapts a 3D reference point with 3D sampling offsets and the pseudo-3D attention weights are computed over multiple radar views. In comparison, the QRFPose of Fig. 8 (b) adapts a 3D reference point with projected 2D sampling offsets and 2D attention weights are separately computed over each radar view [33]. This simple lifting operation may lead to better computational complexity of the pseduo-3D attention over the number of radar views. In the following, we provide a computational complexity analysis for the two types of deformable attention mechanisms, given $V$ radar views:

- Decoupled 2D deformable attention: $\mathcal{O}(8VNN_{\texttt{offset}}d)$, where
  - 3D reference point projected to $V$ 2D radar views: $\mathcal{O}(6VN)$,
  - Offset estimation: $\mathcal{O}(2VNN_{\texttt{offset}}d)$, where 2 is due to the computation of 2D $(x,y)$ offsets,

16

Table 6: Complexity comparison of pseudo-3D vs. decoupled 2D deformable attention.

| Queries ($N$) | Views ($V$) | 2D Att | Pseudo-3D Att | Ratio (3D/2D) | Savings |
|---|---|---|---|---|---|
| 10 | 2 | $160NC$ | $150NC$ | **0.94↓** | **6.25%** |
| 10 | 5 | $400NC$ | $330NC$ | **0.83↓** | **17.5%** |
| 10 | 10 | $800NC$ | $630NC$ | **0.79↓** | **21.3%** |

  – Attention weights: $\mathcal{O}(VNN_{\texttt{offset}}d)$,
  – Feature aggregation: $\mathcal{O}(5VNN_{\texttt{offset}}d)$, where 5 is due to bilinear interpolation and weighted sum,

- Pseudo-3D deformable attention: $\mathcal{O}(6VNN_{\texttt{offset}}d + 3NN_{\texttt{offset}}d)$, where

  – Offset estimation: $\mathcal{O}(3NN_{\texttt{offset}}d)$, where 3 is due to the computation of the 3D $(x, y, z)$ offsets,
  – 3D offset projected to $V$ 2D radar views: $\mathcal{O}(6VN)$,
  – Attention weights: $\mathcal{O}(VNN_{\texttt{offset}}d)$,
  – Feature aggregation: $\mathcal{O}(5VNN_{\texttt{offset}}d)$.

Note that, in the above analysis, $\mathcal{O}(6VN)$ is excluded from the final complexity expressions as $6VN \ll 5VNN_{\texttt{offset}}d$ in practice. Table 6 shows the complexity comparison with a specific number of queries and an increasing number of views. It is observed that the pseudo-3D attention achieves computational savings of $17.5\%$ with $V = 5$ radar views and $21.3\%$ with $V = 10$ radar views, compared to the decoupled 2D attention.

## C   Optional View Mask Module

As an extension of our pseudo-3D deformable attention, we introduce an optional view mask module that aims to put more attention weights on the feature on the more important view with a hard-thresholding approach. The view mask module first computes a view selection mask $\mathbf{M}_{\texttt{attn}}$ from a query $\mathbf{q}$ corresponding to reference points of interest, as:

$$\mathbf{M}_{\texttt{attn}} = \sigma(\lambda \cdot \texttt{FFN}(\mathbf{q})) \in \mathbb{R}^{N_{\texttt{offset}} \times 2}, \qquad (17)$$

where $\sigma$ is the Sigmoid function, and $\lambda(\approx 1e5) \in \mathbb{R}$ makes the sigmoid output very close to 0 or 1 while preserving gradients. The element $m_{i,j}(\approx 1)$ signals that the $j$th view ($j = 0$ to the horizontal and $j = 1$ to the vertical) should retain its share of attention at the $i$th sampling point, whereas $m_{i,j}(\approx 0)$ marks it for suppression.

For each $i$th sampling point, we can consider three patterns to adjust the attention weights $f_{\texttt{attn}}$ for that point $A_{i,0}, A_{i,1}$ and, potentially, other sampling points according to the corresponding row in the view selection mask $[m_{i,0}, m_{i,1}]$.

- $[m_{i,0}, m_{i,1}] = [1, 1]$: use both views, weights unchanged;

- $[m_{i,0}, m_{i,1}] = [0, 1]$ or $[m_{i,0}, m_{i,1}] = [1, 0]$: ignore one view and transfer its weight to the other, e.g., when $[m_{i,0}, m_{i,1}] = [0, 1]$, the adjusted attention weights are $\hat{A}_{i,0} = A_{i,0} + A_{i,1}, \hat{A}_{i,1} = 0$, to to ensure that the view selection decision at one sampling point does not influence the other sampling points;

- $[m_{i,0}, m_{i,1}] = [0, 0]$: ignore this sampling point in both views. Its weight is evenly redistributed so that $\sum_{i,j} \hat{A}_{i,j} = 1$ still holds.

In this way, the view selection mask adaptively changes through training so that the attention mechanism associates with the view that is more informative for each sampling point. We can also utilize this view selection module to control the use of multi-view features by manually setting the values in the view selection mask, which we conduct a relating ablation study in the following Appendix I.
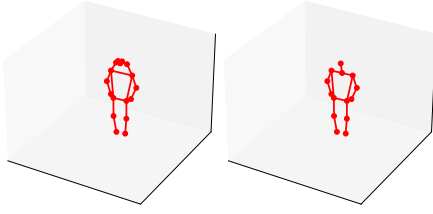
Figure 9: Template keypoints $\mathbf{K}_{\texttt{world}}$ for MMVR (left) and HIBER (right) dataset.
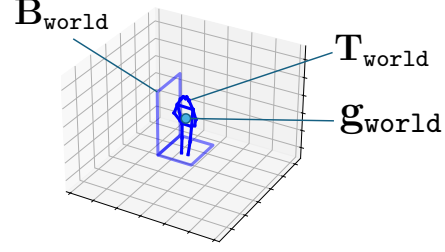
Figure 10: An example of a template pose located in the 3D world coordinate.

