

MDPose: Human Skeletal Motion Reconstruction Using WiFi Micro-Doppler Signatures

Chong Tang, Wenda Li, *Member, IEEE*, Shelly Vishwakarma, Fangzhan Shi, Simon Julier, *Member, IEEE*, and Kevin Chetty, *Member, IEEE*,

Abstract—Motion tracking systems based on optical sensors typically often suffer from issues, such as poor lighting conditions, occlusion, limited coverage, and may raise privacy concerns. More recently, radio frequency (RF)-based approaches using commercial WiFi devices have emerged which offer low-cost ubiquitous sensing whilst preserving privacy. However, the output of an RF sensing system, such as Range-Doppler spectrograms, cannot represent human motion intuitively and usually requires further processing. In this study, MDPose, a novel framework for human skeletal motion reconstruction based on WiFi micro-Doppler signatures, is proposed. It provides an effective solution to track human activities by reconstructing a skeleton model with 17 key points, which can assist with the interpretation of conventional RF sensing outputs in a more understandable way. Specifically, MDPose has various incremental stages to gradually address a series of challenges: First, a denoising algorithm is implemented to remove any unwanted noise that may affect the feature extraction and enhance weak Doppler signatures. Secondly, the convolutional neural network (CNN)-recurrent neural network (RNN) architecture is applied to learn temporal-spatial dependency from clean micro-Doppler signatures and restore key points' velocity information. Finally, a pose optimising mechanism is employed to estimate the initial state of the skeleton and to limit the increase of error. We have conducted comprehensive tests in a variety of environments using numerous subjects with a single receiver radar system to demonstrate the performance of MDPose, and report 29.4mm mean absolute error over all key points positions, which outperforms state-of-the-art RF-based pose estimation systems.

Index Terms—Human Skeletal Motion Reconstruction, Machine Learning, WiFi Sensing Technology

I. INTRODUCTION

WITH the rapid development of the Internet of Things (IoTs) and smart buildings in recent years, users can engage with their surroundings in a more natural way, such as through voice and pose. This requires accurately tracking human movements and interpreting them properly. One of the most promising solutions is based on computer vision and machine learning algorithms[1], [2]. However, camera based systems suffer from issues around lighting conditions, occlusion, coverage and privacy. This motivates the emergence of RF-based techniques using ubiquitous WiFi signals which offer low-cost precise sensing in a variety of scenarios whilst preserving privacy.

Many WiFi-based sensing methods use the amplitude changes of the channel state information (CSI) to analyze the human motion properties, and have validated their feasibility in the fields of indoor localisation[3], activity recognition[4], [5] and healthcare[6]. However, few CSI-based systems can use the phase information due to unsynchronised local oscillators in WiFi networks[7], resulting in a loss of the important Doppler information. More recently, the passive WiFi radar (PWR) system provides an alternative for the WiFi-based sensing tasks. There are many applications that have been successfully demonstrated by PWR systems, including occupancy detection, activity classification and hand gesture recognition[8]–[11]. Different from the CSI-based systems, PWR can express not only the amplitude of returned signals, but also the micro-Doppler frequency shifts in different body parts using micro-Doppler (μ -Doppler) spectrograms, which offers a more comprehensive perception of human movements. Furthermore, it is a passive radar system in which the signal illumination is independent of the radar itself. This implies that no modifications to the existing WiFi network are required, instead the CSI-based systems still need additional network interface cards. Therefore, PWR has the great potential to improve the performance of the current WiFi-based sensing.

However, neither μ -Doppler spectrograms nor CSI amplitudes are intuitive to ordinary users. Even if some general information, such as human activity categories, can be obtained using machine learning algorithms, more details of the motion are still not available. In this case, there is considerable interest in wireless sensing field which aims to reconstruct more detailed human skeletal motions from RF signals. Zhao et al. [12], [13] proposed RF-Pose and RF-Pose3D to achieve human skeleton estimation with only RF signals. The frameworks used RGB frames or motion capture (Mocap) data as the supervision to train CNN-based skeleton reconstruction model. From results, we can see that with the skeletal estimation, wireless sensing systems present comparable performance to traditional camera-based systems and outperform them in extreme conditions, such as through-the-wall, occlusion and dark environments. However, the RF source they use is a frequency modulated continuous wave with a 1.78GHz sweep bandwidth, which requires the construction of additional hardware equipment and may limit the application of the system in some cases. Jiang et al. [14] proposed the WiPose framework to achieve 3D pose estimation using WiFi CSI information. Their model is based on the CNN-RNN architecture and can recursively calculate poses from the estimated quaternion information. The final results demonstrated that WiFi signals

C. Tang, W. Li, S. Vishwakarma, F. Shi, and K. Chetty are with the Department of Security and Crime Science, University College London, UK e-mail: (chong.tang.18@ucl.ac.uk, wenda.li@ucl.ac.uk, s.vishwakarma@ucl.ac.uk, fangzhan.shi.17@ucl.ac.uk, k.chetty@ucl.ac.uk).

S. Julier is with the Department of Computer Science, University College London, UK email: (s.julier@ucl.ac.uk)



Fig. 1: The key points locations on human body

carry sufficient information for human skeleton reconstruction and can also tackle extreme sensing environments. However, as previously stated, the lack of phase information may impede further developments of CSI-based systems. Furthermore, how to initialise the first pose and deal with the accumulation of errors in long-term estimation have not been fully discussed in [14]. Therefore, the system still has potential for improvement.

The pioneering efforts by researchers in this area has inspired us to develop MDPose, a novel framework for human skeletal motion reconstruction using the commercial WiFi network. After removing noise from measured spectrograms, MDPose can extract velocity information from μ -Doppler spectrograms for up to 17 skeletal key points (as shown in Fig. 1) and recursively calculate poses. Furthermore, MDPose has a pose optimising mechanism that can obtain optimisation vectors based on the current pose and the future's motion to address issues related to the initial pose estimation and the long-term error accumulation. In particular, MDPose enables training the velocity estimation network with simulated spectrograms due to our denoising strategy, which can effectively address the training difficulties caused by the lack of measurement data. Finally, various experiments were carried out in different environments and with different subjects, successfully demonstrating the effectiveness and generalisation performance of MDPose.

The rest content of our paper is organized as follows: First, Section II briefly describes the framework. Next, Sections III to V introduce details of three main components of MDPose. Then Section VI explains the dataset, hardware and software we used in experiments, and the training settings of neural networks. Finally, we present experimental results and important discussions in Section VII, and conclude the work in Section VIII.

II. SYSTEM OVERVIEW

MDPose is based on μ -Doppler information and allows the reconstruction of human skeletal movements without any modifications of the existing WiFi network. To address a series of challenges in practical applications, we design a novel framework to handle them step by step.

