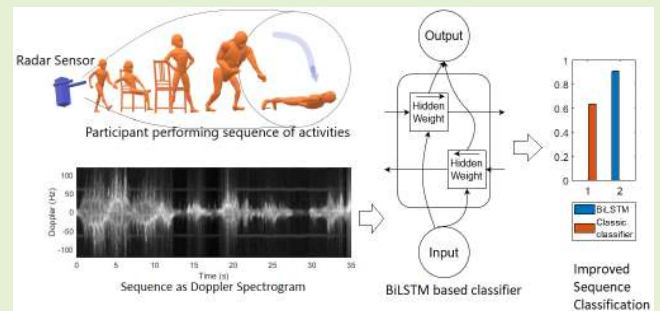


# Continuous Human Activity Classification From FMCW Radar With Bi-LSTM Networks

Aman Shrestha<sup>1</sup>, Student Member, IEEE, Haobo Li<sup>2</sup>, Student Member, IEEE, Julien Le Kernec<sup>1</sup>, Senior Member, IEEE, and Francesco Fioranelli<sup>3</sup>, Senior Member, IEEE

**Abstract**—Recognition of human movements with radar for ambient activity monitoring is a developed area of research that yet presents outstanding challenges to address. In real environments, activities and movements are performed with seamless motion, with continuous transitions between activities of different duration and a large range of dynamic motions, compared with discrete activities of fixed-time lengths which are typically analysed in the literature. This paper proposes a novel approach based on recurrent LSTM and Bi-LSTM network architectures for continuous activity monitoring and classification. This approach uses radar data in the form of a continuous temporal sequence of micro-Doppler or range-time information, differently from other conventional approaches based on convolutional networks that interpret the radar data as images. Experimental radar data involving 15 participants and different sequences of 6 actions are used to validate the proposed approach. It is demonstrated that using the Doppler-domain data together with the Bi-LSTM network and an optimal learning rate can achieve over 90% mean accuracy, whereas range-domain data only achieved approximately 76%. The details of the network architectures, insights in their behaviour as a function of key hyper-parameters such as the learning rate, and a discussion on their performance across are provided in the paper.

**Index Terms**—FMCW radar, micro-Doppler, remote activity monitoring, classification, LSTM and Bi-LSTM networks.



## I. INTRODUCTION

RADAR sensors in the context of short-range human monitoring are becoming increasingly popular, specifically in applications such as activities classification in smart homes within the ambient assisted living framework, recognition of gestures for human-computer interaction, and contactless vital sign monitoring [1], [2]. Broadly speaking, two categories of sensors can be used in all these applications, namely wearable and non-wearable sensors [3]. The former are usually attached to the body parts of the monitored subject with clasps or Velcro-straps, or are worn and carried in pockets. These sensors take fine resolution data from the specific movements of the human torso and limbs, characterized through their acceleration and angular velocity or displacements, or through

Manuscript received June 18, 2020; accepted June 26, 2020. Date of publication July 1, 2020; date of current version October 16, 2020. This work was supported by the U.K. EPSRC under Grant EP/R041679/1. The associate editor coordinating the review of this article and approving it for publication was Prof. Kazuaki Sawada. (Corresponding author: Julien Le Kernec.)

Aman Shrestha, Haobo Li, and Julien Le Kernec are with the James Watt School of Engineering, University of Glasgow, Glasgow G12 8QQ, U.K. (e-mail: julien.lekernec@glasgow.ac.uk).

Francesco Fioranelli is with the MS3 Section, Department of Microelectronics, TU Delft, 2628 Delft, The Netherlands (e-mail: f.fioranelli@tudelft.nl).

Digital Object Identifier 10.1109/JSEN.2020.3006386

direct measurement of the surface temperature, arterial movement for vital signs. The latter are often suggested as an alternative for wearable sensors, being less invasive in terms of interaction and management effort required by the end-users, often older people with forms of cognition impairments, and being independent from battery life duration [3]–[5]. Joint use of wearable and non-wearable sensors in multimodal frameworks have also been investigated, aiming to find the best combinations of relevant information from each sensor to be fused together, in order to achieve better monitoring performances [6]–[10].