## D  Details of Loss Functions

In this section, we give a supplementary explanation of the loss functions in RAPTR provided in Section 4. For simplicity, we consider a single corresponding sample from each set, denoted as $\{\mathbf{b}_{\texttt{world}} \in \mathbf{B}_{\texttt{world}}, \mathbf{p}_{\texttt{image}} \in \mathbf{P}_{\texttt{image}}, \hat{\mathbf{p}}_{\texttt{image}} \in \hat{\mathbf{P}}_{\texttt{image}}, \tilde{\mathbf{p}}_{\texttt{world}} \in \tilde{\mathbf{P}}_{\texttt{world}}, \hat{\mathbf{p}}_{\texttt{world}} \in \hat{\mathbf{P}}_{\texttt{world}}\}$.

**Coarse-grained 3D Loss:**  For 3D gravity loss, we compute the centroid of a 3D BBox label in $\mathbf{b}_{\texttt{world}} = [x_{\texttt{min}}, y_{\texttt{min}}, z_{\texttt{min}}, x_{\texttt{max}}, y_{\texttt{max}}, z_{\texttt{max}}]$ as the 3D gravity center label $\mathbf{g}_{\texttt{world}} \in \mathbb{R}^{1 \times 3}$ as

$$\mathbf{g}_{\texttt{world}} = \big[\frac{x_{\texttt{max}} - x_{\texttt{min}}}{2}, \frac{y_{\texttt{max}} - y_{\texttt{min}}}{2}, \frac{z_{\texttt{max}} - z_{\texttt{min}}}{2}\big]. \tag{18}$$

Given a refined 3D pose estimate $\hat{\mathbf{p}}_{\texttt{world}} = \{(\hat{\mathbf{p}}_{\texttt{world},x}^{(k)}, \hat{\mathbf{p}}_{\texttt{world},y}^{(k)}, \hat{\mathbf{p}}_{\texttt{world},z}^{(k)})\}_{k=1}^{K}$ at the joint decoder, we also collapse it into its centroids as $\hat{\mathbf{g}}_{\texttt{world}} \in \mathbb{R}^{1 \times 3}$ as

$$\hat{\mathbf{g}}_{\texttt{world}} = \big[\frac{1}{K}\sum_k \hat{\mathbf{p}}_{\texttt{world},x}^{(k)}, \frac{1}{K}\sum_k \hat{\mathbf{p}}_{\texttt{world},y}^{(k)}, \frac{1}{K}\sum_k \hat{\mathbf{p}}_{\texttt{world},z}^{(k)}\big], \tag{19}$$

where $\hat{\mathbf{p}}_{\texttt{world},x/y/z}^{(k)}$ is the $x$-, $y$-, and $z$-coordinate for the $k$-the joint. The 3D gravity loss $\mathcal{L}_{\texttt{gravity}}$ is defined as the Euclidean distance between the two 3D gravity centers, $\mathbf{g}_{\texttt{world}}$ and $\hat{\mathbf{g}}_{\texttt{world}}$.

For 3D template loss, we construct a 3D template keypoint label for each $\mathbf{b}_{\texttt{world}}$ using template keypoints at the coordinate origin $\mathbf{K}_{\texttt{world}} \in \mathbb{R}^{K \times 3}$. The template pose $\mathbf{T}_{\texttt{world}} \in \mathbb{R}^{K \times 3}$ is given as $\mathbf{T}_{\texttt{world}} = \mathbf{K}_{\texttt{world}} + \mathbf{1}^\top \mathbf{g}_{\texttt{world}}$. Fig. 9 shows the template keypoints $\mathbf{K}_{\texttt{world}}$ with different numbers of keypoints for the MMVR and HIBER datasets, and Fig. 10 provides an example of locating these template keypoints in the 3D world coordinate based on the location of a 3D BBox label. The 3D template loss $\mathcal{L}_{\texttt{template}}$ is defined as the Euclidean distance between the template poses $\mathbf{T}_{\texttt{world}}$ and the initial 3D pose estimates $\tilde{\mathbf{P}}_{\texttt{world}}$ from the pose decoder.

**Fine-grained 2D Loss:**  Specifically in the fine-grained 2D loss, OKS loss $\mathcal{L}_{\texttt{OKS}}$ is the loss function based on object keypoint similarity (OKS), a metric used to evaluate the accuracy of keypoint estimations taking into account the object scale and keypoint visibility, which is defined as

$$\text{OKS}(\mathbf{p}_{\texttt{image}}, \hat{\mathbf{p}}_{\texttt{image}}) = \sum_k \exp(-\frac{d_k^2}{2s^2\psi_k^2}), \tag{20}$$

where $d_k$ is the distance between the $k$-th estimated joint $\hat{\mathbf{p}}_{\texttt{image}}^{(k)}$ and the corresponding label $\mathbf{p}_{\texttt{image}}^{(k)}$, $s$ is the object scale, and $\psi_k$ is a pre-defined constant for the $k$-th joint. Here, we assume that all keypoint labels are annotated as visible points. Since OKS is a metric in which higher values indicate greater similarity, its negative logarithm $-\log(\text{OKS})$ is taken when used as a loss function $\mathcal{L}_{\texttt{OKS}}$.

