
***mRI*: Multi-modal 3D Human Pose Estimation Dataset using mmWave, RGB-D, and Inertial Sensors**

Sizhe An

University of Wisconsin-Madison
sizhe.an@wisc.edu

Yin Li

University of Wisconsin-Madison
yin.li@wisc.edu

Umit Ogras

University of Wisconsin-Madison
uogras@wisc.edu

Abstract

The ability to estimate 3D human body pose and movement, also known as human pose estimation (HPE), enables many applications for home-based health monitoring, such as remote rehabilitation training. Several possible solutions have emerged using sensors ranging from RGB cameras, depth sensors, millimeter-Wave (mmWave) radars, and wearable inertial sensors. Despite previous efforts on datasets and benchmarks for HPE, few dataset exploits multiple modalities and focuses on home-based health monitoring. To bridge this gap, we present *mRI*¹, a multi-modal 3D human pose estimation dataset with mmWave, RGB-D, and Inertial Sensors. Our dataset consists of over 160k synchronized frames from 20 subjects performing rehabilitation exercises and supports the benchmarks of HPE and action detection. We perform extensive experiments using our dataset and delineate the strength of each modality. We hope that the release of *mRI* can catalyze the research in pose estimation, multi-modal learning, and action understanding, and more importantly facilitate the applications of home-based health monitoring.

1 Introduction

3D Human pose estimation (HPE) refers to detecting and tracking human body parts or key joints (e.g., wrists, shoulders, and knees) in the 3D space. It is a fundamental and crucial task in human activity understanding and movement analysis with numerous application areas, including rehabilitation [45, 34, 7, 6], professional sports [38], augmented/virtual reality, and autonomous driving [29]. In particular, human pose estimation plays an increasingly important role in healthcare applications, such as remote rehabilitation training [40, 20]. The current mainstream rehabilitation treatment involves a physical therapist supervising the patients in person. In contrast, HPE-based health monitoring systems can help clinicians correct patients' movements or instruct them remotely. To this end, multiple datasets have studied HPE with health-related physical movements [6, 45, 34, 7].

Many existing studies rely heavily on processing RGB frames from color cameras for human pose estimation [21, 5, 17, 37, 18, 26]. RGB image and video frames are the most common input types since they offer a non-invasive approach for HPE. However, the image quality depends heavily on the environmental setting, such as light conditions and visibility [3]. Moreover, using image and video data poses significant privacy concerns, especially in a household environment. Finally, the data-intensive nature of real-time video processing requires computationally powerful equipment with high cost and energy consumption.

¹Project page: <http://sizhean.github.io/mri>

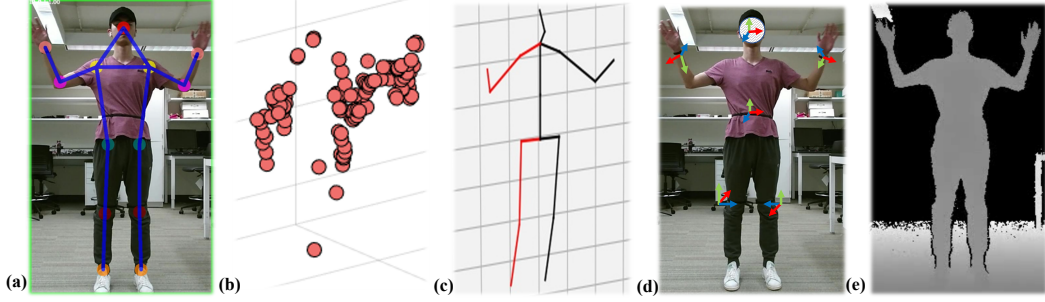


Figure 1: Overview of all modalities and annotations in *mRI* dataset. All sub-figures uses the same sample frame during ‘both upper limb extension’. (a) 2D human keypoints with bounding box on RGB image, (b) 3D mmWave point cloud, (c) 3D human skeletons, (d) IMU rotations, (e) depth image. *mRI* dataset supports **human pose estimation** and **action detection** tasks. With *mRI*, researchers from the fields of machine learning, computer vision, wearable computing can exploit the complementary advantages of **multi-modality**, while clinical and rehabilitation experts can focus on its **healthcare** movements.

Frame quality, privacy, and computational power drawbacks of video processing can be addressed by emerging *complementary sensor modalities*, such as lidar, millimeter wave (mmWave) radar [3, 51, 54], and wearable inertial sensors [48, 46, 47, 49, 31, 52, 2]. The point cloud from lidar overcomes frame quality and privacy challenges. However, it has a high cost and computation power requirements to process the data, making it unsuitable for indoor applications such as rehabilitation. In contrast, mmWave radar can generate high-resolution 3D point clouds of objects while maintaining low cost, privacy, and computational power advantages. Similarly, wearable inertial sensors provide accurate rotation and acceleration information regarding joints with low cost and computational power requirements [46, 47, 49, 2], yet at a price of body worn sensors.

High-quality and large-scale datasets provide a vital foundation for algorithm development. To catalyze research in HPE, this work (*mRI*) combines mmWave radar, RGB-Depth (RGB-D), and Inertial sensors to exploit their complementary advantages. We present a comprehensive 3D human pose estimation dataset performed by 20 human subjects, consisting of more than 160k synchronized frames from three sensing modalities. The contributions and unique aspects of *mRI* are as follows:

- **Multiple Sensing Modalities.** *mRI* consists of mmWave point cloud, RGB frames, depth frames, and inertial signals. The experimental data is captured using a commercial low-power, and a low-cost mmWave radar, two depth cameras, and six high-accuracy inertial measurement units (IMUs). All sensors are temporally synchronized and spatially calibrated. To the best of our knowledge, *mRI* is the first dataset that combines these complementary modalities, as elaborated in Section 2.
- **Healthcare Movements Focus.** We use ten clinically-suggested rehabilitation movements that involve the upper body, lower body, and the major muscles related to human mobility, as described in Section 3.2. These movements are crucial for patients to recover from sequelae of central nervous system disorders, such as Parkinson’s disease (PD) and cerebrovascular diseases (e.g., stroke). Hence, the *mRI* dataset can serve as a reference from healthy subjects, while the experimental methodology can enable future studies with patients.
- **Flexible Data Format and Extensive Benchmarks.** We release the raw synchronized and calibrated sensor data and a comprehensive set of benchmarks for 2D/3D human pose estimation and action detection using multiple modalities (see Section 4). The proposed end-to-end pipeline pre-processes the raw data into the point cloud, features, and 2D/3D keypoints. In addition, all manually-labeled actions annotations and 3D human key points ground truth are released to public, as detailed in Section 3.2.
- **Low-Power & Low-Cost Requirements.** Widespread use of home-based rehabilitation depends critically on the affordability and operating cost of the deployed systems (see Section 3.1). Our *mRI* dataset and findings pave the way to sustainable systems with low-power and low-cost sensors and edge devices. For example, only mmWave radar and IMU sensors can be used in the field after they are trained with all three modalities (including RGB-D) in a clinical environment.

Dataset	Sensing Modalities				# of Subjects	# of Seqs	# of Actions	# of Synced Frames	Annotations		
	RGB	Depth	IMU	mmWave					Action	2DKP	3DKP
COCO [21]	✓	-	-	-	-	-	-	104k	-	✓	-
MPII [5]	✓	-	-	-	-	24k	410	25k	✓	✓	-
MPI-INF-3DHP [26]	✓	-	-	-	8	16	8	1.3M	-	✓	✓
Human3.6M [17]	✓	✓	-	-	11	839	17	3.6M	✓	✓	✓
CMU Panoptic [18]	✓	✓	-	-	8	65	5	154M	-	✓	✓
NTU RGB+D [37]	✓	✓	-	-	40	56k	60	4M	✓	✓	✓
3DPW [46]	✓	-	✓	-	7	60	-	51k	-	✓	✓
MPI08 [31]	✓	-	✓	-	4	24	24	14k	-	-	✓
TNT15 [48]	✓	-	✓	-	1	-	5	14k	✓	-	✓
MoVi [13]	✓	-	✓	-	90	1044	21	712k	✓	✓	✓
RF-Pose [55] [†]	✓	-	-	✓	100	-	1	-	✓	✓	-
RF-Pose3D [54] [†]	✓	-	-	✓	>5	-	5	-	✓	✓	✓
mmPose [36] [†]	-	-	-	✓	2	-	4	40k	✓	-	✓
mmMesh [51] [†]	✓	-	-	✓	20	-	8	3k	✓	-	✓
MARS [3]	-	-	-	✓	4	80	10	40k	✓	-	✓
Reiss et al. [34]	-	-	✓	-	9	-	18	3.6M	✓	-	-
HPTE [7]	✓	✓	-	-	5	240	8	100k	✓	-	✓
EmoPain [8]	✓	-	-	-	50	-	11	33k	✓	-	✓
AHA-3D [6]	✓	-	-	-	21	79	4	170k	✓	-	✓
UI-PRMD [45]	✓	-	-	-	10	100	10	60k	✓	-	✓
<i>mRI</i>	✓	✓	✓	✓	20	300	12	160k	✓	✓	✓

Table 1: Comparison across related datasets. For 2D keypoint annotations, only COCO [21] and MPII [5] are annotated manually, all others are derived from deep models. —: Not report in the paper. †: The dataset is not open-source/available. The first group of rows shows earlier RGB and RGB-D datasets. The middle group of rows presents datasets with emerging sensing modalities such as IMUs and mmWave. The last group of rows lists healthcare-related datasets.