- **Data Collection:** To collect μ -Doppler changes, MDPose uses the PWR system in which the signal illumination i.e. WiFi access point is independent of the radar itself. This implies that no modifications to the existing WiFi network are required. Compared to CSI systems that operate using NIC cards, it can be more easily adapted to a wide range of scenarios.
- **Data Denoising:** Interference and noise may mask important signals and thus affect the extraction of dynamic characteristics. Therefore, the data denoising is the first step in MDPose. In this step, a deep learning-based denoising network is implemented to effectively remove unwanted signals.
- **Velocity Estimation:** There are many choices to describe the motion property, such as the position, quaternion, etc. For the μ -Doppler, it is strongly related to velocity profiles. Therefore, in our case, estimating velocity could have more accurate results compared with using other types of data. So, we propose a CNN-RNN neural network to extract velocity information.
- **Pose Optimising:** Although we can obtain the velocities from the previous step, we do not know the initial pose p_0 so that we cannot calculate the following poses. In this case, we propose a pose optimising mechanism to help estimate a proper p_0 . Then with p_0 and velocities, we can finally calculate the following poses based on the relationship between the velocity and the position. However, unlike the quaternion-based approach[14], there is no strong constraint between the velocities of the different joints, and therefore unrealistic skeletons are likely to occur as time increases. To address the issue of accumulative error, the pose optimising method can be used again, which can recursively optimise the current pose according to the velocities and poses afterwards.

The overall framework of MDPose has been presented in Fig. 2, and more details will be discussed in the following sections.

III. CLEAN UP MICRO-DOPPLER SPECTROGRAM

PWR includes separately located reference and surveillance channels. The reference channel captures the source signals, while the surveillance channels gathers signals reflected off targets but may also receives reflections from environmental clutters, etc. In this section, we introduce the signal processing methods we used to extract the desired μ -Doppler spectrogram, and introduce useful denoising algorithms to clean up the measured data.

A. PWR Signal Processing

If we denote the source signal as $u(t)$, then the signal received by reference channel $S_{ref}(t)$ can be expressed as:

$$S_{ref}(t) = A_{ref}u(t - \tau_{ref}) \quad (1)$$

where $S_{ref}(t)$ is the copy of $u(t)$ with a delay τ_{ref} and amplitude scaling factor A_{ref} . For the surveillance channel,

its received signal S_{sur} can be expressed as:

$$S_{sur}(t) = \sum_{l=1}^N A_l u(t - \tau_l) e^{j2\pi f_l t} + \sum_{m=1}^M A_m u(t - \tau_m) e^{j2\pi f_m t} + A_{DSI} u(t - \tau_{DSI}) + \sum_{n=1}^N A_n u(t - \tau_n) \quad (2)$$

where A and τ still are the amplitude scaling factor and the delay that belong to different returns. And we will have four terms that represent the reflected signals from the L targets, the multipath reflection, the direct signal interference (DSI) and the returns from stationary clutter objects in the surrounding environment, respectively. Specifically, we also have Doppler shift f_l caused by the movement of targets which is also related to the target velocity profiles, and another Doppler shift f_m due to the multipath effect. Next, we can extract μ -Doppler and range information by cross-ambiguity function (CAF) as shown in the following:

$$CAF(\hat{\tau}, \hat{f}_D) = \int_{-\infty}^{\infty} S_{sur}(t) S_{ref}^*(t - \hat{\tau}) e^{-j2\pi \hat{f}_D t} dt \quad (3)$$

where $\hat{\tau}$ and \hat{f}_D are the expected target delay and Doppler shift, respectively, and the $[*]$ is the complex conjugate. This will output CAF results as:

$$CAF(\hat{\tau}, \hat{f}_D) = \sum_{l=1}^N CAF_{tgt}^l + \sum_{m=1}^M CAF_{multipath}^m + CAF_{DSI} + \sum_{n=1}^N CAF_{clutter}^n \quad (4)$$

In Equation 4, we can see that the result does not only contain the desired target Doppler-range information, but also have noise caused by multipath reflections, DSI and environmental clutters. Among them, the DSI has the greatest impact on the results, which may mask reflected signals from targets and cause unwanted peaks at the zero-Doppler frequency. To eliminate it, we then apply the CLEAN algorithm as presented in the below:

$$CAF'(\hat{\tau}, \hat{f}_D) = CAF(\hat{\tau}, \hat{f}_D) - \alpha CAF_{self}(\hat{\tau} - T_k, \hat{f}_D) \quad (5)$$

where $CAF_{self}(\hat{\tau} - T_k, \hat{f}_D)$ is the CAF over the reference channel and α is the maximum absolute value of $CAF(\hat{\tau}, \hat{f}_D)$. Due to the insufficient range resolution, after the CLEAN process, the PWR system only takes μ -Doppler information i.e. velocity profiles from the CAFs and combines them along time axis to generate the final μ -Doppler spectrogram. However, even the DSI can be removed, the multipath noise and the attenuation of the signal with distance may still affect the quality of the spectrogram. Therefore, other effective denoising algorithms are required to further clean up the data.

B. Latent Feature-wise Mapping Network

The latent feature-wise mapping network (FMNet) proposed in our work[15] combines strengths of variational auto-encoder and adversarial auto-encoder that maps noisy measurement latent features to a clean space to achieve the noise reduction. In short, FMNet will initially learn how the various motion attributes are spread over a clean space and how to restore data from the clean space. Given a noisy measured data,

it can then find the latent features in this space that have the closest motion information to it. After that, the decoder structure of FMNet will restore the μ -Doppler spectrogram based on the found clean features, and the interference from environmental clutters and multipath noise can be effectively removed. In [15], we have demonstrated that compared with other networks[16], [17], FMNet has better generalization ability, up to 10% improvement, which still performs well in some challenging scenarios, such as through-the-wall sensing.

Furthermore, to construct the clean space, we used the simulation data matched to the measurement data to train the FMNet, so the final outputs are simulation-like. More details about FMNet can be found in [15].

IV. VELOCITY AND POSE ESTIMATION

To reconstruct the skeletal motion, an intuitive solution might be directly learning positions of each key points from a μ -Doppler spectrogram. However, common WiFi bandwidth is from 20 to 40MHz, corresponding to a range resolution of around 7 to 3.5meters, which is not sufficient to distinguish between different parts within the body range (normally less than 2meters). Moreover, independently considering positions may lead to the discontinuity in reconstructed movements. On the other hand, the μ -Doppler spectrogram is directly related to the velocity profiles: the position of a Doppler bin refers to the direction of a target movement and how fast it moves; the strength of a Doppler bin is determined by radar cross section, distance between the target and receivers, the speed overlay and other complex factors. Therefore, learning velocities is more consistent with the characteristic of μ -Doppler spectrograms. After obtaining velocities, we use a simple update equation to build the time dependency between poses, as the following:

$$p_t = p_{t-1} + v_t \delta t \quad (6)$$

where p and v are pose and velocity vectors, their subscripts represent frame indices and δt is the time interval between two frames.