## E  Datasets

**HIBER:**  HIBER [35] is an open-source multi-view radar dataset for indoor human perception tasks including detection, segmentation, and keypoint estimation. They provide horizontal and vertical radar heatmaps and corresponding labels such as 2D BBoxes, 2D segmentation masks, 2D

Table 7: Hyper parameters for RAPTR

| Name | Notation | Value | |
|---|---|---|---|
| | | **HIBER** | **MMVR** |
| **Data** | | | |
| Radar image resolution | $W, H, D$ | 160, 160, 200 | 128, 128, 256 |
| # of training samples | - | 59000 / 54280 (WALK / MULTI) | 86579 / 190441(P1S1 / P2S1) |
| # of validation samples | - | 6490 / 5900 (WALK / MULTI) | 10538 / 23899 (P1S1 / P2S1) |
| # of test samples | - | 3540 / 3540 (WALK / MULTI) | 10785 / 23458 (P1S1 / P2S1) |
| # of keypoints | $K$ | 14 | 17 |
| **Model params** | | | |
| Backbone | - | ResNet 18 | |
| # of feature scale | $S$ | 3 | |
| Feedforward dimension in Transformer | - | 1024 | |
| # of encoder layers | $L_{\texttt{enc}}$ | 3 | |
| # of pose decoder layers | $L_{\texttt{pose}}$ | 2 | |
| # of joint decoder layers | $L_{\texttt{joint}}$ | 3 | |
| # of deformable sampling offsets | $N_{\texttt{offset}}$ | $K$ / 4 (pose / joint decoder) | |
| # of heads in multi-head attention | - | 8 | |
| Feature dimension | $d$ | 128 | |
| # of input frames | $T$ | 4 | |
| # of pose query | $N$ | 10 | |
| **Training params** | | | |
| Optimizer | - | AdamW | |
| Base learning rate | - | 2e-4 | |
| Weight decay | - | 1e-4 | |
| LR scheduler | - | Cosine Decay | |
| Batch size | - | 32 | |
| Epochs | - | 50 | |
| Gradient clip norm | - | 0.1 | |
| Early stopping patience | - | 5 epochs | |
| **Loss weights** | | | |
| 3D template loss | $\lambda_1$ | 1.0 | |
| 3D gravity loss | $\lambda_2$ | 1.0 | |
| 2D Keypoint loss | $\lambda_3$ | 5.0 | |
| 2D OKS loss | $\lambda_4$ | 1.0 | |
| Class loss | $\lambda_5$ | 1.0 | |
| **Computational Resource** | | | |
| GPU | | NVIDIA A40 | |
| # of workers | | 8 | |
| Approximate training time | | 3 hours / 10K samples | |

keypoints, 3D BBoxes, and 3D keypoints. Among its data environments, WALK and MULTI are currently available. WALK comprises frames that feature a single individual, while frames in MULTI consistently depict two individuals walking concurrently. The frames provided are captured from ten distinct viewpoints within a single room, designated as "view01" through "view10." We use "view02" to "view10" for training, validating, and testing the models, with the data splits provided, and the specific number of frames is listed in Table 7.

**MMVR:** MMVR [24] is a more recent open-source multi-view radar dataset for indoor human perception. They provide horizontal and vertical radar heatmaps and corresponding labels, such as 2D BBoxes, 2D segmentation masks, 2D keypoints, and 3D BBoxes. They collected data from 25 subjects in 6 different scenarios (e.g., open/cluttered office spaces) spanning over 9 days. MMVR consists of 1) P1: single-person scenarios in an open space without any obstacles, and 2) P2: multi-person scenarios in a cluttered office spaces, including sitting postures. P1 is designed to establish fundamental benchmarks for radar-based human perception tasks, while P2 is designed to challenge with more realistic and complicated indoor scenarios and cross-environment, cross-subject generalization. The data split we use is S1 that they provide, and the specific number of frames is listed in Table 7.

# F   Hyper Parameters

In the evaluations presented in Section 5, the model training for all baselines and our RAPTR and its variants share the hyper parameters outlined in Table 7, unless otherwise specified.

# G  Baseline Implementations

We provide the specific implementation for the baseline methods that we use in the evaluation.

**Person-in-WiFi 3D:**  We refer to the official implementation [40] and modify some parts of the code to make them compatible with the datasets that we use. We employ a ResNet backbone to extract multi-scale features from the radar heatmaps, as well as our RAPTR does. We then take the C4 feature map, flatten it as the $N_{\texttt{token}}$ tokens with dimensions of $d$, and feed it into the network. Specifically, $N_{\texttt{token}} = 260$ for the HIBER dataset and 256 for the MMVR dataset. Since the original study uses Wi-Fi channel state information (CSI), which inherently lacks explicit spatial structure, and transforms it into 180 tokens for input, our approach of converting C4 feature maps derived from radar heatmaps into approximately 200 tokens can be reasonably justified in terms of fairness and comparability. Regarding the loss function, we implement the loss as the summation of class loss, 2D keypoint loss, refined 2D keypoint loss, and 3D gravity loss, with loss weights of 1.0, 5.0, 10.0, and 1.0, respectively.