2 Related Work

2.1 3D Human Pose Estimation

Marker-based optical motion capture (MoCap) systems are often used to acquire accurate 3D body pose [17, 8, 45]. Optical MoCap systems require attaching reflective markers to the body and are quite costly, thus are limited to laboratory settings. Recently, MoCap systems based on body-worn IMUs have been developed [46, 31, 49, 48, 13]. They are considerably cheaper yet at a cost of tracking accuracy due to drifting [1]. Our dataset explores using low-cost IMUs for 3D HPE.

Besides marker-based MoCap, Marker-less MoCap has received much attention. Depth cameras are often used for pose estimation [27], yet are limited by their sensing range (within 5 meters). Recent effort has focused on pose estimation using RGB cameras. With the help of machine learning, 3D joints can be estimated from a single RGB image [24], or from several RGB images from different viewing angles captured by multiple cameras [18, 26, 31], or from a sequence of RGB frames within a video [30]. However, RGB cameras are easily affected by poor light conditions, and raise privacy concerns for home-based monitoring. More recently, mmWave-based pose estimation, including radio frequency sensing, has emerged as a promising solution [55, 54, 36, 51, 3]. A mmWave-based solution has demonstrated comparable accuracy to RGB and depth cameras, yet excels at privacy-preserving and long working range. Our dataset includes mmWave for 3D HPE.

Moving forward, the results of 3D HPE can be used by skeleton-based action recognition [23, 37] to localize and recognize actions in time, broadening its applications in health monitoring [22, 28] and human behavior analysis [35]. Our dataset provides action annotations and we evaluate using the estimated pose for temporal action localization [12].

2.2 Datasets for Human Pose Estimation

High-quality datasets with annotations are crucial for the advancement of pose estimation. Table 1 summarizes previous works on HPE datasets and compare them to our *mRI*. Some of the early

effort focuses on 2D HPE (e.g., COCO [21] and MPII [5]), or 3D HPE a single modality (e.g., 3DHP with images, mmPose [36] with mmWave, and MPI08 [31] with IMU). More recent works combines multiple modalities for 3D HPE. For example, Human3.6M [17] contains RGB images and depth maps of 11 professional actors performing 17 daily activities, coupled with ground-truth 3D poses from optical MoCap. RF-Pose3D [54] presents the first study to use radio frequency sensing for 3D HPE, together with a dataset of both RGB images and radio signals. MoVi [13] incorporates both IMU signals and RGB frames, as well as ground-truth 3D poses from MoCap, and presents a benchmark for both 3D HPE and human activity recognition. In comparison to existing dataset,

To the best of our knowledge, *mRI* is the *first HPE dataset with the most comprehensive set of sensing modalities*, including RGB, depth, IMU, and mmWave. In addition, *mRI fills the vacancy of standardized mmWave-based human pose estimation*, as all current mmWave-based HPE datasets are either not open-sourced or without proper keypoints annotations and RGB references.

2.3 Human Pose Estimation for Rehabilitation

HPE promises to capture complex body movement naturally occurring in daily activities or prescribed by clinicians, and thus offers a promising vehicle to inform treatment and to quantify the progress of treatment. Individual sensing modality has been previously considered, including RGB camera [7, 8], depth camera [6, 19, 45], IMUs [34], and MoCap [8, 45]. Reiss et al. [34] presents a dataset monitoring physical activities with three IMUs and a heart rate monitor. The home-based physical therapy exercises (HPTE) dataset [7] uses Kinect to record video and depth streams while users perform eight therapy actions. The EmoPain dataset [8] captures both joint information and face videos to classify the pain level based on the emotion in the rehabilitation movements. The AHA-3D [6] dataset contains 79 skeleton videos recorded by Kinect for four healthcare activities. Similarly, the UI-PRMD [45] dataset captures common physical rehabilitation exercises using the Kinect and Vicon MoCap. Similar to these works, *mRI* focuses on rehabilitation exercises, and provides the most comprehensive set of sensing modalities while remaining competitive in its scale.

3 Dataset

mRI includes 3D point cloud from mmWave, RGB frames and depth maps from RGB-D cameras, joints rotations and accelerations from wearable IMU sensors, as well as annotations of 2D keypoints, 3D joints, and action labels of 12 clinically relevant movements. *mRI* consists of 300 time-series sequences with 160K synchronized frames and more than 5M total data points from all sensors, from 20 subjects. We hope that our dataset will contribute to the multi-modal machine learning community, and facilitate applications of HPE for rehabilitation and other healthcare problems. In what follows we describe the hardware system to capture the data and the data collection process. More details can be found in the supplementary A.3.

3.1 Capturing Multi-Modal Signals for Human Pose Estimation

To capture multi-modal data, we designed a sensor system composed of one mmWave radar, two sets of RGB and depth cameras, and six wearable IMUs. Detailed specifications and features of all sensors are shown in Table 2. The mmWave radar and two Kinect V2 sensors are placed on a desk 2.4 m away from the subject, wearing six IMU sensors, as shown in Figure 2. We now describe the data capturing for each modality and the synchronization across modalities.

	#	Freq.	Con.	Power	Privacy \uparrow	Anti-inter. \uparrow	Intrusive	Output
mmWave [41]	1	10 Hz	Wired	2.1 W	***	***	No	Point cloud
RGB [27]	2	30 Hz	Wired	16 W	*	*	No	RGB frame
Depth [27]	2	30 Hz	Wired	16 W	**	**	No	Depth and infra-red frame
IMU [50]	6	50 Hz	BLE	120 mW	***	***	Yes	Accelerations and quaternions

Table 2: **Comparison across sensors.** #: Number of sensors. Freq.: Sampling frequency. Con.: Type of connection to the host PC. Privacy indicates privacy-preserving ability. Anti-interference represents how much it is affected by environmental factors like non-ideal lighting.

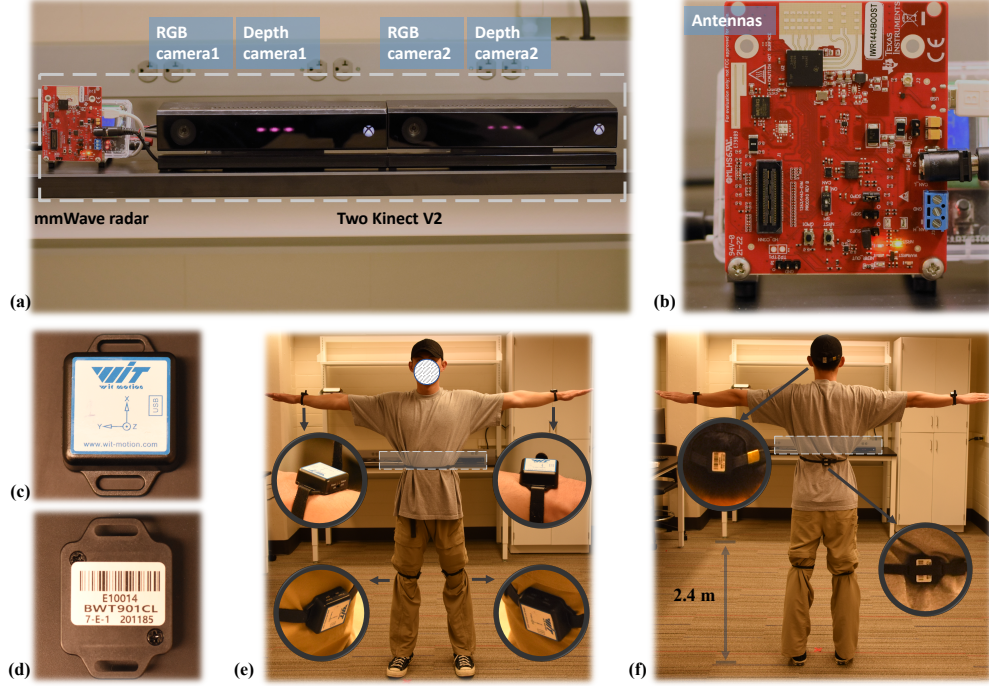


Figure 2: Overview of the experimental setup. (a) shows all non-intrusive sensors, including mmWave radar, two RGB, and depth cameras. (b) shows a zoom-in version of the mmWave radar and its antennas. The front and back views of the IMU are shown in (c) and (d), respectively. (e) and (f) show the front and back view of the subject standing as a “T pose” with six IMUs and zoom-in views of IMUs. The gray dash line boxes in (a), (e), and (f) represent the position of non-intrusive sensors.

Point cloud from mmWave radar. A Texas Instruments (TI) IWR1443 Boost mmWave radar [41] is used to obtain the mmWave point cloud. 3D mmWave point cloud is generated by Frequency Modulated Continuous Wave (FMCW) radar using multiple transmit (Tx) and receiver antennas (Rx) configuration [3, 36, 54]. The radar emits a chirp signal, a sinusoid wave whose frequency increases linearly with time. Then the reflected signals are received at the Rx antenna side. The range, velocity and angle resolutions are computed with the received data using *range FFT*, *Doppler FFT*, and *angle estimation* algorithms, respectively. After the constant false alarm rate (CFAR) algorithm eliminates the noise, a point cloud capturing object shape and movement is constructed as

$$P_i = (x_i, y_i, z_i, d_i, I_i), i \in \mathbb{Z}, 1 \leq i \leq N \quad (1)$$

where x_i, y_i, z_i are the spatial coordinates of the point, d_i represents the Doppler velocity, I_i denotes the signal intensity, and $N = 64$ represents the total number of points in a given frame. To further increase the density of the point cloud, we follow [4] to fuse points from three consecutive frames, i.e., increasing the number of points per frame from 64 to 192. See more details in the Supplement A.4.