A. Neural Network

To learn velocity sequence from a spectrogram, we develop a CNN-RNN architecture, as presented in the velocity estimation network in Fig. 2. We first use 1-dimensional convolutional layers over each frame to extract spatial features. These features are then concatenated along time axis and fed into a two-layer Long Short-Term Memory (LSTM) RNN to learn their time dependency. Finally, the outputs of each time steps are passed through fully-connected layers (FCs) to generate the velocity sequence.

More specifically, we set the kernel size of convolutional layers to 5×5 for the fine-grained feature extraction. Each convolutional layer is followed by a batch normalisation and rectified linear unit to accelerate training speed and add non-linearity to outputs of this layer. The final layer does not use any non-linear activation function and directly outputs estimated velocity values.

For training the network, we aim to minimize the difference between estimated velocities and the ground-truth velocities

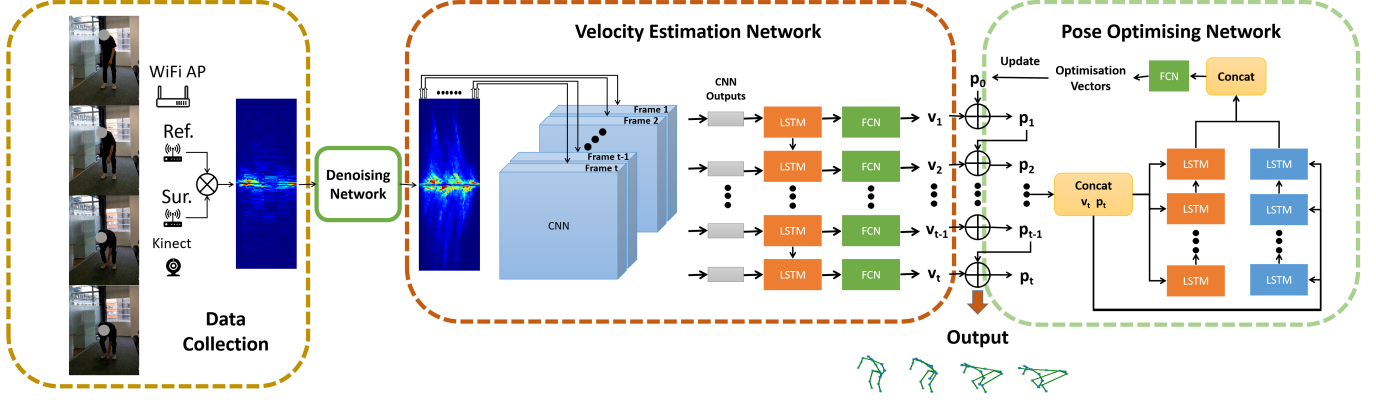


Fig. 2: The architecture details of MDPose

for each frame and key point. Therefore, the loss function can be defined as:

$$Loss = \frac{1}{T} \frac{1}{N} \sum_{t=1}^T \sum_{i=1}^N |v_t^{i'} - v_t^i| \quad (7)$$

where T is the length of frames, N is the number of key points, $v_t^{i'}$ and v_t^i are estimated and real 3D velocity of i^{th} key point at frame t .

B. Training Data

As a data-hungry algorithm, training a deep neural network often requires large volumes of data to retain strong generalisation performance. However, collecting experimental data for radar sensing systems is time-consuming and laborious. In Section III, we propose to use FMNet to remove environmental and multipath noise from measurement spectrograms. Furthermore, due to the denoised data is simulation-like, it is feasible to only use simulation spectrograms as the training data, and then apply the trained model to the denoised spectrograms. We have demonstrated in [15] significant classification improvement in practical applications by using this training strategy. Furthermore, we can easily produce a large amount of simulated data, so this method can greatly ease issues caused by the insufficient measured data. Meanwhile, to the best of the authors' knowledge, this is the first implementation of this approach.

V. POSE OPTIMISING MECHANISM

Once we have extracted velocities from μ -Doppler spectrograms, new questions arise: what is the initial pose p_0 for deriving the subsequent poses and how can we limit the accumulation of positional errors over time? For the first question, due to the restrictions of WiFi bandwidth, it cannot provide sufficient range resolution to determine positions of different body parts directly from the radar returns. However, a faulty p_0 may result in a sequence of illogical postural shifts. This can provide us clues about whether the current p_0 needs to be optimised. Furthermore, if we can make the correction periodically during a long-term motion reconstruction, we can effectively avoid the accumulation of positional errors. In this section, a bi-directional LSTM-based pose optimising mechanism will be introduced.

A. Optimisation Vector

Although conventional supervised learning methods have ability to directly predict p_0 from μ -Doppler or velocity profiles, it may lead to over-fitting and requires a significant amount of training data to maintain good generalisation performance. Due to the uncertainty of p_0 , we can first initialise it based on some empirical guesses, such as a person standing before sitting, or a universal initialisation like T-pose, indicated as p_0' . The challenge then becomes determining a vector p_{diff} such that $P_0 = P_0' + P_{diff}$. But again, because of possible over-fitting issue, directly learning p_{diff} is risky. Therefore, we decompose the problem and gradually approach the optimal p_0 by learning optimisation vectors (\hat{OV}) that is defined as:

$$\hat{OV}_i = \frac{\overrightarrow{p_0'(i)p_0(i)}}{\|p_0'(i)p_0(i)\|} \quad (8)$$

where $\|*\|$ is the norm of a vector. This vector can let us optimise p_0 by moving it with a short distance in the appropriate direction, and enable us to approach the optimal position by repeating this process. This method has a high fault-tolerance and can effectively prevent over-fitting problems.

B. Learn Optimisation Vector

Given $V = [v_1, v_2, \dots, v_t]$ and an initial pose, the pose sequence $P = [p_1, p_2, \dots, p_t]$ can be calculated with Equation 6. As discussed at the beginning, this sequence might be illogical because of a erroneous p_0 but it could provide clues about how to find \hat{OV} . Therefore, we construct an optimising model that includes a bi-directional LSTM and FCs to estimate \hat{OV} based on V and P , as shown in the pose optimising network in Fig. 2.

For creating training dataset, we randomly select a state from the entire Mocap measurement as p_0' and generate a P sequence. Then we define the input features of the network as the element-wise concatenation of V and P , $\begin{bmatrix} v_1 & \dots & v_t \\ p_1 & \dots & p_t \end{bmatrix}$, while the corresponding label is \hat{OV} . Furthermore, the

objective function is defined as the following:

$$\min Loss = \frac{1}{N} \sum_{i=1}^N (1 - \cos(\frac{\hat{O}\hat{V}_i' \cdot \hat{O}\hat{V}_i}{\|\hat{O}\hat{V}_i'\| \|\hat{O}\hat{V}_i\|})) + \frac{1}{N} \sum_{i=1}^N (1 - \|\hat{O}\hat{V}_i'\|)^2 \quad (9)$$

where $[\cdot]$ refers to dot product, $\hat{O}\hat{V}_i$ is the predicted vector for the i^{th} joint. For the first term, we aim to minimize the angle between the $\hat{O}\hat{V}_i'$ and $\hat{O}\hat{V}_i$, and the second term is used to limit the length of the vector.