**HRRadarPose:**  We refer to the official implementation [11] and modify some parts of the code to make them compatible with the datasets that we use. First, we exclusively utilize horizontal view radar heatmaps, excluding vertical view heatmaps. This decision stems from the disparity in angular resolution between the elevation and azimuth axes reported in the HRRadarPose paper. The resolution of the elevation axis is 18 degrees, while the resolution of the azimuth axis is 1.4 degrees, and only the azimuth resolution is comparable to that of HIBER and MMVR (1.3 degrees). We presume that we could solely use the horizontal view heatmaps while ensuring fairness in our evaluations. In addition, we expand the original codes to multi-person scenarios by implementing Non-Maximum Suppression (NMS) on the predictions. Regarding the loss function, we implement the loss as the summation of heatmap loss, 2D keypoint loss, and 3D gravity loss, with loss weights of 5.0, 1.0, and 1.0, respectively.

**QRFPose:**  Currently, the authors of the paper have not released official codes. Therefore, we independently replicated the implementation based on the architectures and parameters outlined in the paper. We verify that our implementation replicates performance similar to that of the original report using 3D keypoint labels. To ensure a fair comparison, we set the number of Transformer decoder layers to 5, which is equivalent to the total number of layers in the pose decoder and the joint decoder in our RAPTR model. Although the original implementation uses RLE loss [16] as the keypoint regression loss function, we employ the conventional Euclidean distance loss in our implementation so that we can integrate the loss with the 3D gravity loss in a more balanced way. Regarding the loss function, we implement the loss as the summation of class loss, 2D keypoint loss, and 3D gravity loss, with loss weights of 1.0, 5.0, and 1.0, respectively.

# H  Metrics

For simplicity, we omit the subscripts that indicate the coordinate system in which the keypoints or BBoxes are defined (`radar`, `world`) in this section. In addition, $\mathbf{p} \in \mathbf{P}$, $\mathbf{b} \in \mathbf{B}$, and $\hat{\mathbf{p}} \in \hat{\mathbf{P}}$ denote the corresponding samples taken from the 3D keypoint labels, the 3D BBox labels, and the 3D pose estimates, respectively.

**MPJPE:**  We employ Mean Per Joint Position Error (MPJPE) as the performance metric to evaluate the 3D pose estimation capabilities of the models. Given a 3D keypoint label $\mathbf{p} = \{\mathbf{p}^{(k)}|(x^{(k)}, y^{(k)}, z^{(k)})\}_{k=1}^{K}$ and the corresponding estimate $\hat{\mathbf{p}} = \{\hat{\mathbf{p}}^{(k)}|(\hat{x}^{(k)}, \hat{y}^{(k)}, \hat{z}^{(k)})\}_{k=1}^{K}$. MPJPE is defined as:

$$\text{MPJPE} = \frac{1}{K} \sum_{k=1}^{K} \|\mathbf{p}^{(k)} - \hat{\mathbf{p}}^{(k)}\|_2. \tag{21}$$

The unit for MPJPE that we use is the centimeter in the world coordinate system. We also evaluate MPJPE along each axis: horizontal (h), vertical (v), and depth (d).

**3D BBox-based Metrics for MMVR:** For the evaluation of the MMVR dataset, due to the absence of 3D keypoint labels in the dataset, we approximate the pose estimation performance in a different way from that for the HIBER dataset. Specifically, we calculate 1) the distance between the center of the 3D pose estimate $\hat{\mathbf{P}}$ and the 3D BBox label $\mathbf{B}$, and 2) the absolute error in the edge lengths along each axis of the box. Specifically, we first construct a 3D BBox that encloses the estimated 3D keypoints as $\hat{\mathbf{b}} = [\min(\mathbf{p}_x), \min(\mathbf{p}_y), \min(\mathbf{p}_z), \max(\mathbf{p}_x), \max(\mathbf{p}_y), \max(\mathbf{p}_z)]$ where $\mathbf{p}_x, \mathbf{p}_y, \mathbf{p}_z$ is the set of $x$-, $y$- and $z$-coordinates of the estimated keypoints. We then calculate the center coordinate of the 3D BBox label $\mathbf{b} = [x_{\texttt{min}}, y_{\texttt{min}}, z_{\texttt{min}}, x_{\texttt{max}}, y_{\texttt{max}}, z_{\texttt{max}}]$ and $\hat{\mathbf{b}}$ as

$$\mathbf{g} = [\frac{x_{\texttt{max}} - x_{\texttt{min}}}{2}, \frac{y_{\texttt{max}} - y_{\texttt{min}}}{2}, \frac{z_{\texttt{max}} - z_{\texttt{min}}}{2}],$$
$$\hat{\mathbf{g}} = [\frac{\max(\mathbf{p}_x) - \min(\mathbf{p}_x)}{2}, \frac{\max(\mathbf{p}_y) - \min(\mathbf{p}_y)}{2}, \frac{\max(\mathbf{p}_z) - \min(\mathbf{p}_z)}{2}]. \tag{22}$$

The center distance between the BBoxes is the Euclidean distance between $\mathbf{g}$ and $\hat{\mathbf{g}}$. We also calculate the edge lengths of $\mathbf{b}$ and $\hat{\mathbf{b}}$ as

$$\mathbf{l} = (x_{\texttt{max}} - x_{\texttt{min}}, y_{\texttt{max}} - y_{\texttt{min}}, z_{\texttt{max}} - z_{\texttt{min}}),$$
$$\hat{\mathbf{l}} = (\max(\mathbf{p}_x) - \min(\mathbf{p}_x), \max(\mathbf{p}_y) - \min(\mathbf{p}_y), \max(\mathbf{p}_z) - \min(\mathbf{p}_z)), \tag{23}$$

and we calculate the absolute error of the edge length along each axis.