The radar is connected to the host PC through the UART interface. We modify a Matlab Runtime implementation from TI [43] for the data acquisition. The sampling rate is set to 10 Hz since it is sufficient for measuring human movement (the frequency of most voluntary human movements spans from 0.6 to 8 Hz [14]).

RGB and depth frames from RGB-D cameras. Two Microsoft Kinect V2 [27] sensors are used to capture RGB and depth frames. Kinect V2 has a high precision color camera and infra-red camera, generating color and depth frame with a resolution of 1920×1080 and 512×424 , respectively. We modified the software from libfreenect2² to generate aligned color, depth, and infra-red frames with the global timestamp from two Kinect V2 sensors. We calibrate the two cameras using the Matlab camera calibration toolbox [25]. The center of the RGB camera 1, as shown in Figure 2 (a) is selected as the origin of the world coordinate system.

²<https://github.com/OpenKinect/libfreenect2>

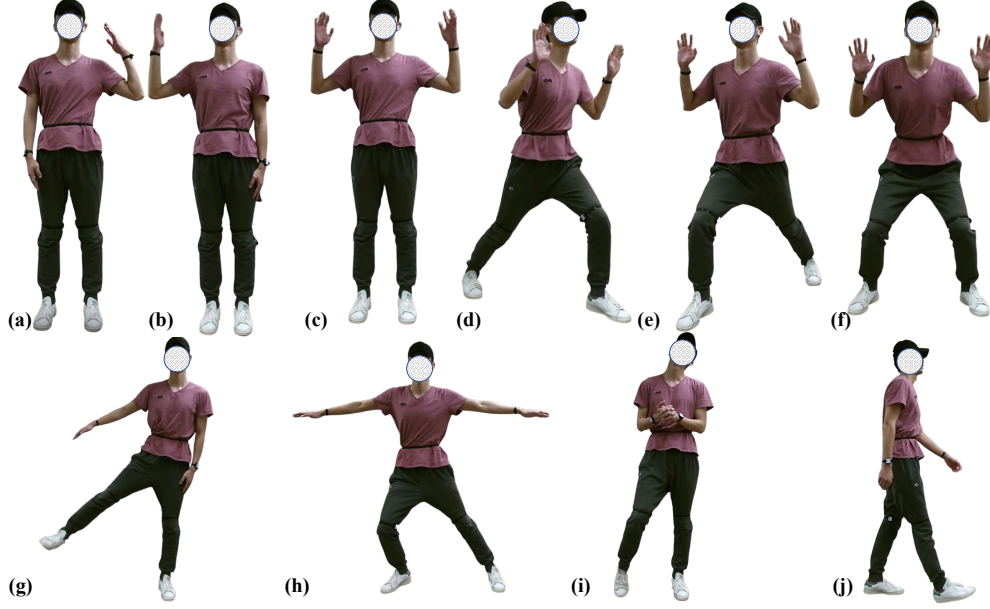


Figure 3: Overview of all movements in *mRI*, as described in Section 3.2. The mirror movements of (g) and (h) are not shown due to limited space.

Joints rotations and accelerations from wearable IMUs. Six Wit-motion BWT901CL IMUs [50] are used to capture the rotation and acceleration of the human joints. In our experiments, the IMUs are tightly attached to the left wrist, right wrist, left knee, right knee, head, and pelvis of the subject to capture the complete information about the human body, as shown in Figure 2(e) and (f). Each IMU contains a 3-axis accelerometer, 3-axis gyrometer, and 3-axis magnetometer as the sensing unit. The raw output data from the sensors are accelerations, angular velocity, Euler angle, and magnet field values. Based on these values, we extract joint quaternion and 3-axis acceleration following [16, 52] as they fully specify the body pose and movement. The IMUs connect to the host PC via a USB-HID device using the BLE protocol with a sampling rate of 50 Hz (see Table 2), ensuring low-latency data transmission with multiple devices.

Sensors synchronization. All sensors are connected to the same host PC, allowing global timestamps from the host attached to each data point from different sensors. We then synchronize all data points using these global timestamps. Since mmWave radar has the lowest sampling rate, we use its timestamp as the basic timestamp. For each timestamp in mmWave radar, we find the timestamp in other sensors with the minimum absolute difference between itself and the mmWave timestamp and align them. The time difference between sensors is less than 5 ms with the proposed time alignment method. Finally, the synchronized data across all modalities have the same number of data points.

3.2 Data Collection, Annotation, and Visualization

Rehabilitation exercises. We consider 12 movements related to rehabilitation exercises covering the entire human body. The first ten rehabilitation movements are modified from [45, 3]. Figure 3 shows all movements: (a) left upper limb extension, (b) right upper limb extension, (c) both upper limb extension, (d) left front lunge, (e) right front lunge, (f) squat, (g) left side lunge, right side lunge, (h) left limb extension, and right limb extension. The 11th and 12th movement are stretching and relaxing in free forms (i), and walking in a straight line (j), respectively. These two movements are meant to increase the diversity of the dataset, as the 11th movement is determined by each subject and the 12th movement features a global displacement. The duration of each type of movement is around one minute per subject. To calibrate the IMUs, we require the subject to perform a “T Pose” at the beginning of each recording.

Participant recruitment and consent. To conduct human subject study, we obtained an approval from the IRB at the university. Our participants were recruited locally and all experiments were

carried out in a laboratory setting. Before each session, a researcher introduces the research goal, experiment procedure, and potential risk via both verb communication and video tutorials. The participant is free to raise questions before he or she sign the consent form, and is free to withdraw from the study at any time. We refer more details to our Ethic statements.

20 healthy participants consented and managed to perform the study. There are 13 males and 7 females, with an average age of 24.1 ± 4.4 and a height of 175.6 ± 9.3 cm.

Obtaining human body pose. We now describe how we derive 2D keypoints and 3D joints given our sensor data. Without using MoCap, our solution is a combination of 2D keypoint detection (body parts), 3D triangulation (joints), and an optimization-based refinement.

- First, we use HRNet [39] (with bounding boxes from Mask RCNN [15]) to detect 2D keypoints of human body parts in all RGB frames from both cameras.
- Next, we triangulate two sets of 2D keypoints captured at the same time yet from different cameras, using camera parameters obtained via camera calibration. The results are a set of 3D body joints (17 in total following COCO format).
- Finally, we refine the 3D joints in each video by solving an optimization problem. Our optimization minimizes 2D reprojection error, imposes equal bone length constraint for all frames, and enforces temporal smoothness of the 3D joints.

Specifically, our refinement step solves the following optimization problem

$$\min_{\{\mathbf{p}_i\}} \sum_{i=1}^{\mathbb{Z}} \left(\|P^l \mathbf{p}_i - \mathbf{q}_i^l\| + \|P^r \mathbf{p}_i - \mathbf{q}_i^r\| \right) + \sum_j^{\text{bonelist}} \|\mathbf{B}_j - \text{median}(\mathbf{B})\| + \sum_{i=1}^{\mathbb{Z}-1} \|\mathbf{p}_{i+1} - \mathbf{p}_i\|, \quad (2)$$

where $\{\mathbf{p}_i\}$ is the set of 3D joints of size \mathbb{Z} , \mathbf{q}_i^l and \mathbf{q}_i^r are the 2D keypoints from the left and right camera, respectively. P^l and P^r are the camera projection matrix for the left and right camera, respectively. $\{\mathbf{B}_j\}$ is a set of bone length defined by connecting a subset of the joints (e.g., wrist to elbow, elbow to shoulder). The first term represents the re-projection errors of the two cameras. The second term enforces equal bone length across all frames in the same video (i.e., the same subject). And the third term imposes temporal smoothness of the 3D joint coordinates. More details, including both quantitative and qualitative results, can be found in the supplement. After the optimization, we re-project the 3D joints to 2D and thus update the 2D keypoints.

Keypoints quality. To validate the reliability of the obtained 3D joints, we report the reprojection error of the derived 3D joints by comparing their 2D projections to human annotated 2D keypoints. Specifically, we randomly sample 50 video frames from our dataset, manually annotate the 2D keypoints for each frame, and calculate the error between the projected 3D joints and the annotated 2D keypoints, following [5]. The mean absolute percentage error (MAPE) is 1.5%, and the percentage of correct keypoints thresholded at 50% of the head segment length (PCKh) is 98.9. More details and visualization can be found in the supplement A.2.

Annotating actions in videos. We also provide annotations of the 12 movements for each video. The multi-media annotation tool ELAN [9] is employed to annotate the videos. For each video sequence, we manually label the start and end timestamp and the category of the 12 different movements.

4 Evaluation and Benchmarks

We introduce a standardized evaluation pipeline of using our dataset for 3D human pose estimation and human action detection. We use latest models to benchmark the performance of each modality and discuss their results.

4.1 3D Human Pose Estimation

Our main benchmark is 3D HPE. We now describe our experiment protocol, evaluation metrics, and the method we used, followed by the presentation of our results.