In practice, we can select numerous p_0' for the same velocity sequence to generate more training data, which can highly enhance the model's generalisation performance.

C. Algorithm Overview

Based on this network, the pose optimising mechanism is proposed as shown in Algorithm 1. We also define a hyperparameter, optimization rate (optr), to adjust the step of each optimisation, which is 0.01 in our experiments.

Algorithm 1 Pose optimising mechanism

Require: the previous initial pose, $p_0'^{-}$; velocity sequence, V
while not converged **do**
 $P \leftarrow V$ and $p_0'^{-}$
 Network input \leftarrow concatenate P and V
 $\hat{O}\hat{V}' \leftarrow$ optimising network
 Updated initial pose $p_0'^{+} \leftarrow p_0'^{-} + \hat{O}\hat{V}' \times optr$
end while

Apart from the initialization of the first frame, we can also apply this mechanism periodically during a long-term skeletal reconstruction. In this case, rather than beginning from a random pose, we will use the pose from the previous frame as a starting point and optimise it, which can significantly lessen the effect due to positional error accumulation.

VI. DATA COLLECTION AND EXPERIMENTAL DETAILS

A. Dataset

To demonstrate the performance of MDPose, we carried out various experiments to collect data from 7 subjects and 3 different environments. In the dataset, activities include walking towards the receiver (W+), walking away the receiver (W-), turn-around (TR), sit-down (SD), stand-up (SU), aggressive hitting (HT), passive covering (CV), pick-up (PU) and body rotation (BR). We also recorded some combined movements, such as $W+ \rightarrow TR \rightarrow SD$ and $SU \rightarrow W- \rightarrow PU$. Each is completed within 5 to 10 seconds and the length of the whole measurement is around 190 minutes, resulting in over 2000 individual activities (also called 2000 samples). Furthermore, we deployed a Kinect sensor to simultaneously collect the ground-truth Mocap data for labelling the measurement data and generating simulation data. Therefore, every measured spectrograms has associated matched velocity profiles and simulated spectrograms. Some data examples are shown in Fig. 4.

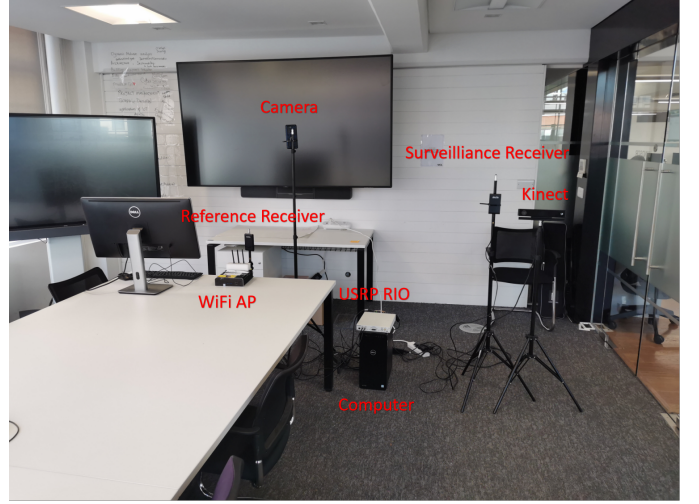


Fig. 3: The experimental setup of PWR system

We can see that the dataset contains various activities, such as movements with continuous changes (W+, W-), quick and strenuous movements (HT, CV), movements with the similar μ -Doppler pattern (SD, HT), etc. This can provide comprehensive scenarios to validate the robustness of the framework. Furthermore, when we split the data into training and testing sets, apart from common splitting strategies i.e. splitting the whole dataset with 70%/60%/50% split rates, we divided the dataset with the split rate of 23%. This training set only has samples from one person and one environment with the length of around 44 minutes, while the testing set contains data from different subjects and different environments. This extreme situation is to test whether MDPose still performs well with different subjects and in new environments.

B. SimHumalator

To generate the simulation spectrograms based on Mocap data, we used SimHumalator, which is an open source simulation tool that can produce human μ -Doppler spectrograms for the PWR sensing scenarios. It can be downloaded for free from <https://uws1.co.uk/> and more details about SimHumalator are presented in [18]. During data generation phase, SimHumalator uses the same WiFi standard and PWR parameter settings as the real experiments to guarantee that the same motion information is conveyed in the measured and simulated data.

C. PWR Experimental Setup

As shown in Fig. 3, the PWR system uses one surveillance channel that is placed with the Kinect sensor together to simultaneously capture subjects' radar returns and velocity information, and uses one reference channel that is placed close to the WiFi AP to receive reference signals. We have implemented the PWR system in three different experimental environments: a broad living room, a narrow living room and a broad meeting room. In each case, the geometrical layout of the surveillance and reference channels are not identical but are relatively similar. This is to prevent having too much difference

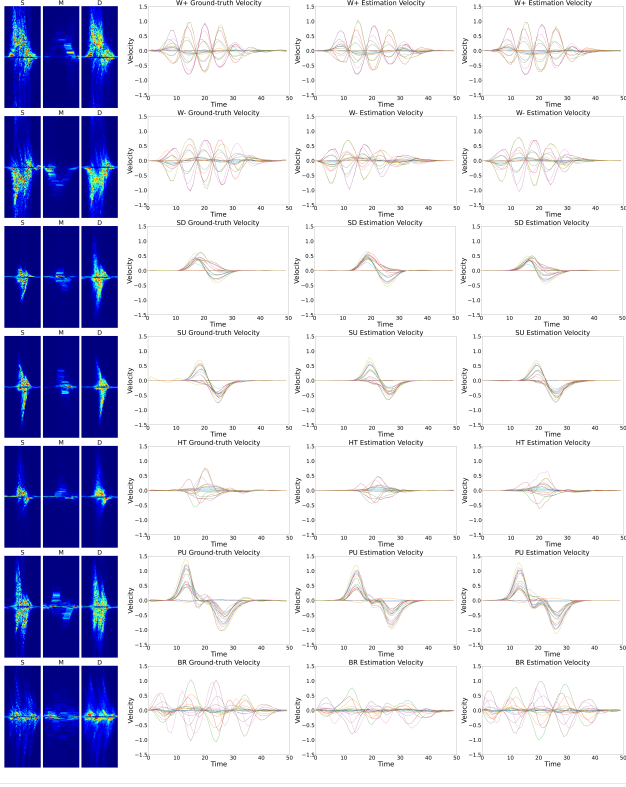


Fig. 4: The examples of velocity estimation results of different activities. From the left to right, they are simulation spectrogram, measurement spectrogram, denoised spectrogram, ground-truth velocity plot for 17 key points, MDPose (M)-estimated velocity plot and MDPose (D)-estimated velocity plot.

in the measurement data. On the other hand, a NI USRP RIO software-defined radio is connected to the two channels for the channel synchronization and real-time signal acquisition.