# I  Additional Ablation Studies and Visualization

To validate the effectiveness of RAPTR, we conduct additional ablation studies. Unless otherwise specified, we conduct the studies with the hyper parameters in Table 7.

## I.1  Numerical Results

**Additional Results for MMVR on P2S1:** Table 8 shows the evaluation results for the MMVR dataset on P2S1. P2S1 includes cluttered indoor scenarios with multiple subjects, which is thus more challenging than P1S1. RAPTR outperforms baselines and shows improvements in center distance by 71.54%, 85.28%, and 69.47% compared to Person-in-WiFi 3D, QRFPose, and HRRadarPose, respectively. We defer the qualitative evaluation for MMVR P2S1 to Appendix I.2.

Table 8: 3D pose estimation performance on MMVR (P2S1).

| Method | Center distance (cm) | Edge length error (cm) | | |
|---|---|---|---|---|
| | | (h) | (v) | (d) |
| Person-in-WiFi 3D | 103.43 | 48.29 | 112.75 | 152.88 |
| QRFPose | 200.03 | 115.80 | 126.14 | 335.30 |
| HRRadarPose | 96.43 | 32.19 | 51.02 | 175.04 |
| RAPTR (ours) | **29.44** | **18.74** | **27.14** | **40.29** |

**Full 3D Supervision:** We compare the performance of our RAPTR under (i) full supervision with fine-grained 3D keypoint labels and (ii) weak supervision with 2D keypoint and 3D BBox labels for HIBER (MULTI). Table 9 shows the performance comparison in MPJPE. Under full 3D supervision, the RAPTR architecture achieves an MPJPE of 8.93 cm. Even when trained under weak supervision, MPJPE increases by only about 10 cm, indicating that our structured loss design with two-stage decoding approach effectively learns reliable 3D body structures.

Table 9: RAPTR performance with full and weak supervision (HIBER MULTI).

| Method | Head | Neck | Shoulder | Elbow | Wrist | Hip | Knee | Ankle | MPJPE | (h) | (v) | (d) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| RAPTR (Full 3D supervision) | 7.90 | 5.28 | 6.41 | 7.94 | 12.43 | 5.75 | 8.80 | 14.65 | 8.93 | 4.97 | 2.86 | 5.12 |
| RAPTR (Weak supervision) | 18.39 | 13.13 | 16.44 | 20.12 | 24.62 | 15.01 | 17.76 | 23.22 | 18.99 | 7.80 | 4.38 | 14.54 |

Table 10: Effect of 3D Templates and their scales on the performance (HIBER MULTI).

| 3D Template | MPJPE |
|---|---|
| Standing (scale=1) | $18.99 \pm 0.16$ |
| Standing (scale=0.5) | $20.11 \pm 0.54$ |
| Sitting (scale=1) | $20.84 \pm 0.91$ |
| Standing (learned scale) | $23.13 \pm 0.33$ |

**Effect of 3D Templates and Their Scales:**    We evaluate the impact of 3D templates and their scale on the final MPJPE performance. Specifically, we experiment with

- A **standing pose** scaled by two factors: $0.5\times$ and $1\times$,
- A **sitting pose** of a $1.6\,\mathrm{m}$-tall person,
- A **learned scaling factor** applied to the standard standing pose.

Table 10 suggests that the choice of 3D template has minor impacts on the final MPJPE, likely due to the refinement capability of the second-stage joint decoder, as long as the first-stage decoder generates a reasonable, human-like initial pose.

**Effect of Loss Weighting Factors:**    We assess the RAPTR performance under varying loss weighting factors $\lambda_i$. Three configurations are evaluated: 1) equal weights for all loss terms, 2) increased weights on 3D losse terms, and 3) increased weights on 2D keypoint loss. Table 11 provides the performance comparison among these settings.

When all weighting factors are set to $1.0$, RAPTR achieves an average MPJPE o $19.91\,\mathrm{cm}$. Increasing the weights of the 3D losses to $5.0$ degrades performance, resulting in an MPJPE of $24.04\,\mathrm{cm}$. In contrast, emphasizing the 2D keypoint loss yields the best performance with an MPJPE of $18.99\,\mathrm{cm}$. These results suggest that appropriately balancing the loss terms, particularly by increasing the weight of the 2D keypoint loss, plays a crucial role in enhancing joint localization accuracy.

Table 11: Effect of loss weighting factors on the RAPTR performance (HIBER MULTI).

| $\lambda_1$ 3D template | $\lambda_2$ 3D gravity | $\lambda_3$ 2D keypoint | $\lambda_4$ 2D OKS | $\lambda_5$ class | MPJPE |
|---|---|---|---|---|---|
| 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | $19.91 \pm 0.65$ |
| 5.0 | 5.0 | 1.0 | 1.0 | 1.0 | $24.04 \pm 0.75$ |
| 1.0 | 1.0 | 5.0 | 1.0 | 1.0 | $\mathbf{18.99 \pm 0.16}$ |

**Impact of View Selection:**    To investigate how the use of features from horizontal and vertical views affects performance, we evaluate the performance of 3D pose estimation by configuring the view selection mask $\mathbf{M}_{\mathtt{attn}} \in \mathbb{R}^{N_{\mathtt{offset}} \times 2}$, integrated into pseudo-3D deformable attention as described in Appendix C, according to several predefined patterns. Specifically, given that the first column of the view selection mask $\mathbf{M}_{\mathtt{attn}}[:, 0]$ and the second column $\mathbf{M}_{\mathtt{attn}}[:, 1]$ correspond to horizontal and vertical views, respectively, we set the values in the mask as