Among non-wearable sensors, radar has attracted much attention recently as a possible alternative to video-cameras, thanks to its insensitivity to light conditions and easy integration into the end-users' home environment, as modern radar systems can be designed to look like a normal Wi-Fi router. Furthermore, radar may offer less privacy issues than cameras, as plain images or videos of the end-users and their private environments are not collected [11]–[15].

The work in [15] represented one of the first publications in the research field of radar-based human activities classification, where a set of specifically designed features extracted from micro-Doppler spectrogram was used in conjunction with a Support Vector Machine (SVM) classifier. More recently, the development of deep learning [16] and related classification

methods based on neural networks, has attracted significant interest also for their application to radar-based monitoring of human activities [2]. Their main advantage is the possibility to extract salient features automatically within the network, without explicit inputs or fine-tuning of parameters by the human operators that might miss important information and design a feature set prone to overfitting.

Numerous contributions in the literature used Deep Convolutional Neural Networks (DCNNs) to process the radar data as images. The work in [17] used DCNNs for classification of specific individuals and groups of individuals based on their walking gait. Comparison with conventional supervised-learning classifiers such as Naïve Bayes and SVM were provided, demonstrating better performances when using the DCNN. A DCNN was also used in [14] for human gait recognition, exploiting a dual-channel architecture where the network had two separate branches at the input, in order to accept spectrograms calculated with different temporal resolutions. A specifically designed DCNN was also used in [11] to identify specific individuals in different rooms based on their walking gait, with the additional complexity of the subjects following free-form, unconstrained trajectories.

In [18] DCNNs were used to classify human activities from their spectrograms, and in [19] a novel DCNN architecture was proposed to specifically account for the diversity induced by the different aspect angles on the radar signatures of human movements, especially with respect to their Doppler signature. Modifications to the conventional architectures of DCNNs were proposed in [2], [20], [21], in particular exploiting Convolutional Auto-Encoders (CAEs) to perform unsupervised pre-training of the weights of the network. CAEs and DCNNs for classification were also combined with a novel technique to augment the amount of available data in the training set by using Kinect-based motion caption simulations, enhanced by a diversification technique to improve the fidelity of the simulated synthetic data. Attempting to combine simultaneous classification of human activities and identification of specific individuals from their movements, [22] proposed a deep multi-task convolutional network validated using simulated human micro-Doppler data generated from the Carnegie Mellon MOCAP dataset.

Recent contributions in the literature have explored the usage of GANs, Generative Adversarial Networks [23]–[26], to address the need of a very large amount of data for training deep neural networks for classification, as it is a significant challenge to gather a lot of experimental radar data. GANs have been shown to be an effective tool to generate synthetic radar data starting from a relatively small set of experimental radar data, although there remain outstanding research challenges to evaluate the fidelity and reliability of such synthetic data, and their best usage to improve classification performances. The work in [23] used GANs to generate synthetic radar signatures for walking gaits at different speed, and [24] applied a similar approach to data of six human actions, including movements other than simply walking. Notably, the work in [25] proposed a novel approach to use the adversarial learning of GANs combined with a PCA-based (Principal Component Analysis) kinematic sifting approach

to reject the synthetic radar samples that present unrealistic data, i.e. data with artefacts that would not be realistically present in experimental data. Although not presenting any classification work on radar data, the investigation in [27] is of interest to show how micro-Doppler signatures of pedestrians (plus other automotive targets) appear when using special waveforms based on 512-bit Golay codes that enable joint radar-communication functionalities.

All the above papers that apply deep learning and deep neural networks to the classification of radar data have used convolutional neural networks in various architectures. All have in common the interpretation of the radar data as 2D images, i.e. matrices of pixels, typically Doppler-time patterns in the spectrograms. Even in cases when range-Doppler plots are used [28], this framework of processing the radar data as images remains in place.