D. Training Networks

All neural networks are implemented with Pytorch[19] library on NVIDIA Quadro RTX 4000 graphics processing unit, and the architecture details and the training hyper-parameter settings of each networks have been listed in Table I.

VII. EXPERIMENTAL RESULTS

In this section, a thorough evaluation of MDPose will be carried out both quantitatively and qualitatively. Comparisons with state-of-the-art approaches will also be made.

A. Velocity Estimation

As introduced in Section III, spectrograms will be first denoised with FMNet, which aims to remove unwanted interference and reconstruct clear motion details. Furthermore, it enables training the network with the simulation data. To demonstrate the advantages of this strategy, we compared its velocity estimation results with those obtained from networks trained directly with measurement data. We denote simulated,

measured and denoised spectrograms as S , M and D , respectively.

In Fig. 4, the first column presents the matched S , M and D of different activities. From M s, we can observe the overall direction of movement as well as the magnitude of velocities of different body parts. For example, in $W+$, the subject approaches the receiver from a distance. Therefore, the overall magnitude of the velocity increases first and then decreases. Meanwhile, when the person gets closer, the power of the received signal rises, making the Doppler pattern more visible. Such properties are important for applications like localisation and ranging, however, when retrieving human motion, the weak signal might be masked by noise, resulting in information loss. In contrast, the effects of signal attenuation and noise are limited in S . So, we can observe a distinct and complete Doppler pattern, which is beneficial for the human skeletal reconstruction task. And this is also why we made denoising and data enhancement as the first step of MDPose. From D s, we can observe that FMNet can effectively clean M s and generate data that is close to corresponding S s. Although some local regions remain blurred, the enhanced features can significantly improve the performance of velocity estimates, which can be demonstrated from velocity plots.

The second to fourth columns in Fig. 4 are plots of the ground-truth velocity, the estimated velocity based on M and the estimated velocity based on D , respectively. Overall, MDPose can successfully extract velocity information from both M and D , and most results have consistent patterns with the ground-truth. However, we can see that the D -based estimation has better performance than the M -based estimation, most notably in the plots for $W-$ and BR . We can observe that in the weak signal sections, the corresponding M -based estimations have considerably lower velocity amplitudes than they should, which is because the motion features are masked by the background noise and the model fails to extract them. On the other hand, after being processed by FMNet, D -based estimations do not have the same issue and all of them are qualitatively close to the ground-truth plots. Furthermore, since the D -based results are obtained using a CNN-RNN model trained on S , the difference between the D - and M -based performances might be more evident in the case of a limited training set.

For the quantitative analysis, we present the velocity errors of each key point in Table II. We compared the mean absolute errors between M - and D -based results. It shows that the D -based results have less errors than the M -based results for most of key points, and even with some exceptions (3, 6, 7, 9), their errors are still at a low level. For some points, the D -based results can achieve maximum $4.1mm/frame$ (i.e. around $41mm/s$) improvement, and this difference would significantly affect the long-term motion estimation. Finally, the overall error of D -based MDPose is $6.8mm/frame$ which is $0.9mm/frame$ better than M -based MDPose.

To sum up, we have demonstrated that MDPose can extract velocity properties from μ -Doppler spectrograms. Meanwhile, due to clearer features, the D -based estimation normally outperforms the M -based estimation. However, D cannot completely restore the details in S , and there are still some

Network	Layer	Parameter Setting	BatchNorm	Activation	Training Setting
Velocity Estimation	Conv1d	filters 32, kernel size 5x5, stride 2x2, padding 0	Yes	ReLU	optimizer: Adam learning rate: 0,001 batch size: 64
	Conv1d	filters 64, kernel size 5x5, stride 2x2, padding 0	Yes	ReLU	
	Conv1d	filters 64, kernel size 5x5, stride 2x2, padding 0	Yes	ReLU	
	LSTM	hidden size 64, number of layers 2, bidirectional True	No	-	
	Linear	input size 128, output size 128, bias True	Yes	ReLU	
	Linear	input size 128, output size 51, bias True	Yes	ReLU	
Optimisation Vector Estimation	Linear	input size 51, output size 51, bias True	No	-	optimizer: Adam learning rate: 0,001 batch size: 128
	LSTM	hidden size 51, number of layers 2, bidirectional True	No	-	
	Linear	input size 102, output size 256, bias True	No	ReLU	
	Linear	input size 256, output size 51, bias True	No	Tanh	

TABLE I: The architecture and training details of the two networks

Joints	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	Overall
M	0	1.6	6.0	8.5	1.9	5.6	8	2.8	4.5	7.7	14.5	15.8	7.6	14.8	18.5	5.7	7.6	7.7
D	0	1.5	6.1	8.5	1.5	6.0	8.4	2.6	4.9	7.2	10.5	11.7	7	12.4	15.8	5.1	6.9	6.8

TABLE II: Velocity estimation errors of each key points (unit: mm/frame)

ambiguities. In this case, we may further improve the model's robustness by adding noise to S during the network training. On the other hand, we believe other kinds of denoising algorithms can also improve the estimation performance and will be the subject of future research.

B. Pose Estimation

After obtaining the velocity sequence, we can use the pose optimising method to initialise the start pose and calculate the subsequent sequence of poses.. We presented some frames in Fig. 5 for the qualitative analysis.

Video frames are illustrated on the first row which are followed by the matched ground-truth poses (blue), M -based poses (yellow) and D -based poses (green). Meanwhile, we have circled the areas where the error is greater than 100mm. We observe that MDPose can accurately calculate the pose sequence from the velocity and initial pose estimations. For D -based results, most of them are consistent with the ground-truth poses, and matched with the real video recordings. Even if there are relatively large errors in some frames, they do not affect the overall performance and are still within acceptable limits. Moreover, comparing M - and D -based results, we can observe that the failure of the velocity estimation due to the low quality μ -Doppler spectrogram significantly affects the pose estimation, especially in the third activity BR. This once again illustrates the importance of denoising and feature enhancement before feeding spectrograms into the network. Additionally, we can observe a phenomena: since the subsequent sequences are derived using the prior pose and velocity, large errors in the earlier pose are likely to be inherited by the future sequences, resulting in a continuous accumulation of errors, such as the last two frames of M -based results in W+, the last two frames of D -based results in BR and all frames of M -based results in BR. This, therefore, will have an impact on MDPose performance in terms of long-term skeletal motion estimation. In this case, we can address this issue with the optimising mechanism, which will be discussed in Section VII-C.

Apart from the qualitative results, we also presented the quantitative analysis in Table III. It thoroughly compares pose errors of M - and D -based MDPose results over different

activities and different key points. Again, the D -based results are superior compared to the M -based results in most cases. Particularly for CV and BR, key points 13 and 15 of M -based results have more than 200mm errors, which have a significant impact on skeletal reconstruction. By contrast, in these cases, the D -based estimations maintain the low level of errors. Furthermore, by comparing overall errors of different activities and key points, all of the D -based results have less errors than the M -based results, with the mean absolute error of around 29mm, while it is around 44mm for the M -based results.