- **Both Views**: all values in $\mathbf{M}_{\mathtt{attn}}$ to 1 so that the attention weight matrix $f_{\mathtt{attn}}$ is used as is;
- **Horizontal View Only**: $\mathbf{M}_{\mathtt{attn}}[:, 0] = \mathbf{1}$ and $\mathbf{M}_{\mathtt{attn}}[:, 1] = \mathbf{0}$ so that the features sampled from the horizontal feature map are aggregated and those of the vertical feature are omitted;
- **Vertical View Only**: The reversed one of "horizontal only": $\mathbf{M}_{\mathtt{attn}}[:, 0] = \mathbf{0}$ and $\mathbf{M}_{\mathtt{attn}}[:, 1] = \mathbf{1}$;
- **Random Mask**: The values in the mask are randomly assigned for each mini-batch step in the training process;
- **Adaptive View Selection**: The mask values are adaptively determined by the corresponding queries, described in Appendix C.

Table 12 shows the performance comparison in MPJPE. When using only a single view as input, restricting the input to either the horizontal or vertical view leads to a noticeable increase in MPJPE,

Table 12: Effect of view selection patterns on the performance.

|  | MPJPE | (h) | (v) | (d) |
|---|---|---|---|---|
| Both Views | 20.31 ± 0.34 | 8.44 ± 0.34 | 5.02 ± 0.05 | 15.21 ± 0.67 |
| Horizontal View Only | 20.55 ± 2.05 | 8.74 ± 0.67 | 5.42 ± 0.82 | 14.97 ± 1.99 |
| Vertical View Only | 23.78 ± 0.77 | 12.24 ± 1.45 | 4.66 ± 0.14 | 16.61 ± 1.15 |
| Random Mask | 21.18 ± 0.33 | 8.75 ± 0.78 | 4.96 ± 0.72 | 16.04 ± 0.10 |
| Adaptive View Selection (ours) | 18.99 ± 0.16 | 7.80 ± 0.31 | 4.38 ± 0.25 | 14.54 ± 0.13 |

averaging $20.55\,\mathrm{cm}$ and $23.78\,\mathrm{cm}$, respectively. In contrast, under the multi-view input setting, the adaptive view selection strategy achieves the lowest average MPJPE of $18.99\,\mathrm{cm}$, outperforming both the "random mask" and "both views" configurations.

**Analysis of Attention Weight Assignment:** We investigate the contribution of each radar view by analyzing the distribution of attention weights and view-selection patterns in the pose and joint decoders. Table 13 shows the joint-wise breakdown of attention weight assignment. Note that weights in the pose decoder are normalized over all joints and views, while those in the joint decoder are normalized per joint. Due to averaging over the test set, values may not sum to exactly 1. In the pose decoder, both horizontal and vertical views receive relatively balanced attention across joints. On the other hand, in the joint decoder, vertical views consistently receive higher weights. Specifically, the edge joints exhibit larger disparities, such as 25.07% difference for the knee and 18.39% for the wrist. Since the pose decoder merely aligns the estimates to template poses, it only requires a rough estimate of the overall pose structure, and thus the attention weights are almost evenly distributed for both views. In contrast, the joint decoder is tasked with refining each joint, which makes the vertical view more critical, as it provides more comprehensive visibility across all joints.

Table 13: Attention weight assignment in the pose/joint decoders ($\times$1e-2). The larger value indicates the more attention is weighted on the view for each joint. While the pose decoder assigns balanced attention across joints, the joint decoder makes the vertical view more critical with larger weights.

|  |  | Head | Neck | Shoulder | Elbow | Wrist | Hip | Knee | Ankle |
|---|---|---|---|---|---|---|---|---|---|
| Pose | Horizontal | **4.59** | 3.07 | **3.53** | 3.07 | 2.82 | 3.56 | 2.30 | 1.57 |
|  | Vertical | 4.01 | **3.35** | 2.33 | **3.74** | **3.68** | **4.68** | **2.36** | **3.63** |
| Joint | Horizontal | 34.08 | 40.32 | 41.96 | 40.18 | 31.45 | 38.40 | 32.06 | 35.72 |
|  | Vertical | **50.20** | **46.83** | **46.31** | **46.84** | **49.84** | **49.94** | **57.13** | **47.69** |

In RAPTR, we adopt a view selection approach in which the view selection mask is adaptively determined by the corresponding queries, as described in Appendix C, and Table 14 shows the joint-wise breakdown of view selection assignments in the pose and joint decoders. "Omit Both", "Use Horizontal", "Use Vertical", and "Use Both" are represented by the view selection mask values $[0, 0]$, $[1, 0]$, $[0, 1]$, and $[1, 1]$ for each corresponding row in the mask, respectively. While the pose decoder shows no prominent value imbalance across the assigned patterns, the joint decoder tends to rely more on both views, with high "Use Both" ratios observed for almost all joints, such as 31.06% for shoulder, 32.48% for hip, and 32.44% for neck. This trend reflects that the joint decoder prefers to integrate information from both views when refining keypoint locations.