Experiment protocol. We consider two settings of data splits. **Setting 1 (S1 Random Split):** A random split of 80% and 20% of all data is used as the training and testing set, respectively. **Setting 2 (S2 Split by Subjects):** A randomly selected subset (80%, i.e., 16 out of 20) of the subjects is

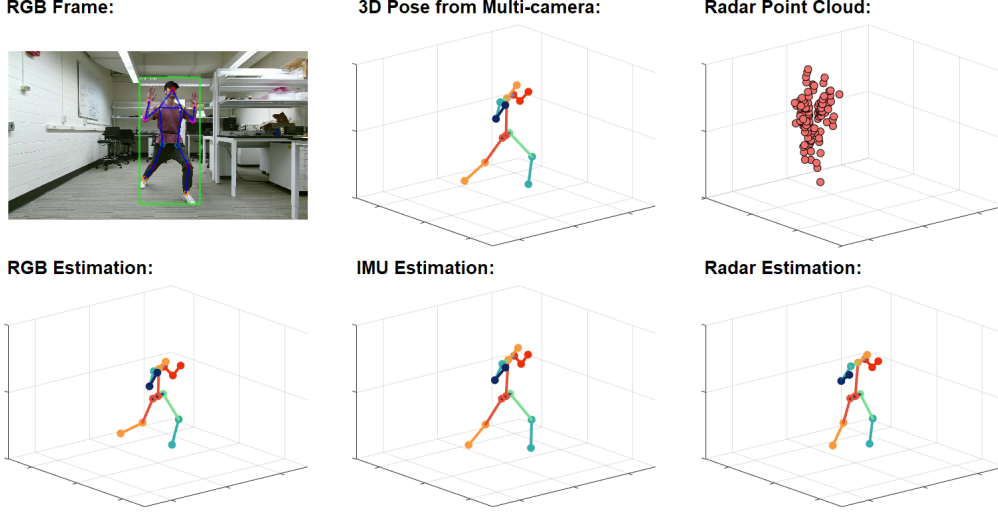


Figure 4: Visualization of sample pose data and results during left front lunge. Top row (from left to right): an RGB frame with detected human bounding box and 2D keypoints, the refined 3D pose derived from two cameras, and the 3D point cloud from mmWave radar. Bottom row (from left to right) shows the estimated 3D pose from a single RGB camera, IMU signals, and mmWave radar.

used for training, while the rest are for testing. S1 mimics a case where personalized HPE model is possible, while S2 evaluates across-subject generalization. For each setting, we randomly sample three splits and report the averaged results. More details are provided in the supplement A.3.

Further, we also define two evaluation protocols based on the design of our movements, as mentioned in Section 3.2. **Protocol 1 (P1)** consists of all 12 movements, including stretching and relaxing in free forms and walking. While **Protocol 2 (P2)** only considers the first ten rehabilitation movements. Such protocols help us investigate the robustness of the model in terms of fixed/free form movement.

Evaluation metrics. We adopt Mean Per Joint Position Error (MPJPE) and Procrustes Analysis MPJPE (PA-MPJPE), widely used in human body pose estimation [17], as the main metrics. MPJPE represents the mean Euclidean distance between ground truth and prediction for all joints. MPJPE is calculated after aligning the root joints (the pelvis) of the estimated and ground truth 3D pose. PA-MPJPE is MPJPE after being aligned to the ground truth by the Procrustes method [10], a similarity transformation including rotation, translation, and scaling. We also report additional metrics such as joint angles provided in the supplement.

Methods. We conduct 3D human pose estimation using mmWave, RGB, and IMUs separately using latest methods. Here we briefly introduce the methods considered in our evaluation and refer to our supplement for more implementation details.

- **mmWave:** We use the data processing pipeline and model from [3] that learns a convolutional neural network on the 5D point cloud to regress the 3D joints. The model is trained from scratch on our dataset, and outputs the 3D joints in the global coordinates system.
- **RGB:** We adopt the model from [30], where 2D keypoints from a sequence of frames are “lifted” into 3D joints (in the camera coordinate system) using a convolutional neural network. We use the pre-trained model from [30]. As the pre-trained model outputs a different set of joint, we only evaluate on a subset that intersects with our set of joints.
- **IMUs:** We employ the feature processing method from [52], with a convolutional neural network trained to regress rotations relative to a root joint (e.g., pelvis) using data from IMUs. The model is trained from scratch on our dataset.

Results and discussion. Table 3 shows the 3D HPE results for mmWave, RGB, and IMUs. Under S1 and P1, mmWave-based HPE achieves 163 and 94 mm for MPJPE and PA-MPJPE, respectively. The metrics are further reduced to 125 and 74 mm for P2. IMU-based HPE obtains MPJPE and

		Protocol 1		Protocol 2	
Modality	Setting	MPJPE (mm)↓	PA-MPJPE (mm)↓	MPJPE (mm)↓	PA-MPJPE (mm)↓
mmWave	S1	163.3±9.1	94.1±3.6	125.1±2.4	74.1±1.0
	S2	186.6±23.8	97.3±7.8	126.6±11.3	75.0±7.1
RGB	S1	116.9±0.1	66.8±0.2	115.0±0.1	64.4±0.1
	S2	120.1±3.7	67.5±1.9	118.4±3.8	64.7±1.4
IMUs	S1	80.2 ±12.6	51.9 ±1.9	40.9 ±1.0	28.4 ±0.9
	S2	147.4±18.4	74.5±5.9	94.3±13.8	54.0±4.9

Table 3: 3D human pose estimation results for mmWave, RGB, and IMUs. We report the mean and standard deviation of joint errors averaged across multiple splits under both our settings (**S1** & **S2**).

PA-MPJPE of 87 and 60 mm, respectively, under **S1** and **P1**. Figure 4 shows visualization comparison of estimation across different modalities.

Under **S2**, mmWave-based HPE performs similarly to **S1**, while IMU-based HPE obtains worse results than **S1**. This is because the sensing data from IMU is more fine-grained on the joints while mmWave grasps more information about body trunk, which is not too subject-specific. As a result, the IMU-based model is more sensitive to different subjects. We can observe that for all modalities **S2** yields higher standard deviations than **S1** since the difference between subjects is much more significant than random split, between train and test set. Similarly, **P1** yields higher standard deviations than **P2** since all movements in **P2** are fixed positions, which makes the model learning the keypoints distribution easier.

RGB-based HPE achieve 116 and 66 mm MPJPE and PA-MPJPE for **P1** under **S1**. Both data-split yield similar results. To compare, the same model achieves 36 mm PA-MPJPE on Human3.6M dataset. However, the model is trained and evaluated on Human3.6M while it is only evaluated on *mRI* without any fine-tuning. We leave fine-tuning the model on *mRI* as future work. In summary, all modalities perform reasonably well on our dataset.

Result visualization. We further visualize sample results of 3D pose estimation from different modalities in Figure 4. Additional examples can be found in the supplement A.5.

4.2 Skeleton-based Action Detection

Moving forward, we explore using the estimated 3D joints for temporal action detection in untrimmed videos — the simultaneous localization and recognition of action instance in time. Specifically, given an input untrimmed video, temporal action localization seeks to predict a set of action instances with varying size. Each instance is defined by its onset, offset, and action labels.

Experiment protocol. We consider the more challenging setting **S2**, where a model is tasked to detect actions performed by subjects not presented in the training set. Here each movement type defines one action category. Similar to our HPE experiments, we evaluate on both **P1** (12 categories) and **P2** (10 categories focusing on rehabilitation exercises). Importantly, we consider using individual modalities and all combinations of these modalities (e.g., RGB+IMU or RGB+mmWave). To combine multiple modalities, 3D joint data from each modality at every time step is concatenated, and the resulting sequence is fed into the model.

Evaluation metrics. Following prior work [12], we report the mean average precision (*mAP*) at multiple temporal intersection over union (tIoU) thresholds ([0.5:0.05:0.95]). tIoU is defined as the intersection over union between two temporal windows, i.e., the 1D Jaccard index. Given a tIoU threshold (e.g., 0.75), *mAP* computes the mean of average prevision across all action categories. An average *mAP* is also reported by averaging the *mAP* scores across all tIoUs.

Method. We make use of a latest method — ActionFormer [53] for temporal action detection. ActionFormer develops a Transformer based model and achieves state-of-the-art results across action detection benchmarks. Specifically, we feed the model with a sequence of estimated 3D poses from different modalities at a sampling rate of 2 Hz, and train the model from scratch on our dataset.

Results and discussion. Table 4 summarizes the results from three modalities and their combinations averaged across all splits. Overall, all modalities perform fairly well, with *mAP* scores around 90%. Under **P1**, IMUs data have the best results with 93.4% *mAP*, and outperform the RGB frames and

Modality	$mAP\uparrow$							
	Protocol 1				Protocol 2			
	tIoU=0.50	tIoU=0.75	tIoU=0.95	average	tIoU=0.50	tIoU=0.75	tIoU=0.95	average
mmWave (W)	98.22 \pm 3.08	97.59 \pm 4.17	29.02 \pm 6.31	87.04 \pm 4.89	99.00 \pm 1.73	97.21 \pm 2.41	31.11 \pm 15.34	87.55 \pm 3.61
RGB (R)	100.00 \pm 0.00	99.14 \pm 0.75	44.80 \pm 10.55	91.56 \pm 2.08	100.00 \pm 0.00	100.00 \pm 0.00	59.87 \pm 8.12	95.07 \pm 1.46
IMUs (I)	100.00 \pm 0.00	100.00 \pm 0.00	53.55 \pm 12.39	93.46 \pm 2.30	100.00 \pm 0.00	100.00 \pm 0.00	60.13 \pm 6.82	94.89 \pm 1.39
W+R	100.00 \pm 0.00	100.00 \pm 0.00	55.71 \pm 11.20	94.17 \pm 1.58	100.00 \pm 0.00	100.00 \pm 0.00	59.89 \pm 15.18	95.09 \pm 2.14
W+I	100.00 \pm 0.00	100.00 \pm 0.00	56.53 \pm 12.23	94.38 \pm 1.70	100.00 \pm 0.00	100.00 \pm 0.00	62.42 \pm 5.65	95.26 \pm 1.08
I+R	100.00 \pm 0.00	99.61 \pm 0.67	60.10 \pm 11.97	94.54 \pm 1.45	100.00 \pm 0.00	100.00 \pm 0.00	61.10 \pm 8.46	94.80 \pm 1.28
W+R+I	100.00 \pm 0.00	100.00 \pm 0.00	60.62 \pm 8.42	94.88 \pm 1.75	100.00 \pm 0.00	100.00 \pm 0.00	66.16 \pm 10.89	95.83 \pm 1.50

Table 4: Action detection results with mmWave (**W**), RGB (**R**), IMUs (**I**), and their combinations. We report the mean and standard deviation of mAP averaged across 3 splits under our setting 2 (**S2**).

radar signals by 1.9% and 6.4%, respectively. Under **P2**, both IMUs data and RGB frames perform equally well with improved mAP (around 95%). The RGB frames achieve a major improvement when evaluated under **P2**. It is interesting to cross reference the results of HPE and action detection. While RGB frames have lower joint errors under **S2** and **P1**, they have slightly worse results on action detection. On the other hand, IMUs data perform consistently well on action detection in **P1** and **P2**.