Compared to the above state of the art methods, we investigate in this paper recurrent neural networks that interpret radar data as a temporal series and characterize the time-varying nature of a sequence of human activities and movements. In particular, we use Long Short Term Memory networks in their Bidirectional implementation (Bi-LSTM). Long Short Term Memory (LSTM) [29] is a recurrent neural network that can learn temporal dependencies between samples at separated time steps in a sequential data stream. LSTMs have been promoted as an ideal solution for temporally variant data for many applications, ranging from text and speech detection, audio processing, natural language processing and translation, up to finance and cell-biology [30]–[33].

However, LSTM and especially Bi-LSTM have been minimally discussed in the literature as a stand-alone tool for radar-based human activities classification, and represent an under-explored approach if compared with the DCNNs mentioned in previous paragraphs. In [34] an LSTM was used to classify the walking gait of small groups of people vs individual persons in an outdoor scenario. In [35], [36] recurrent networks have been used to classify six different human activities; specifically an LSTM was used in [36] and a stacked GRU (Gated Recurrent Unit) network, based on a simplified architecture of the LSTM cell, was used in [35]. It should be noted these data were collected as separated “snapshots”, i.e. separate radar recordings for each individual activity, thus missing to capture the natural transitions between each activity and the previous and following activity. This conventional snapshot data collection was also applied in [37] and [38], where LSTM networks were applied respectively to raw IQ radar data and to range profiles to classify separated human activities.

Summarising, to the best of our knowledge, so far very few works in the literature have investigated the use of LSTM networks, let alone Bi-LSTMs, for radar-based classification of human activities; when these have been used, the data referring to the classes of interest were collected as separated radar recordings, thus without capturing the realistic transitions between human movements. In this paper on the contrary we analyse continuous sequences of human activities involving a relatively large group of subjects, and exploring different combinations of the activities and therefore

inter-activity transitions. The main contributions are as follows:

- We analyse realistic, continuous sequences of human activities. Within them, natural transitions between the different actions can happen at any time, with unconstrained duration for each activity and for the transitional period in which the body parts reposition themselves appropriately in order to perform the following action. This represents a novel element in data acquisition towards an enhanced realism of the captured scene [2];
- We propose stacked Bidirectional LSTM networks as a novel deep learning tool alternative to DCNNs, to perform radar-based classification of these continuous sequences of human activities. Bi-LSTM are inherently suitable for such analysis, because they can capture both temporal forward and backward correlated information within the radar data, specifically the kinematic constraints and characteristics that relate each recorded activity to the previous and the following actions. Insights on the effects on the performance related to choices in data pre-processing and key hyperparameters (e.g. learning rate) are provided.
- We base our analysis on experimental data collected using a C-band radar and involving 15 participants performing different combinations of 6 activities. This enables to validate the proposed approach on a relatively large set of participants and with sequences presenting different transitions from an activity to another.

The remainder of this paper is organized as follows. Section II describes the experimental setup with the radar, data collection, and overall methodology. Section III presents a description of the results obtained with LSTM and Bi-LSTM networks used for different data domains (range/Doppler) and offers some insight on optimizing performances. Finally, section IV concludes the paper and outlines possible future work.

## II. METHODOLOGY

In this section the overall methodology of the proposed approach is presented. This includes the discussion on the motivations for using Bi-LSTM networks, and the description of the experimental setup, data collection, and data pre-processing.

### A. Motivation for Continuous Activities and Bi-LSTM Networks

As discussed in the introduction, the research focus on human activity detection with radar has been on discrete separated activities, which are typically performed and recorded one at a time. For the analysis of continuous activities, discrete data samples can be sequentially concatenated as in [36], [39], but this does not capture the full realism of unconstrained human movements, where the duration of each action can change, and the inter-activity transitions can happen at any time.