Furthermore, we also compared MDPose results with the state-of-the-art results, as shown in Table IV. We can observe that MDPose (D) reduces the error by around 14.2mm when compared to RFPose3D (CSI), considerably improving the reconstruction accuracy. However, the MDPose (D) still performs slightly worse than the WiPose (CSI), with a difference of 1.1mm. This is because the current single surveillance-channel setup of the MDPose does not compare favourably with the multi-receiver WiPose. With the addition of multiple channels in the future, we believe there is great potential to improve the performance of the MDPose. Additionally, although the performance of MDPose (M) is somewhat worse relative to the other results, it can still achieve comparable performance to RFPose3D (CSI), which could also indicate the effectiveness of the MDPose framework.

C. Pose Optimising

In our work, the pose optimising mechanism has two uses: initial pose estimation and error reduction in a long-term estimation.

For the initial pose estimation, we first randomly initialize a pose as the frame 0, then the optimising mechanism can help it gradually approach the ideal starting pose, and finally it can be used to calculate pose sequence with the estimated velocities. We presented four examples in Fig. 6. The blue skeleton with blue nodes is the ground-truth initial pose while the green skeleton with orange nodes is the estimated initial pose of each optimisation epoch. In practice, we used the same pose as the initialisation, as shown in the first column in Fig. 6. We totally have 50 optimisation epochs and the rest of

Joints		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	Overall
W+	M	0	11	21	18	13	29	27	17	21	24	35	35	31	12	21	29	23	22
	D	0	10	21	29	9	22	25	12	13	13	23	15	23	20	11	20	31	17
W-	M	0	14	24	31	11	25	36	13	14	30	29	30	32	46	46	24	28	23
	D	0	12	16	31	6	23	36	13	13	14	24	21	28	26	42	20	25	21
SD	M	0	10	54	54	61	28	45	61	25	26	32	50	33	33	33	30	30	34
	D	0	5	32	35	9	3	53	26	28	24	38	41	22	33	38	30	26	26
SU	M	0	13	24	28	8	25	24	23	25	25	27	27	27	21	20	28	25	22
	D	0	4	28	24	17	27	20	16	18	13	22	25	19	17	16	22	16	18
HT	M	0	41	20	42	44	27	52	31	25	47	44	133	95	41	118	18	27	47
	D	0	42	20	28	48	25	39	19	20	25	32	57	16	31	113	19	20	17
CV	M	0	8	70	120	13	57	91	20	37	37	53	158	281	146	237	39	65	84
	D	0	20	74	117	12	46	83	24	35	43	52	71	60	46	82	37	59	51
PU	M	0	15	31	35	22	27	26	32	39	39	40	45	39	43	44	47	53	52
	D	0	21	30	25	11	29	26	21	31	30	33	45	41	49	50	40	40	31
BR	M	0	63	54	36	71	29	52	30	26	42	92	147	112	153	205	37	57	71
	D	0	21	16	17	15	25	19	11	22	25	80	85	42	92	100	29	36	37
Overall	M	0	22	37	46	30	31	44	28	27	34	44	78	81	62	91	32	39	44
	D	0	17	29	38	16	25	41	18	23	23	38	45	31	39	57	27	32	29

TABLE III: Pose estimation error Comparison (unit: mm)

RFPose3D (CSI)	WiPose (CSI)	MDPose (M)	MDPose (D)
43.6	28.3	44.3	29.4

TABLE IV: Pose reconstruction error comparison with the state-of-the-art results

Training Rate	23%	50%	60%	70%
RFPose3D (CSI)	-	52.4	48.6	43.6
WiPose (CSI)	-	34.6	32.1	28.3
MDPose (D)	57.1	40.2	35.7	29.4

TABLE V: Pose reconstruction comparison with different training rates

columns present the results after processing every 10 epochs. For the first row, the activity is SU, so the ideal initial pose is sitting-down. Therefore, when we initialise it with standing-up pose, the leg parts have large difference. However, as the optimisation proceeds, we can observe the legs progressively being adjusted to the proper positions. At the same time, the arms have also been slightly optimised to a better position. For the remaining points, they reach stability in the early epochs and are not affected by the other points. We observe similar characteristics in the rest of activities. Furthermore, for the third activity W+, the initialisation is quite close to the ground-truth pose. We can see that the optimisation algorithm still has a good handle in this case, and each point remains stable after a slight adjustment.

We quantified the mean absolute error between the estimated initial attitude and the ground truth initial attitude for the above activities, as shown in Fig. 7. From the plot, we can see that the errors start at relatively high values and gradually decrease with the optimisation process. Most of them can eventually be less than 20mm.

On the other hand, the optimising mechanism can also be used in a long-term skeletal motion estimation to reduce the effect of the accumulative error. To test its performance, we collected a long μ -Doppler spectrogram with around 350 frames (i.e. 35 seconds). The sequence of activity is $SU \rightarrow W+ \rightarrow PU \rightarrow BR \rightarrow W- \rightarrow SD \rightarrow SU$. From the blue line in Fig. 8, we can observe that without optimisation, previous errors can continually affect later estimations, resulting in the errors becoming increasingly larger. This phenomenon matches the discussion in Section and Fig. 5. By contrast, after applying optimising mechanism every 50 frames as pointed by green arrows, the error can be periodically reduced so that the

overall error fluctuates within an acceptable range.

According to above analysis, we have demonstrated that the proposed optimising mechanism is an effective method, which is not only an essential step in the skeletal reconstruction, but also plays an important role in further enhancing the quality of the estimations. Also by introducing the concept of $\hat{O}V$, the difficulty of training the model and the risk of overfitting are greatly reduced.

D. Other Evaluation

Generalisation Performance: to validate the generalisation ability of MDPose under different indoor sensing scenarios, we have three experimental environments and 7 subjects. We presented some examples of two different subjects in Fig. 9. As we can observed, the MDPose can still successfully estimate poses of different subjects, and the results are consistent with the video frames and the ground-truth poses.