Table 14: The joint-wise breakdown of view-selection assignments in the pose/joint decoders (%). The joint decoder relies on both views more than the pose decoder, especially for edge joints.

|  |  | Head | Neck | Shoulder | Elbow | Wrist | Hip | Knee | Ankle |
|---|---|---|---|---|---|---|---|---|---|
| Pose | Omit Both | 25.81 | 21.09 | 23.59 | 22.50 | 26.16 | 18.99 | 25.31 | **28.42** |
|  | Use Horizontal | 22.61 | 25.47 | **27.15** | **31.84** | 22.74 | 23.29 | 19.35 | 26.55 |
|  | Use Vertical | **29.65** | 26.70 | 22.72 | 20.79 | 24.43 | 21.76 | 22.42 | 23.33 |
|  | Use Both | 21.92 | **26.74** | 26.54 | 24.86 | **26.67** | **35.96** | **32.92** | 21.70 |
| Joint | Omit Both | 20.73 | 17.87 | 19.94 | 19.24 | 24.20 | 21.54 | 19.13 | 23.48 |
|  | Use Horizontal | 22.75 | 21.55 | 22.94 | 23.18 | 21.24 | 25.86 | 23.17 | 25.25 |
|  | Use Vertical | 27.99 | 28.14 | 26.06 | 27.69 | **28.20** | 20.12 | 26.52 | 24.98 |
|  | Use Both | **28.53** | **32.44** | **31.06** | **29.88** | 26.37 | **32.48** | **31.18** | 26.29 |

### I.2 Additional Visualizations

**Visualization of Pseudo-3D Deformable Attention:** Fig. 11 shows the visualization of pseudo-3D deformable attention at the last layers in the pose and joint decoders. For each view (horizontal on the left, vertical on the right), close-up regions around the bright radar reflections corresponding to subjects are extracted and visualized. The red dots on the plots indicate the sampling locations selected by deformable attention. The pseudo-3D deformable attention mechanism primarily samples features from regions surrounding human subjects. Specifically, in the vertical view, the sampling points are clearly divided and distributed by joint. Moreover, as the sampling offsets are computed across the $x$, $y$, and $z$ axes at once in our pseudo-3D deformable attention, the sampling locations maintain consistent alignment in the depth direction across the views for the same subject.
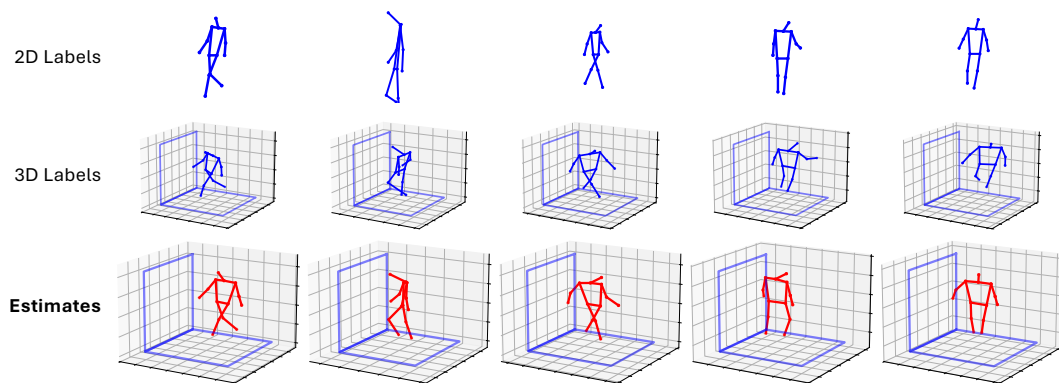


Figure 11: Visualization of pseudo-3D deformable attention. The attention mechanism samples features around the bright signals caused by body reflection, represented as red dots on the plots.

**Additional Visualization Cases:** We provide visualizations of the RAPTR estimation results for more cases in Fig. 12. We present visualizations of 2D and 3D labels and 3D pose estimates for HIBER WALK, MULTI, MMVR P1S1, and P2S1, arranged from top to bottom.
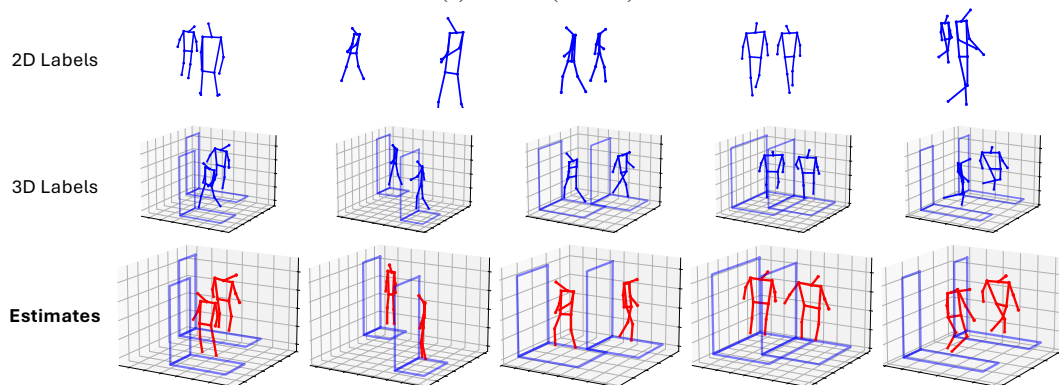
For the HIBER dataset, RAPTR maintains a stable estimation quality by capturing overall body orientation and limb articulation in (a), where only a single subject is presented, and (b) with multiple subjects. Although minor inaccuracies are observed in some cases, RAPTR preserves the spatial arrangement and relative depth of the subjects. As can be observed in the figure, the annotated 3D BBoxes are often significantly larger than the actual human body size, which originates from the dataset itself. Since RAPTR utilizes only the center of the 3D BBoxes as reference for coarse-grained 3D cue, it remains largely unaffected by such inaccuracies in box scale.