To evaluate this more realistic scenario, the data set analyzed in this work includes continuous activities performed in a natural manner by the participants. This also captures the diversity in sequential order and transitions between the different

activities. Examples of radar spectrograms of these continuous activities are shown in the following sections. Continuous recordings of radar data can appear in time as sequences of range profiles, often stacked together to form range-time matrices, or micro-Doppler spectrograms [40]. The majority of the works in the literature would interpret these radar data as 2D images or 3D “cubes of voxels” and process them with methods inspired by the image processing community, such as convolutional neural networks or auto-encoders. In this framework, a sliding window of fixed length could also be applied across the sequence of radar data to extract images of individuals or sub-sets of activities. However, in a realistic sequence of human movements, there is no fixed duration of each individual action and the transitions between actions can happen at random times. Therefore, rather than images, these continuous radar data appear more similar to sequences of speech or audio signals where individual words or patterns can appear at any time and with unconstrained duration.

For this reason, the recurrent neural network architectures inspired by the work in the audio/speech processing community are explored in this paper. Specifically, we focus on Bidirectional LSTM.

The main property of the LSTM is the memory capability to capture the long-term dependency between data separated by a significant number of time steps [29]. This is relevant in speech, where two strongly correlated words can be separated by other words (e.g. auxiliary verb and past participles in Germanic languages, nouns and adjectives where many adjectives are used). Radar data can resemble speech, as different actions performed at different time steps are correlated by human kinematics (e.g. one can stand up only after sitting down, but a variable amount of time can separate these two actions). However, speech or audio data do not encode any kinematic information or constraint, that are instead the main feature of the radar data and what radar-based classification algorithms aim to understand.

Then, bidirectionality is the capability of correlating the data processed at a given timestep with both data from past and future timesteps [41]. This is again an essential property in speech/language processing to capture the relation between different words in a long sentence, but also a relevant capability in radar-based activity classification to capture the kinematic constraints of human movements (e.g. an action performed at a given timestep is related to previous actions and can constrain future actions).

### B. Experimental Setup and Data Collection

The data from 15 participants (14 male and 1 female) aged 21-35 years were collected at the University of Glasgow in July 2018. The participants recreated daily life activities/movements in this room, where an activity zone was set for them to perform their movements. This area, along with the radar setup, is shown in Fig. 1. One male participant provided data twice, and these data from the repeated recordings have been used as the validation set for the networks, as having a validation set improves generalization of the classifier model [42]. The experimental setup with the surrounding clutter and



Fig. 1. Environment of the experimental setup with radar antennas mounted on tripods and sources of static clutter (furniture) shown.



Fig. 2. A pictorial list of activities; these six activities were performed in a different order in three different continuous sequences.

furniture is shown in Fig. 1. In this data collection campaign, in addition to the FMCW radar system, three wearables were used to record the activity data for multimodal fusion purposes, but the analysis of their data is beyond the scope of this paper. The FMCW radar was operating at 5.8GHz with 400MHz instantaneous bandwidth and 1ms chirp duration.

The data include six human activities: walking (A1), sitting on a chair (A2), standing up (A3), bending to pick up an object (A4), drinking a glass of water (A5) and simulating a frontal fall (A6). These activities are shown in Fig. 2. This shows the individual activities recorded as discrete actions, but they were performed as continuous sequences, i.e., with each action performed one after the other with varying duration and unconstrained transitions between them. The total duration of each sequence was 35 seconds, and three different sequence orders were recorded for each participant, namely:

- A1: A2: A3: A4: A5: A6
- A5: A4: A2: A3: A1: A6
- A4: A5: A1: A2: A3: A6

For the 15 participants, the three different sequences of continuous activities provided 45 different recordings for the main data-set, plus 3 recordings for the validation set with a repeated participant.

### C. Training and Testing Set Composition

A total of 48 different sequences were collected, where 3 sequences were repeated recordings performed by one subject, and set aside as the validation test. Out of the remaining 45 sequences, the testing test always included one of the 45 sequences, repeating the process 45 times to test all sequences. For the training set, two different approaches were followed to investigate the effect on the classification performance of prior knowledge/data about a specific participant. In other words,

In the approach labeled as “New,” the two sequences belonging to the subject under test were removed from the

training set, leaving a total of 43 sequences for training. In this case, the test subject is unknown to the classifier, as if a new person joining the experiment. In the approach labeled as “Known Prior,” two random sequences out of 45 were removed from the training set, leaving 43 sequences for training for consistency with the previous case. However, the two sequences performed by the test subject were purposely kept in the training set. In this case, the classifier did have some knowledge of the test subject through the two sequences in the training set, although the order of the activities and the related transitions were different. These two approaches were tested to evaluate any difference that prior knowledge of an individual human subject would provide to the classification algorithm and its performance.