Further combining the modalities results in a noticeable performance boost. It is probably not surprising that using all three modalities yields the best results, outperforming the best single modality by 1.4% (**P1**) and 0.8% (**P2**) in average mAP and with most gains in mAP under tIoU=0.95 (+7.1% for **P1** and +6.0% for **P2**). Fusing any of the two modalities leads to improved performance than the best of the constituting modality, except the combination of IMUs+RGB under **P2**. These results demonstrate a first step towards multi-modal learning with our dataset.

mmWave radar is less invasive than IMU sensors and offers better privacy than RGB cameras. While in its infancy for human sensing, this modality presents an emerging solution for home-based health monitoring. Part of our goal in this paper is to explore mmWave radar for human sensing by comparing its performance to other common modalities. The results indicate that mmWave radar leads to compelling performance for both human pose estimation and action localization. While its results are worse than those with RGB cameras or IMU sensors, mmWave radar might still be preferred for privacy-sensitive applications.

5 Ethics Statement

The human subject studies reported in this paper was reviewed and approved by the IRB committee at the University of Wisconsin-Madison. Each participant was informed about the research project and signed the consent form. The data has been de-identified with facial information blurred in all videos and participant ID anonymized, and made publicly available to facilitate future research.

To the best of the authors’ knowledge, this work does not disadvantage any person directly. The authors do acknowledge that any pose estimation and activity recognition method can potentially be used with malicious intent, such as tracking user movements. If the human pose estimation/human activity understanding algorithms are directly used to make decisions for patients, potential failures in the classification would affect the users’ quality of life. Therefore, the data and insights on patient activity must be verified by health professionals before making any decisions.

6 Conclusion and future work

In this paper, we proposed *mRI*— a multi-modal 3D human pose estimation dataset of rehabilitation exercises performed by 20 subjects, consisting of more than 160k synchronized frames. *mRI* combines mmWave, RGB-D, and IMUs as sensing modalities, and thus provides the most comprehensive benchmark to date for pose estimation and action detection. We described the creation of our dataset and demonstrated extensive benchmarks using our dataset. Our results help to understand the advantages of individual sensing modalities in the context of home-based health monitoring. We hope that *mRI* can catalyze the research including but not limited to pose estimation, multi-modal learning, and action understanding, thus facilitating critical applications in healthcare. We envision a variety of meaningful future work leveraging our dataset, drawing attention from communities including machine learning, computer vision, wearable computing, multi-modal sensing, and healthcare.

Acknowledgments and Disclosure of Funding

This research was funded by NSF CAREER award CNS-2114499.

References

- [1] N. Ahmad, R. A. R. Ghazilla, N. M. Khairi, and V. Kasi. Reviews on various inertial measurement unit (imu) sensor applications. *International Journal of Signal Processing Systems*, 1(2):256–262, 2013.
- [2] S. An, G. Bhat, S. Gumussoy, and U. Ogras. Transfer learning for human activity recognition using representational analysis of neural networks. *arXiv preprint arXiv:2012.04479*, 2020.
- [3] S. An and U. Y. Ogras. Mars: mmwave-based assistive rehabilitation system for smart healthcare. *ACM Transactions on Embedded Computing Systems (TECS)*, 20(5s):1–22, 2021.
- [4] S. An and U. Y. Ogras. Fast and scalable human pose estimation using mmwave point cloud. In *Proceedings of the 59th ACM/IEEE Design Automation Conference*, page 889–894, 2022.
- [5] M. Andriluka, L. Pishchulin, P. Gehler, and B. Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *Proceedings of the IEEE Conference on computer Vision and Pattern Recognition*, pages 3686–3693, 2014.
- [6] J. Antunes, A. Bernardino, A. Smailagic, and D. P. Siewiorek. Aha-3d: A labelled dataset for senior fitness exercise recognition and segmentation from 3d skeletal data. In *Prof. of The British Machine Vision Conference (BMVC)*, page 332, 2018.
- [7] I. Ar and Y. S. Akgul. A computerized recognition system for the home-based physiotherapy exercises using an rgbd camera. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 22(6):1160–1171, 2014.
- [8] M. S. Aung et al. The automatic detection of chronic pain-related expression: requirements, challenges and the multimodal emopain dataset. *IEEE Transactions on Affective Computing*, 7(4):435–451, 2015.
- [9] H. Brugman, A. Russel, and X. Nijmegen. Annotating multi-media/multi-modal resources with elan. In *LREC*, pages 2065–2068, 2004.
- [10] K. Daniilidis. Pose from 3d point correspondences: The procrustes problem - pose estimation, 2022.
- [11] V. Dham. Programming chirp parameters in ti radar devices. *Application Report SWRA553, Texas Instruments*, 2017.
- [12] B. G. Fabian Caba Heilbron, Victor Escorcia and J. C. Niebles. ActivityNet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 961–970, 2015.
- [13] S. Ghorbani, K. Mahdavian, A. Thaler, K. Kording, D. J. Cook, G. Blohm, and N. F. Troje. Movi: A large multi-purpose human motion and video dataset. *Plos one*, 16(6):e0253157, 2021.
- [14] A. Godfrey, R. Conway, D. Meagher, and G. ÓLaighin. Direct measurement of human movement by accelerometry. *Medical engineering & physics*, 30(10):1364–1386, 2008.
- [15] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask R-CNN. In *Proc. of IEEE Intl. Conf. on Computer Vision*, pages 2961–2969, 2017.
- [16] Y. Huang, M. Kaufmann, E. Aksan, M. J. Black, O. Hilliges, and G. Pons-Moll. Deep inertial poser: Learning to reconstruct human pose from sparse inertial measurements in real time. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 37:185:1–185:15, Nov. 2018. First two authors contributed equally.
- [17] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu. Human3. 6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE transactions on pattern analysis and machine intelligence*, 36(7):1325–1339, 2013.
- [18] H. Joo, H. Liu, L. Tan, L. Gui, B. Nabbe, I. Matthews, T. Kanade, S. Nobuhara, and Y. Sheikh. Panoptic studio: A massively multiview system for social motion capture. In *The IEEE International Conference on Computer Vision (ICCV)*, 2015.
- [19] D. Leightley, J. Darby, B. Li, J. S. McPhee, and M. H. Yap. Human activity recognition for physical rehabilitation. In *2013 IEEE International Conference on Systems, Man, and Cybernetics*, pages 261–266. IEEE, 2013.
- [20] Y. Li, C. Wang, Y. Cao, B. Liu, J. Tan, and Y. Luo. Human pose estimation based in-home lower body rehabilitation system. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2020.
- [21] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.