For the MMVR dataset, RAPTR has to deal with a diverse range of body configurations, including seated and crouched poses, and increased subject variability. In the P1S1 setting (c), RAPTR consistently estimates plausible 3D pose estimations that are well aligned with 2D keypoint labels, even for non-standard upright posture like spreading arms. On the other hand, in the P2S1 setting (d), RAPTR demonstrates reasonable performance for more complex body configurations with multiple subjects: it often captures the overall structure and spatial arrangement of each individual. In some cases, RAPTR even reconstructs plausible limb configurations where the 2D labels are inaccurate or incomplete, such as the seated person's legs in the leftmost example of (d). This suggests that the architectural design, which first estimates initial poses using the pose decoder with a template pose, then refines joint positions via the joint decoder, effectively preserves a human-body prior throughout the process. However, there are visible imperfections in some cases, such as over-extended limbs or inaccurate limb orientations, illustrating the difficulty of 3D pose estimation in scenarios with occlusion and extreme articulation.
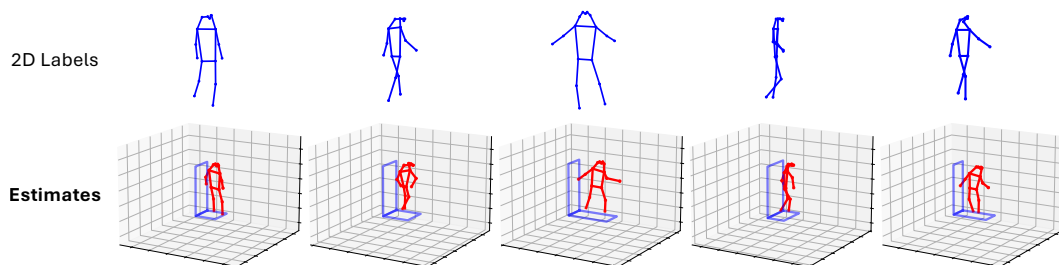
**Failure Cases:** We provide failure cases for the HIBER MULTI dataset in Fig. 13. In (a) and (b), RAPTR locates pose estimations significantly away from the 3D BBox labels, indicating errors in both subject localization and depth reasoning. In (c), the estimated poses appear severely distorted, failing to preserve the anatomical structure of the body. Despite being spatially close to the correct locations, the joint configurations are implausible, indicating a breakdown in fine-grained pose refinement. In (d), RAPTR fails to correctly associate the estimated poses with the 3D BBoxes, leading to redundant
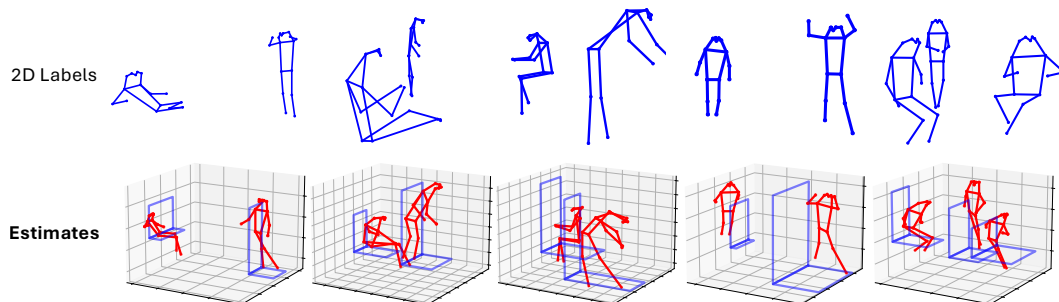
(a) HIBER (WALK)

(b) HIBER (MULTI)

(c) MMVR (P1S1)

(d) MMVR (P2S1)

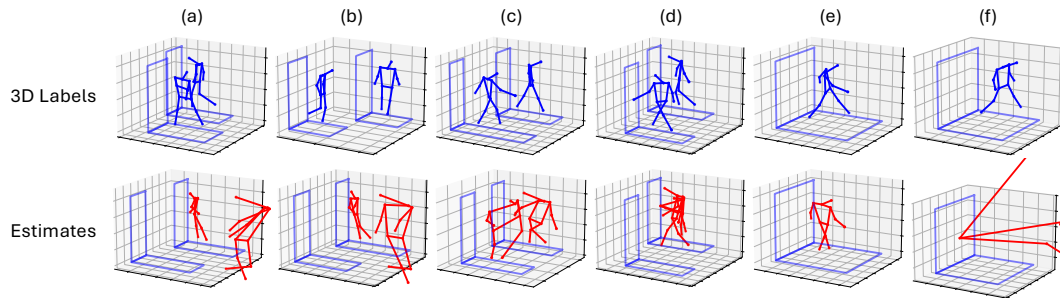Figure 12: Visualization for RAPTR 3D pose estimation.

Figure 13: Failure cases for the HIBER (MULTI) dataset.

estimates that multiple predictions are assigned to the same individual, degrading the quality of the final estimation. In (e), RAPTR fails to recover the correct body orientation. While the label pose is stepping forward with the left leg in the back-right direction, the estimated pose incorrectly keeps the body facing forward and places the right leg in an unnatural cross-step position. In (f), RAPTR produces a completely corrupted estimate with no apparent correspondence to the human pose. These failures are frequently observed when the two individuals overlap on the 2D image plane and the 2D keypoint labels themselves exhibit pose inaccuracies, indicating the limitation of our RAPTR, which relies on the 2D keypoint labels to accurately represent human shapes.