### D. Radar Data Processing and Representations

For radar data to be used as inputs to the classifier, firstly, a Fourier transform is performed on the matrix of raw radar returns to generate range-time profiles (RP). To remove static clutter, a moving target indicator (notch filter) is then applied, and then using specific range bins where the target is performing the activities, a Short-Time Fourier Transform (STFT) is applied to find the Doppler-time pattern to characterize the micro-Doppler signatures [40]. In this experiment, a 0.2s Hamming window and an overlap factor of 95% is used to generate the micro-Doppler spectrograms. The LSTM and Bi-LSTM neural networks can take either range profiles or spectrograms as inputs and the effect of using them will be evaluated in the sections below.

### E. Machine Learning Library

While other works have used various libraries [15], [36], this work utilized the Deep learning toolbox included in MATLAB 2018B. Additionally, MATLAB was also used for radar signal processing and for manually labeling the ground truth data from physical observation during the experiments.

## III. EXPERIMENTAL RESULTS AND PERFORMANCE ANALYSIS

In this section experimental results using different LSTM network architectures are presented, with discussions on changes in performances due to the format of input data used (e.g. spectrograms vs range-time plots), and on significant hyperparameters of the networks (e.g. learning rate).

In this section, lower case symbols will denote vectors, e.g.  $x$ , whereas matrices are denoted by upper case letters  $H$ . An arrow pointing right, e.g.  $\vec{H}_t$  indicates the scalar or vector in the next time step whereas an arrow pointing left, e.g.  $\overleftarrow{H}_t$  indicates the scalar or vector from the previous time step.  $\odot$  denotes the Hadamard product, an element-wise product of two vectors.

### A. Doppler LSTM

The first network investigated is a two-stage stacked LSTM network, which is referred to as *Doppler-LSTM* and serves as a baseline for the spectrogram-based results. As the name

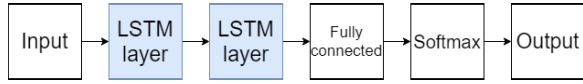


Fig. 3. The network architecture of *Doppler-LSTM* comprising: • An input layer that takes a segment of the spectrogram (250 Doppler bins in each time bin, which is equivalent to one observation) and sends it to the first hidden layer, • Two stacked LSTM layers that extract and update the salient features in the input data, • A fully connected layer that connects the activations of the different LSTM layers necessary for classification, • A softmax that computes the probability distribution of the data belonging to a specified output class, • An output layer that outputs the class label based on the Softmax distribution. Note that the arrows indicate the temporal direction of the recurrent LSTMs. In this case, since a standard LSTM layer is used, only forward based recurrence is considered.

suggests, the input to this network is the spectrogram, which contains micro-Doppler information and is fed into the network as a sequence of different vectors time bin after time bin. This network only implements a forward based dependency in analyzing the data from the sequential timesteps.

Fig. 3 shows a simplified block diagram of the proposed architecture of the *Doppler-LSTM* network. The inner workings of the gates in an individual LSTM layer are given by equations 1-4 [29].

$$f_t = \sigma_g(W_f x_t + R_f h_{t-1} + b_f) \quad (1)$$

$$i_t = \sigma_g(W_i x_t + R_i h_{t-1} + b_i) \quad (2)$$

$$g_t = \sigma_c(W_g x_t + R_g h_{t-1} + b_g) \quad (3)$$

$$o_t = \sigma_g(W_o x_t + R_o h_{t-1} + b_o) \quad (4)$$

Equation 1 shows the operation of the forget gate and is based on an activation function applied to the sum of the weighted input (with weight  $W$  and input  $x$ ) with the product of recurrent weights  $R$  and the hidden states  $h$  from the previous time instance, plus a bias term  $b$ . The other gates perform similarly, with differences arising from the input and recurrent weights, as well as the bias being unique to each gate.