- [22] Y. Ling and H. Wang. Unsupervised human activity segmentation applying smartphone sensor for healthcare. In *2015 IEEE 12th Intl Conf on Ubiquitous Intelligence and Computing and 2015 IEEE 12th Intl Conf on Autonomic and Trusted Computing and 2015 IEEE 15th Intl Conf on Scalable Computing and Communications and Its Associated Workshops (UIC-ATC-ScalCom)*, pages 1730–1734. IEEE, 2015.
- [23] J. Liu, A. Shahroudy, M. Perez, G. Wang, L.-Y. Duan, and A. C. Kot. Ntu rgb+d 120: A large-scale benchmark for 3d human activity understanding. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(10):2684–2701, 2020.
- [24] J. Martinez, R. Hossain, J. Romero, and J. J. Little. A simple yet effective baseline for 3d human pose estimation. In *Proceedings of the IEEE international conference on computer vision*, pages 2640–2649, 2017.
- [25] Mathworks. Using the Stereo Camera Calibrator App. <https://www.mathworks.com/help/vision/ug/using-the-stereo-camera-calibrator-app.html>, 2018.
- [26] D. Mehta, H. Rhodin, D. Casas, P. Fua, O. Sotnychenko, W. Xu, and C. Theobalt. Monocular 3d human pose estimation in the wild using improved cnn supervision. In *3D Vision (3DV), 2017 Fifth International Conference on*. IEEE, 2017.
- [27] Microsoft. Kinect sensor. <https://developer.microsoft.com/en-us/windows/kinect/> accessed 29 Sep. 2020, 2014.
- [28] H. F. Nweke, Y. W. Teh, G. Mujtaba, and M. A. Al-Garadi. Data fusion and multiple classifier systems for human activity detection and health monitoring: Review and open research directions. *Information Fusion*, 46:147–170, 2019.
- [29] E. Odemakinde. Human pose estimation with deep learning - ultimate overview in 2021, Sep 2021.
- [30] D. Pavlo, C. Feichtenhofer, D. Grangier, and M. Auli. 3d human pose estimation in video with temporal convolutions and semi-supervised training. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [31] G. Pons-Moll, A. Baak, T. Helten, M. Müller, H.-P. Seidel, and B. Rosenhahn. Multisensor-fusion for 3d full-body human motion capture. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2010.
- [32] Pytorch. Pytorch Mobile. <https://pytorch.org/mobile/home/> accessed 8 Jul. 2021, 2022.
- [33] S. Rao. Introduction to mmwave sensing: Fmcw radars. *Texas Instruments (TI) mmWave Training Series*, 2017.
- [34] A. Reiss and D. Stricker. Creating and benchmarking a new dataset for physical activity monitoring. In *Proceedings of the 5th International Conference on Pervasive Technologies Related to Assistive Environments*, pages 1–8, 2012.
- [35] M. Rohrbach, S. Amin, M. Andriluka, and B. Schiele. A database for fine grained activity detection of cooking activities. In *2012 IEEE conference on computer vision and pattern recognition*, pages 1194–1201. IEEE, 2012.
- [36] A. Sengupta, F. Jin, R. Zhang, and S. Cao. Mm-pose: Real-time human skeletal posture estimation using mmwave radars and cnns. *IEEE Sensors Journal*, 20(17):10032–10044, 2020.
- [37] A. Shahroudy, J. Liu, T.-T. Ng, and G. Wang. Ntu rgb+d: A large scale dataset for 3d human activity analysis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1010–1019, 2016.
- [38] C. SIMON-AL-ARAJI. Bringing ai to the nba, 2019.
- [39] K. Sun, B. Xiao, D. Liu, and J. Wang. Deep high-resolution representation learning for human pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5693–5703, 2019.
- [40] T. Tao, X. Yang, J. Xu, W. Wang, S. Zhang, M. Li, and G. Xu. Trajectory planning of upper limb rehabilitation robot based on human pose estimation. In *2020 17th International Conference on Ubiquitous Robots (UR)*, pages 333–338. IEEE, 2020.
- [41] Texas Instruments. IWR1443BOOST. <https://www.ti.com/tool/IWR1443BOOST> accessed 29 Sep. 2020, 2014.
- [42] Texas Instruments. mmWave tutorial. <https://www.ti.com/lit/pdf/swra553> accessed 29 Sep. 2020, 2014.
- [43] Texas Instruments. Zone Occupancy. <https://www.ti.com/lit/pdf/tiduea7> accessed 8 Apr. 2021, 2018.
- [44] Texas Instruments. mmWave fundamentals. <https://www.ti.com/lit/spyy005> accessed 8 Apr. 2021, 2020.

- [45] A. Vakanski, H.-p. Jun, D. Paul, and R. Baker. A data set of human body movements for physical rehabilitation exercises. *Data*, 3(1):2, 2018.
- [46] T. von Marcard, R. Henschel, M. Black, B. Rosenhahn, and G. Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *European Conference on Computer Vision (ECCV)*, sep 2018.
- [47] T. von Marcard, R. Henschel, M. J. Black, B. Rosenhahn, and G. Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 601–617, 2018.
- [48] T. von Marcard, G. Pons-Moll, and B. Rosenhahn. Human pose estimation from video and imus. *Transactions on Pattern Analysis and Machine Intelligence*, 38(8):1533–1547, Jan. 2016.
- [49] T. Von Marcard, B. Rosenhahn, M. J. Black, and G. Pons-Moll. Sparse inertial poser: Automatic 3d human pose estimation from sparse IMUs. In *Computer Graphics Forum*, volume 36-2, pages 349–360. Wiley Online Library, 2017.
- [50] Wit-motions. BWT901CL. <https://www.wit-motion.com/9-axis/witmotion-bluetooth-2-0-mult.html> accessed 8 Apr. 2021, 2021.
- [51] H. Xue, Y. Ju, C. Miao, Y. Wang, S. Wang, A. Zhang, and L. Su. mmmesh: Towards 3d real-time dynamic human mesh construction using millimeter-wave. In *Proceedings of the 19th Annual International Conference on Mobile Systems, Applications, and Services*, pages 269–282, 2021.
- [52] X. Yi, Y. Zhou, and F. Xu. Transpose: Real-time 3d human translation and pose estimation with six inertial sensors. *ACM Transactions on Graphics*, 40(4), 08 2021.
- [53] C. Zhang, J. Wu, and Y. Li. Actionformer: Localizing moments of actions with transformers. In *European Conference on Computer Vision*, 2022.
- [54] M. Zhao et al. Rf-based 3d skeletons. In *Proc. of Conf. of the ACM Special Interest Group on Data Communication*, pages 267–281, 2018.
- [55] M. Zhao, T. Li, M. Abu Alsheikh, Y. Tian, H. Zhao, A. Torralba, and D. Katabi. Through-wall human pose estimation using radio signals. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7356–7365, 2018.

Paper checklist

For all authors:

- Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#)
- Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
- Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) , please check Section 5 (**Ethic statements**).
- Did you describe the limitations of your work? [\[Yes\]](#) , as a dataset paper, all our algorithms and models reported in the paper are just baseline.

If you ran experiments:

- Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#) , please check our project page. However, we only show demo first as the data and code will involve the privacy protocol as stated in Section 5. Post-processing data will delay the dataset release date but We will upload the data once we finish it before the main conference.
- Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#) , in Section 4 and supplementary materials.
- Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[Yes\]](#) , in Section 4 and supplementary materials.
- Did you include the amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) , in supplementary materials.

If you are using existing assets (e.g., code, data, models) or curating/releasing new assets:

- If your work uses existing assets, did you cite the creators? [\[Yes\]](#)
- Did you mention the license of the assets? [\[Yes\]](#) All material published is made available under the following Creative Commons license: Attribution-NonCommercial 4.0 International (CC BY-NC 4.0). We mention this in the supplement.
- Did you include any new assets either in the supplemental material or as a URL? [\[Yes\]](#)
- Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [\[Yes\]](#) , we used others’ algorithm and pre-trained models.
- Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[Yes\]](#) , please check Section 5.

If you used crowdsourcing or conducted research with human subjects:

- Did you include the full text of instructions given to participants and screenshots, if applicable? [\[Yes\]](#) , in **Participant recruitment and consent** of Section 3 and Section 5.
- Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[Yes\]](#) , we partially describe it in **Participant recruitment and consent** of Section 3 and Section 5. The full IRB approvals will be released after we confirm with school.
- Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[Yes\]](#) , those are specified in IRB approvals. Each subject gets 20 dollars Amazon gift card after completing the experiments. In total we spent 400 dollars incentives.

A Supplementary Materials

This document complements the main paper by describing: (1) results of pose estimation using additional metrics related to rehabilitation (A.1); (2) an analysis of our 3D pose refinement used to obtain ground-truth pose for our dataset (A.2); (3) details of our implementation and benchmark (A.3); (4) details of mmWave imaging (A.4); and (5) visualization of our pose estimation results (A.5).

For sections, figures, tables, and equations, we use numbers (e.g., Table 1) to refer to the main paper and capital letters (e.g., Table A) to refer to this supplement.

A.1 Further Analysis of Pose Estimation Results

We report additional evaluation metric, the mean average error (MAE) of joints angle to supplement our main results in the paper (using MPJPE and PA-MPJPE). The metric is widely considered to evaluate rehabilitation-specific movements — a main focus of our dataset. We only consider **Protocol 2** here since it has all rehabilitation-related movements.

Joints angle. We use the joint coordinates estimated by our models to find the angles between critical joints. *mRI* focuses on the four commonly used joint angles: left & right elbow angles and left & right knee angles. The elbow angle is found using the shoulder, elbow, and wrist positions. First, we obtain the bone length between the shoulder and elbow and the length between the elbow and wrist using joint coordinates. Then, the angle is calculated using triangulation from the law of cosines. Similarly, the knee angles are obtained using the hip, knee, and ankle positions. The ground truth angle is computed using the refined ground truth 3D coordinates, and we calculated MAE between the ground truth and each modality. Table A shows detailed results of joints angle MAE. We observe that under **S1**, RGB modality yields below 10° for all joints, while mmWave and IMUs lead to larger errors regarding the elbow angles ($>10^\circ$). This behavior is observed since the movement of the upper limbs is larger than that of the lower limbs for most movements. The setting of **S2** yields higher errors than under **S1** since **S2** requires the model to generalize to unseen subjects, which is arguably more challenging.

Modality	Setting	Left elbow	Left knee	Right elbow	Right knee
mmWave	S1	18.7 ± 0.2	2.9 ± 0.1	16.0 ± 0.2	3.2 ± 0.1
	S2	24.5 ± 2.3	10.4 ± 1.3	22.9 ± 2.9	11.6 ± 1.3
RGB	S1	9.0 ± 0.1	8.3 ± 0.1	9.3 ± 0.1	7.7 ± 0.1
	S2	11.5 ± 0.6	14.8 ± 1.6	11.1 ± 0.7	14.0 ± 1.5
IMUs	S1	7.9 ± 0.1	2.6 ± 0.1	11.3 ± 0.2	2.4 ± 0.1
	S2	8.4 ± 1.0	5.6 ± 0.2	9.7 ± 0.8	5.7 ± 0.2

Table A: MAE of joints angle ($^\circ$) for mmWave, RGB, and IMUs. We report the mean and standard deviation of MAE averaged across multiple splits under both our settings (**S1** & **S2**).

A.2 Analysis of 3D Joints Refinement and Quality

We provide further analysis of our 3D pose refinement used to obtain ground-truth poses for the dataset. There are three terms co-optimized in the objective function given in Equation 2. The first term represents the reprojection errors of the two cameras. The second term enforces equal bone length across all frames in the same video (i.e., the same subject). Finally, the third term imposes temporal smoothness of the 3D joint coordinates. Figure A shows an example of naive 3D pose with *poor quality* (obtained from direct triangulation), and the refined 3D pose after our optimization. We can observe that the refined pose is more stable as only the right arm moves while the lower body parts hardly move, which is the case in reality. Overall, the average objective decreases from 176 to 83, more than 50% for all subjects.