Training Rate: we also demonstrated how the error changes with the decrease of the training rate and compared MDPose (D) with the state-of-the-art results. We can see that the mean absolute error increases as the amount of training data decreases in all frameworks. Among them, MDPose outperforms RFPose3D (CSI) in all cases and is comparable to WiPose (CSI). In general, MDPose errors are maintained to a low level with a little impact on skeletal reconstruction task. Meanwhile, this result already can provide adequate and intuitive motion information to WiFi-based sensing scenarios. Furthermore, we also tested the robustness of MDPose by using only one subject's data for training and using other subjects' data for testing, and finally obtained the mean absolute error of

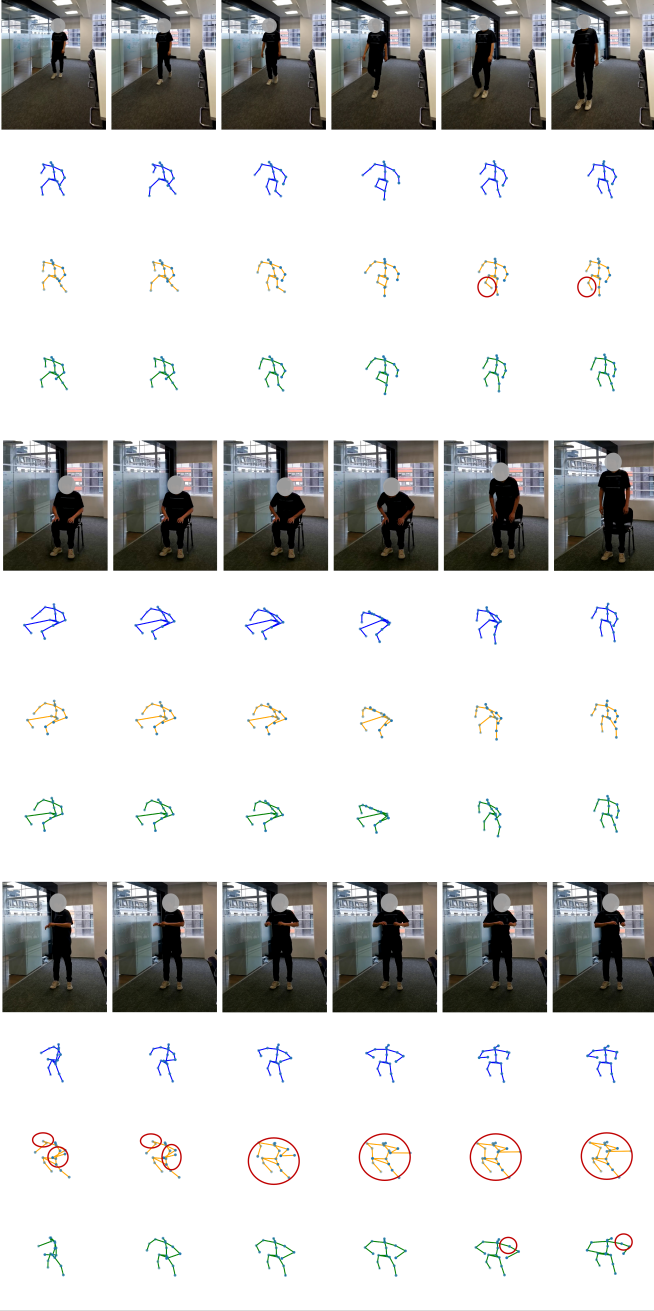


Fig. 5: The examples of pose estimation results of different activities: the first activity is W+, the second activity is SU and the third activity is BR; the blue skeleton is the ground-truth pose, the yellow one is M -based MDPose result and the green one is D -based MDPose result.

57.1mm. This result is close to the error of RFPose3D (CSI) at a training rate of 50%, and is still an acceptable value for most scenarios, which once again validates the good generalisation performance and robustness of MDPose.

Runtime: Table VI presents the quantitative analysis of the efficiency of MDPose. We used NVIDIA Quadro RTX 4000 graphics processing unit and 10 frames (i.e. around 1 second) μ -Doppler spectrogram to test the time taken by the different components of MDPose. As we can see, the denoising and

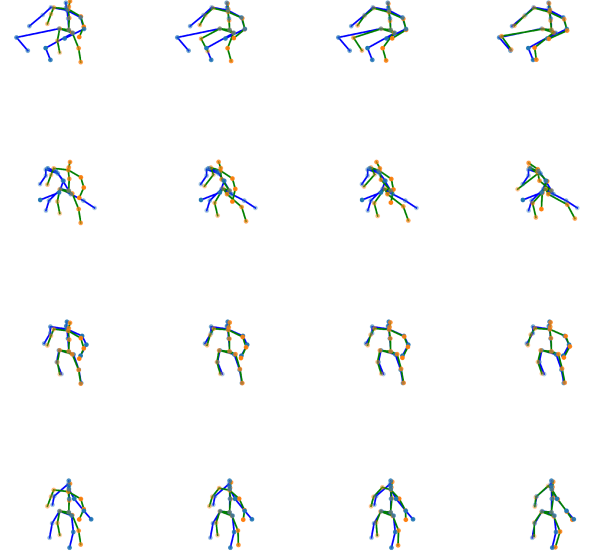


Fig. 6: The examples of pose optimising mechanism for initialising p_0

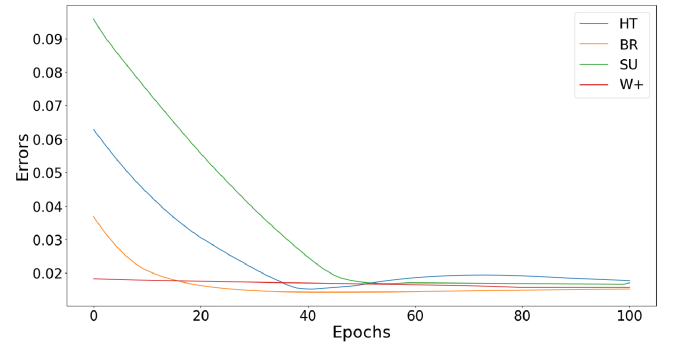


Fig. 7: The position error changes with the increase of optimisation epochs (y-axis unit: m)

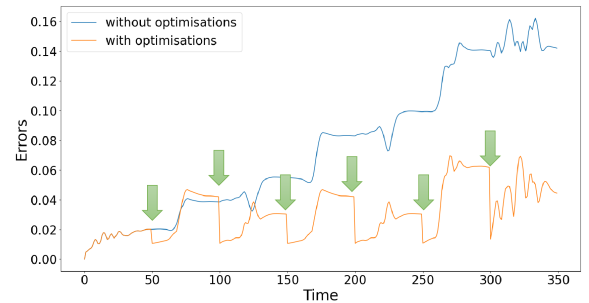


Fig. 8: The examples of pose optimising mechanism for optimising a long-term estimation: green arrows represent when we apply the optimising iteration