The sigmoid activation function, cell state output, and hidden states output are represented by equations 5, 6 and 7.

$$\sigma(x) = (1 + e^{-x})^{-1} \quad (5)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t \quad (6)$$

$$h_t = o_t \odot \sigma_c(c_t) \quad (7)$$

To better understand the use of LSTM based networks for time-dependent data, the LSTM cell behavior can be empirically described. Fig. 4 presents a sketch of the LSTM cell showing the two outputs at the current timestep  $t$ , namely the hidden state  $h_t$  and the cell state  $C_t$ .  $h_{t-1}$  is the hidden state at the previous timestep and  $C_{t-1}$  the cell state at the previous timestep. These two signals, together with the input data at the current timestep  $X_t$ , are the input signals to the LSTM cell. This implies that the outputs at the current timestep depend on the hidden state and cell state from previous timesteps, therefore, utilizing the memory of the network, as also described in the previous set of equations.

Four components control the two outputs:

- $f$  is the forget gate which resets the state of the cell making it forget prior information from the previous cell state;

TABLE I

SIZE AND PROPERTY OF LAYERS USED IN *Doppler-LSTM* NETWORK

| Layer  | Size | Properties   |
|--------|------|--|
| Input  | 240  | based on the frequency bins of the input spectrogram |
| LSTM   | 2400 | number selected to store large sequences in memory   |
| LSTM   | 2400 | number selected to store large sequences in memory   |
| FC     | 6    | Based on the number of possible output classes       |
| Output | 1    | Single output  |

- $g$  is the cell candidate which provides input to the cell state keeping memorable or recurrent information and providing it to the cell state;
- $i$  is the input gate which co-ordinates with  $g$ ;
- $o$  is the output gate to control the addition of the cell state to the hidden state.

The original recurrent neural network architectures, before the development of LSTM, did not have states. Therefore temporally pertinent information across many timesteps was not retained; the cell state changed this, as longer time-based dependencies could now be memorized.

In terms of radar data, this means that the information on human movements can be memorized and correlated over a relatively long time. In the Doppler-time representation (spectrograms), an activity in a sequence of movements is perceived by the radar as a specific pattern of active Doppler bins over time. The network can learn this pattern in its internal parameters to recognize this activity even when it has different lengths of ‘activation’ or delays.

With the temporal dependencies accounted for, the level of abstraction in the input data should be assessed, as spectrograms can be considered a mixture of multi-tones where the micro-Doppler movements induce different Doppler frequency components depending on the movements of individual body parts. Using multiple layers has been suggested as the primary method of detecting higher-level abstractions from the input domain [36]. Therefore as shown in fig 3, the proposed Doppler-LSTM network has two stacked LSTM layers so that these higher-level abstractions can be identified by the network.

The neural networks have 2410 hidden cells for both of the LSTM layers and a learning rate of  $2^{-4}$ . Of the editable hyperparameters, the learning rate is of significance as it is the key contributor to vanishing and exploding gradient problems [16]. In section III-G, we show that this hyperparameter affects radar data significantly and can offset good architectural decisions if an incorrect learning rate is used. The state activation function is the hyperbolic tangent, while the gate activation function is a sigmoid (rectified linear unit is not commonly used with LSTM as it can cause exploding or vanishing gradients.) and gradient descent optimizer is “Adam.” Table I shows a summary of the size and properties of the layers of the *Doppler-LSTM*. This network was then trained and tested with the procedure described in Section II-C.