To validate the reliability of the obtained 3D joints, we further annotate a subset of the whole dataset and calculate the error. Specifically, we manually annotate 2D keypoints of the images, randomly sampled from subjects and movements. Then, we obtain the re-project 2D keypoints using refined 3D joints and camera parameters via camera calibration. Finally, we calculate the mean absolute

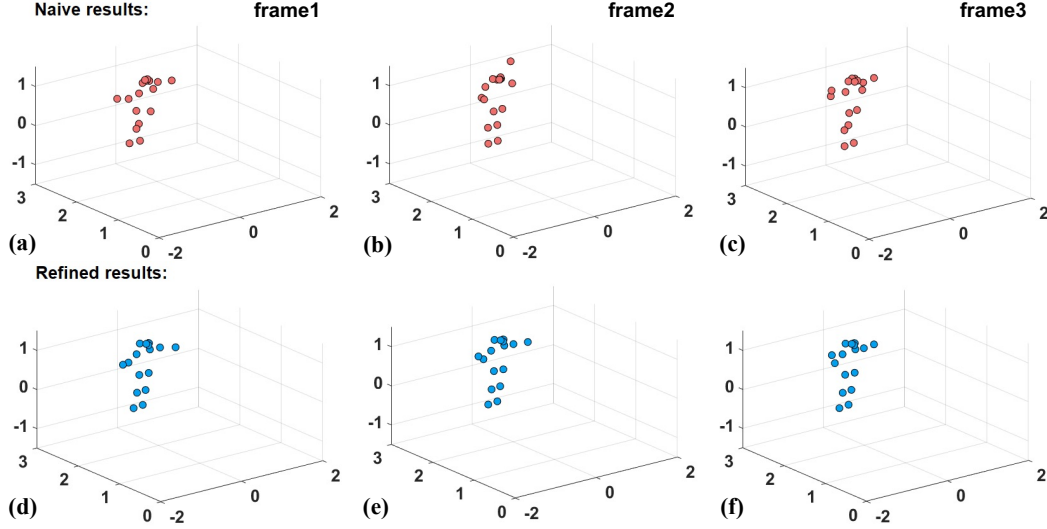


Figure A: An example of comparisons between poor naive and refined 3D pose ground truth. The subject is bending the right arm in three continuous frames from left to right. The first row shows naive results and the second row shows refined results.

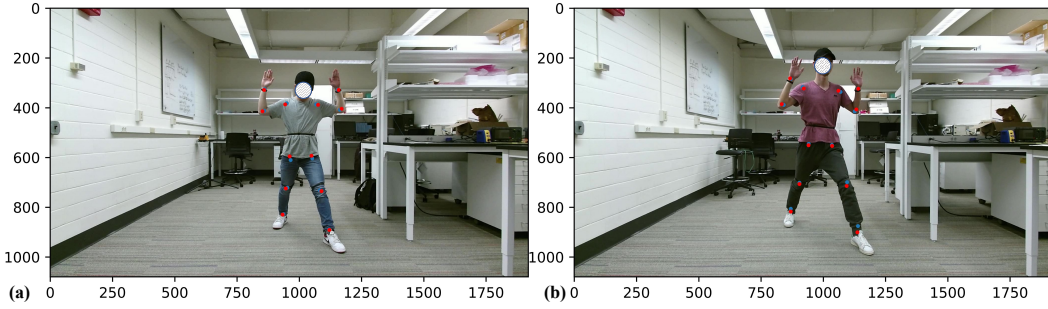


Figure B: A comparison of manual annotated 2D keypoints and re-projected 2D keypoints from refined 3D joints. Blue dots represent manual annotations and red dots show the re-projection keypoints.

percentage error (MAPE), and the percentage of correct keypoints threshold at 50% of the head segment length (PCKh) between the 2D keypoints from the model and the re-projection. The MAPE is 1.5%, and PCKh is 98.92. These quantitative results show that the proposed method of obtaining 3D joints is reasonably accurate. Figure B compares manual annotated 2D keypoints and re-projected 2D keypoints from refined 3D joints. Blue dots represent manual annotations, and red dots show the re-projection keypoints. We can observe that keypoints from the two methods almost overlap.

A.3 Benchmark and Implementation Details

We now describe implementation details of methods considered in our benchmark. We use PyTorch [32] to implement all our models. Intel Xeon Gold 6242R @ 80x 4.1GHz and NVIDIA GeForce RTX 3090 are used to train these models. The code and pre-trained models will be open-sourced to facilitate the research area³. Both raw data and synchronized data are released to the public as well. All material published is made available under the following Creative Commons license: Attribution-NonCommercial 4.0 International (CC BY-NC 4.0).

Data-split. For **Setting 1 (S1 Random Split)**, we set three different random seeds to split the data to 80% and 20% as training and testing set, respectively. For **Setting 2 (S2 Split by Subjects)**, we

³<https://sizhean.github.io/mri>

Symbol	Description	Values	Symbol	Description	Values
f_c	Starting frequency	77 GHz	θ_{res}	Angle resolution	9.55°
T_c	Chirp signal duration	32 μ s	N_{RX}	No. of RX antennas	4
B	Bandwidth	3.20 GHz	N_P	Maximum points detectable per frame	64
S	Slope of chirp signal	100 MHz/ μ s	N_{TX}	No. of TX antennas	3
N	No. of chirps per frame	96	v_{res}	Velocity resolution	0.35 m/s
d_{res}	Range resolution	4.69 cm	v_{max}	Maximum Velocity	5.69 m/s

Table B: List of major parameters and variables related to mmWave and their values for mmWave point cloud generation.

selected subset of the subjects to split the data, generated by three random seeds as well. Three different splits we used in the paper are shown as follows: (1) [17, 13, 11, 15], (2) [9, 7, 20, 8], and (3) [3, 16, 7, 2]. For example, [17, 13, 11, 15] means that subject 17, 13, 11, 15 are used for testing and the rest for training.

mmWave-HPE. We follow [3] for the implementation for mmWave-HPE model. The input layer of the CNN takes the stacked 5-channel feature tensors. Two consecutive convolution layers follow the input layer with 16 and 32 channels, respectively. After the convolutions, the output is fed to the first fully-connected (FC) layer with 512 neurons. The final output of CNN contains 51 neurons, representing 3D coordinates for the 17 joints. All activation functions are Relu except for the final FC layer, where we use linear activation. Dropout layers are used after the convolution and fully connected layers to avoid excessive dependency on specific neurons. The model converges within around 50 epochs with early stopping settings.

RGB-HPE. We employ HRNet-W32 [39] (with bounding boxes from Mask RCNN [15]) to detect 2D keypoints of human body parts in all RGB frames from both cameras. W32 in HRNet represents the width of the high-resolution nets in the last three stages. The pre-trained model from [30] is utilized for 3D joints estimation. It “lift” 2D keypoints from a sequence of frames into 3D joints.

IMU-HPE. We follow [52] for IMU calibration, normalization, and features generation. Each IMU has its own coordinate system. As a result, two steps are needed to make the output compatible with neural network models. First, *calibration*: transforming the raw inertial measurements into the same reference frame. Second, *normalization*: transforming the leaf joint inertia into the root’s space and scaling it to a suitable size for the network input. This method calculates the transition matrices for each sensor before capturing the movements, and it requires subjects to perform a ‘T pose’ before the experiments. The feature tensors extracted and transformed by this method capture the joint rotation and acceleration effectively such that multilayer perceptron (MLP) or CNN can regress the 3D joints with these features. We use a similar model as mmWave-HPE, except the input tensors are only 1-channel feature tensors for IMUs. The model converges within around 30 epochs with early stopping settings.

Skeleton-based action detection. We re-purpose an existing model [53] for the skeleton-based action detection. Specifically, the model takes a sequence of estimated 3D poses from individual modality as inputs. These poses are further encoded into a feature pyramid using a multi-scale transformer. Shared classification and regression heads check the feature pyramid, thus producing an action candidate at every timestamp.

A.4 mmWave Imaging

We follow [3] for the mmWave point cloud generation including software and hardware setup, data pre-process, and follow [4] for fusing the continuous frames point cloud to reduce the effect of sparsity. For the comprehensive details and math derivation of mmWave imaging background, please refer to [33, 42, 44, 11]. Figure C shows a sample input frame from different views. The red marker represents the radar location. Figure C(a) shows that point positions in 3D view, while the other plots show the front view, side view, and top view. Specifically, Figure C(a) illustrates the Doppler velocity, which indicates the relative velocity from the detected point to the radar. The colors in the figures represent the energy intensity of the reflected signals. Table B lists the key parameters we used to generate the mmWave point cloud.