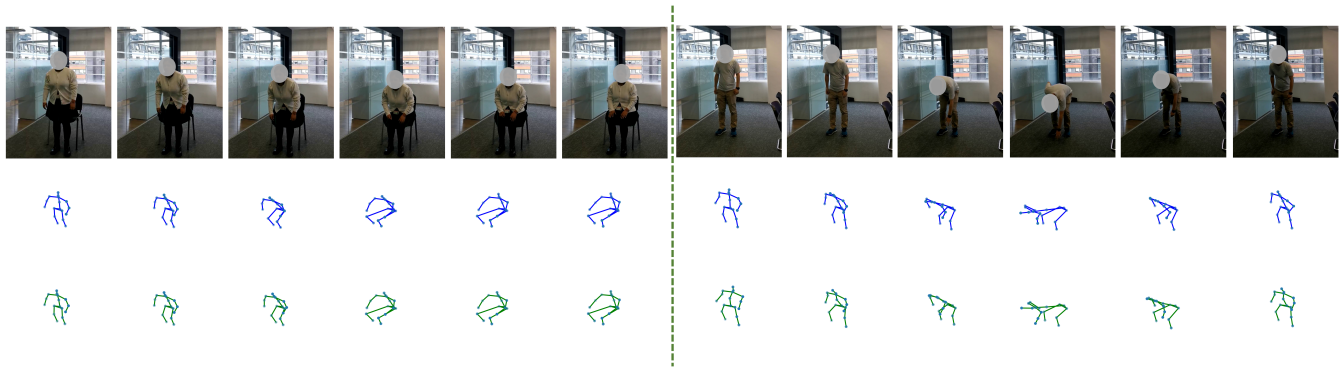


Fig. 9: The examples of pose estimation results of different subjects

Frameworks	Denoising	CNN-RNN	Optimisation	Total
RFPose3D	-	-	-	0.029
WiPose	-	-	-	0.016
MDPose	0.001	0.003	0.03	0.034

TABLE VI: Runtime analysis (unit: second)

velocity estimation only used around 0.004 seconds to obtain velocity properties. However, due to the lack of the initial pose, we have to iteratively use the optimising mechanism to estimate it, resulting in around 0.03 seconds run-time.. Therefore, compared with other methods, MDPose may spend more time to estimate poses. However, this result can still demonstrate that MDPose is an efficient method which has ability to efficiently process long-term μ -Doppler information.

VIII. CONCLUSION

In this paper, we presented a novel WiFi-based human skeletal motion reconstruction framework, MDPose, to effectively extract human motion from μ -Doppler information. It has two main phases to achieve the task: the CNN-RNN-based velocity estimation and initial pose estimation (or pose optimising mechanism). Additionally, we highly recommend using an appropriate denoising algorithm to remove interference and enhance Doppler features. Herein, we developed FMNet to clean up noisy spectrograms, and we believe that other effective denoising methods can also improve the performance of MDPose. From experimental results, we have demonstrated that each component of MDPose works well: the denoising network can provide clean features for the later steps; the velocity estimation network can effectively extract velocity properties of up to 17 body key points; the pose optimising mechanism not only helps to initialise the first pose, but also reduces the impact of accumulative errors on long-term estimations. Overall, MDPose has achieved the state-of-the-art performance with the estimation error of 29.4mm.

For future development of the research, there is still a lot of room for improvement with MDPose. Currently, we only use the single surveillance channel to collect Doppler information, which is not sensitive to velocity that varies along a path parallel to the receiver. Therefore, in our experiments, we only focused on processing activities perpendicular to the

receiver, but this has limitations for some realistic scenarios. To address this issue, our future plan is to use a multi-receiver system to more comprehensively capture the motion properties. Furthermore, we will also explore the multi-people detection in our future development.

ACKNOWLEDGMENTS

This work is part of the OPERA project funded by the UK Engineering and Physical Sciences Research Council (EPSRC), Grant No: EP/R018677/1.

REFERENCES

- [1] C. Zhang and Y. Tian, "Rgb-d camera-based daily living activity recognition," *Journal of computer vision and image processing*, vol. 2, no. 4, p. 12, 2012.
- [2] E. Ristani and C. Tomasi, "Features for multi-target multi-camera tracking and re-identification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 6036–6046.
- [3] Q. Gao, J. Wang, X. Ma, X. Feng, and H. Wang, "Csi-based device-free wireless localization and activity recognition using radio image features," *IEEE Transactions on Vehicular Technology*, vol. 66, no. 11, pp. 10 346–10 356, 2017.
- [4] W. Wang, A. X. Liu, M. Shahzad, K. Ling, and S. Lu, "Device-free human activity recognition using commercial wifi devices," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 5, pp. 1118–1131, 2017.
- [5] —, "Understanding and modeling of wifi signal based human activity recognition," in *Proceedings of the 21st annual international conference on mobile computing and networking*, 2015, pp. 65–76.
- [6] B. Tan, Q. Chen, K. Chetty, K. Woodbridge, W. Li, and R. Piechocki, "Exploiting wifi channel state information for residential healthcare informatics," *IEEE Communications Magazine*, vol. 56, no. 5, pp. 130–137, 2018.
- [7] S. Tewes and A. Sezgin, "Ws-wifi: Wired synchronization for csi extraction on cots-wifi-transceivers," *IEEE Internet of Things Journal*, vol. 8, no. 11, pp. 9099–9108, 2021.
- [8] C. Tang, W. Li, S. Vishwakarma, K. Chetty, S. Julier, and K. Woodbridge, "Occupancy detection and people counting using wifi passive radar," in *2020 IEEE Radar Conference (RadarConf20)*. IEEE, 2020, pp. 1–6.
- [9] W. Li, B. Tan, and R. Piechocki, "Passive radar for opportunistic monitoring in e-health applications," *IEEE journal of translational engineering in health and medicine*, vol. 6, pp. 1–10, 2018.
- [10] S. Z. Gurbuz, A. C. Gurbuz, C. Crawford, and D. Griffin, "Radar-based methods and apparatus for communication and interpretation of sign languages," Oct. 22 2020, uS Patent App. 16/850,664.
- [11] R. Palamà, F. Fioranelli, M. Ritchie, M. Inggs, S. Lewis, and H. Griffiths, "Measurements and discrimination of drones and birds with a multi-frequency multistatic radar system," *IET Radar, Sonar & Navigation*, 2021.

- [12] 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*, 2018, pp. 7356–7365.
- [13] M. Zhao, Y. Tian, H. Zhao, M. A. Alsheikh, T. Li, R. Hristov, Z. Kabelac, D. Katabi, and A. Torralba, "Rf-based 3d skeletons," in *Proceedings of the 2018 Conference of the ACM Special Interest Group on Data Communication*, 2018, pp. 267–281.
- [14] W. Jiang, H. Xue, C. Miao, S. Wang, S. Lin, C. Tian, S. Murali, H. Hu, Z. Sun, and L. Su, "Towards 3d human pose construction using wifi," in *Proceedings of the 26th Annual International Conference on Mobile Computing and Networking*, 2020, pp. 1–14.
- [15] C. Tang, W. Li, S. Vishwakarma, F. Shi, S. Julier, and K. Chetty, "Fmnet: Latent feature-wise mapping network for cleaning up noisy micro-doppler spectrogram," *IEEE Transactions on Geoscience and Remote Sensing*, 2021.
- [16] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [17] A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, and B. Frey, "Adversarial autoencoders," *arXiv preprint arXiv:1511.05644*, 2015.
- [18] S. Vishwakarma, W. Li, C. Tang, K. Woodbridge, R. S. Adve, and K. Chetty, "Simhumalator: An open source end-to-end radar simulator for human activity recognition," *IEEE Aerospace and Electronic Systems Magazine*, 2021.
- [19] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," 2017.