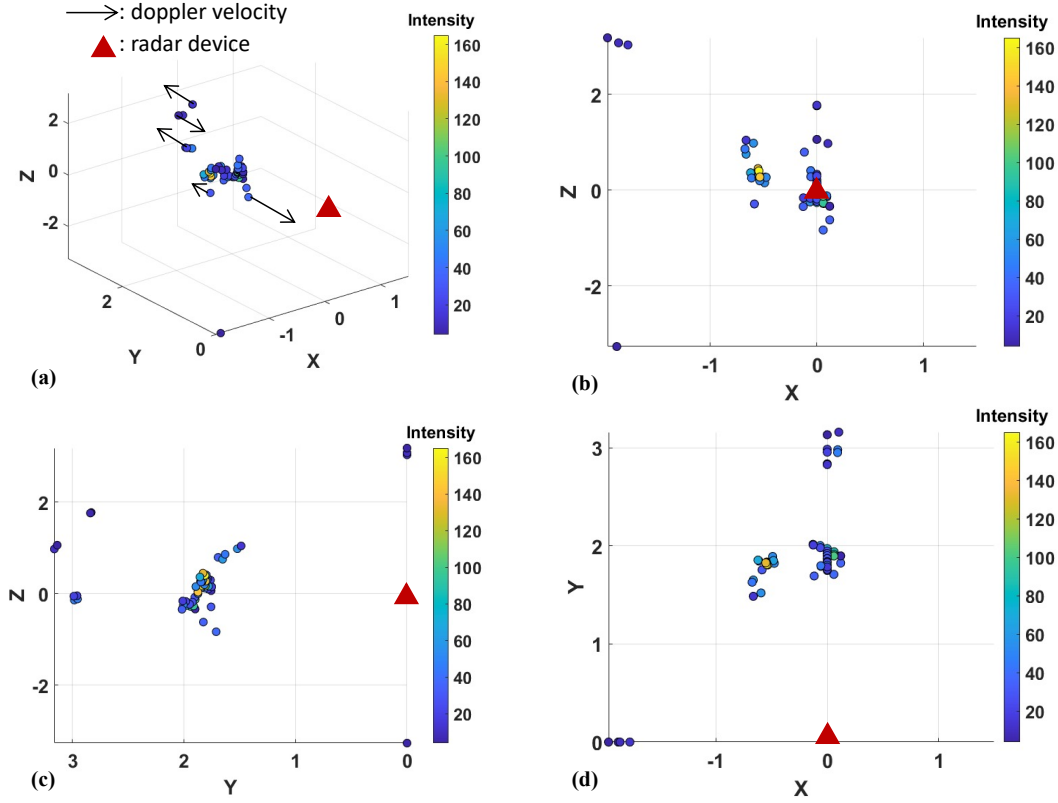


Figure C: mmWave point cloud representation for one frame. (a), (b), (c), and (d) shows the 3D view, front view, side view, and top view, respectively.

A.5 Additional Visualization

Figure D shows one subject performing right side lunge. Figure E and F demonstrate pose estimation results from different camera pose. The results are displayed with the RGB frame from the camera, the refined 3D pose, and the 3D point cloud from mmWave radar. The first row, from left to right: RGB frame with detected human bounding box and 2D keypoints, the refined 3D pose from multiple cameras, and mmWave radar point cloud signal. The second row, from left to right: estimated 3D pose from a single RGB camera, IMU signals, and mmWave radar point cloud. The captions include the action label and four commonly used joint angles: left & right elbow angles and left & right knee angles.

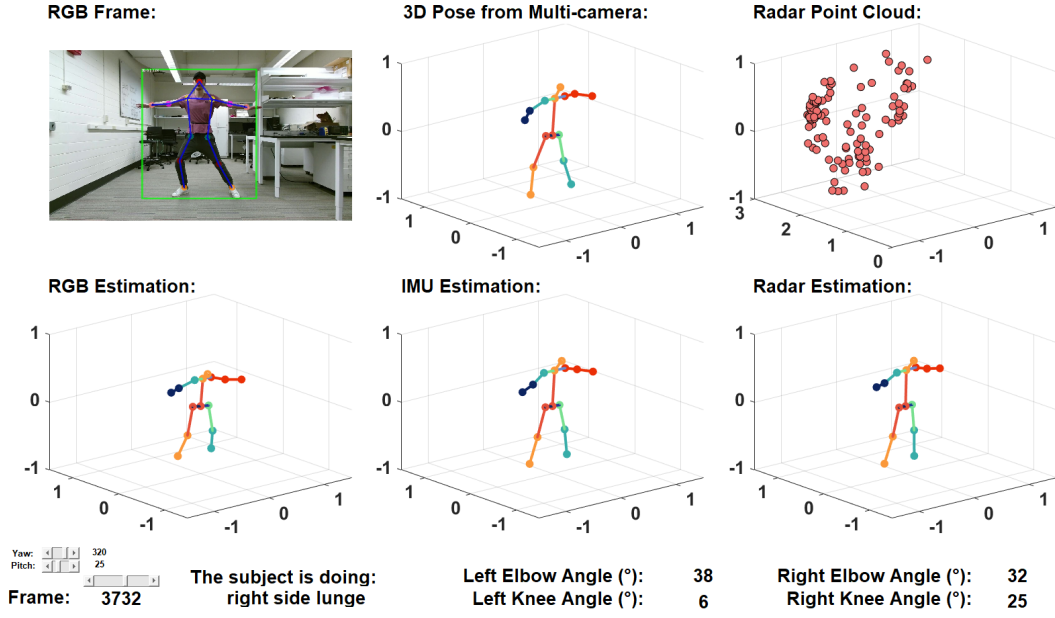


Figure D: Visualization of body poses from one subject performing right side lunge. The units in axes are meter.

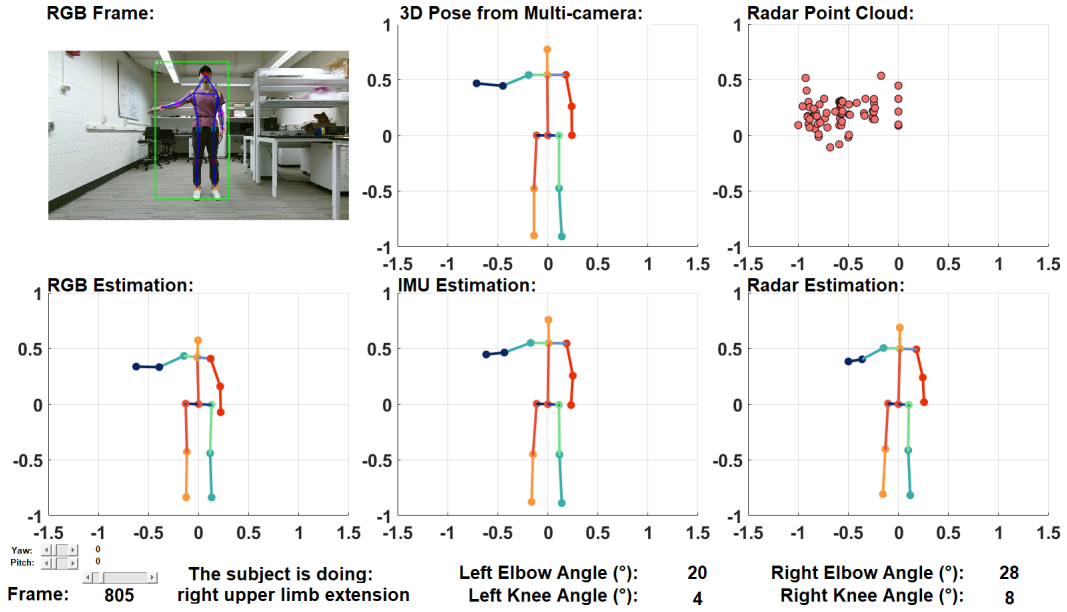


Figure E: Dataset visualization when $yaw = 0^\circ$, $pitch = 0^\circ$. The units in axes are meter.

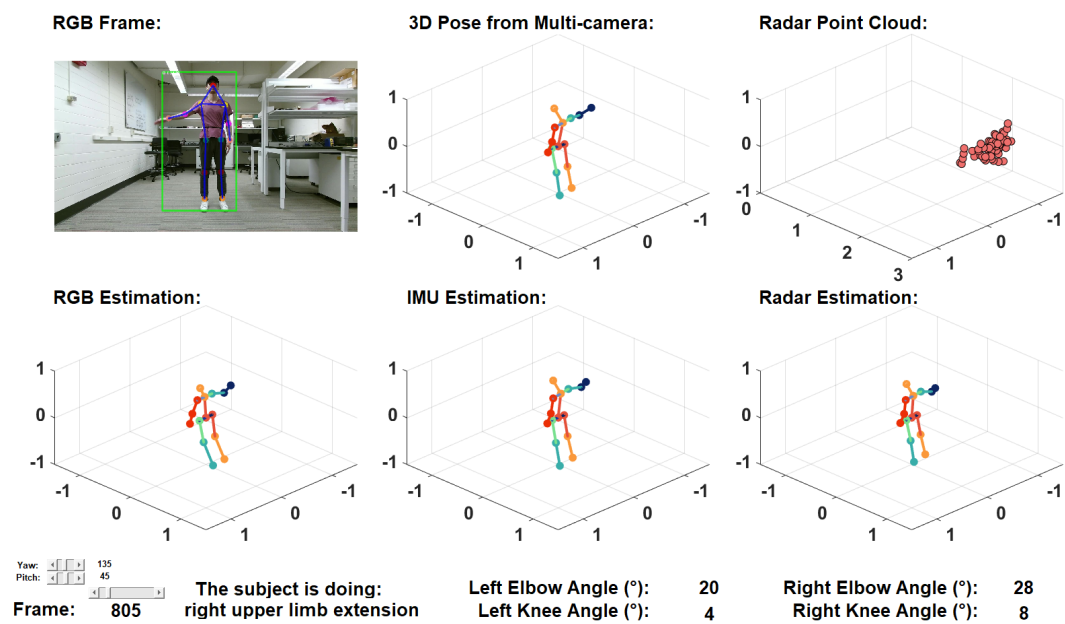


Figure F: Dataset visualization when $yaw = 135^\circ$, $pitch = 45^\circ$. The units in axes are meter.