## B. Doppler Bi-LSTM

The second proposed network referred to from hereon as *Doppler Bi-LSTM*, is a modification of the first one and includes: an input layer, an LSTM layer, a BiLSTM layer, a Softmax layer, and a classification layer. The Bidirectional

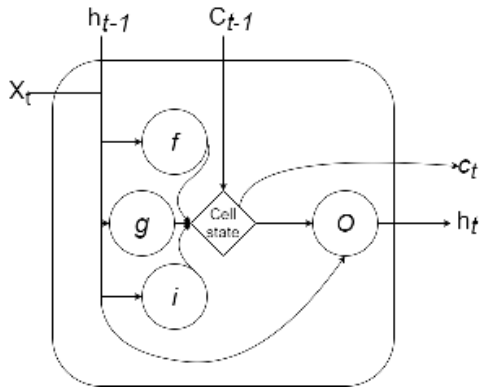


Fig. 4. Overview of the LSTM cell used in *Doppler-LSTM*. It is composed of: •  $f$  forget gate: Control to forget cell state •  $i$  input gate: Control to update cell state •  $g$  cell candidate: Control for information to be added to cell state •  $o$  output gate: Control to add cell state to hidden state. Note that the arrows indicate only forward-based temporal information flow from a timestep  $t-1$  to the following  $t$ .

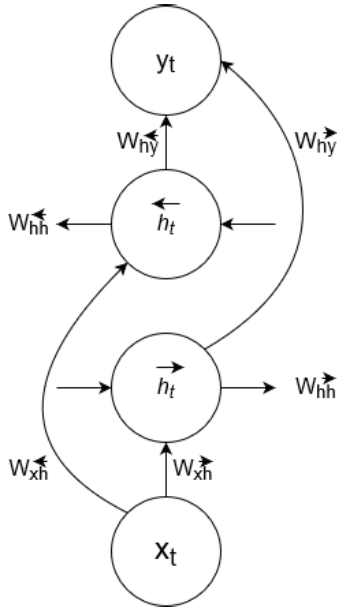


Fig. 5. Interconnections and weight transfers in a Bi-LSTM cell used in the *Doppler Bi-LSTM*. The arrows show the propagation of the information hidden and cells states between the layers.  $X_t$  is the input,  $h_t$  is the hidden state with its forward or backward directionality,  $W_{nn}$  indicates the weights linking hidden states and outputs/inputs and  $Y_t$  is the output.

LSTM cell is the main modification of this network and its details are shown in Fig. 5, whereas Fig. 6 shows the block diagram of all the layers of the proposed network.

Similar to the first network, *Doppler Bi-LSTM* accepts spectrograms as inputs. Differing from the previous Doppler-LSTM, this network processes the forward time-based dependencies first in the initial LSTM layer, and then searches for bidirectional, forward and backward, dependencies in the extracted temporal features. The capability of characterizing and memorizing these forward and backward dependencies in the sequences of data is critical for this network and its performance, as in the sequence of human activities, there are explicit dependencies and kinematic constraints on the order of possible actions.

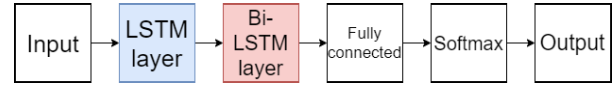


Fig. 6. Network architecture of *Doppler Bi-LSTM*. The key difference between the previous architecture is the presence of a bidirectional layer.

TABLE II

SIZE AND PROPERTY OF LAYERS USED IN *Doppler Bi-LSTM* NETWORK. OPTIMAL LEARNING RATE OF  $1E-4$  IS USED

| Layer   | Size | Properties   |
|---------|------|--|
| Input   | 240  | based on the frequency bins of the input spectrogram |
| LSTM    | 2400 | number selected to store large sequences in memory   |
| Bi-LSTM | 2400 | number selected to store large sequences in memory   |
| FC      | 6    | Based on the number of possible output classes       |
| Output  | 1    | Single output  |

The main equations for a Bi-LSTM cell unit [32] are as follows:

$$\vec{h}_t = \tanh(W_{x\vec{h}} X_t + W_{h\vec{h}} \vec{h}_{t+1} + b_{\vec{h}}) \quad (8)$$

$$\overleftarrow{h}_t = \tanh(W_{x\overleftarrow{h}} X_t + W_{h\overleftarrow{h}} \overleftarrow{h}_{t+1} + b_{\overleftarrow{h}}) \quad (9)$$

$$y_t = W_{\vec{H}y} \vec{h}_t + W_{\overleftarrow{H}y} \overleftarrow{h}_t + b_y \quad (10)$$

The main difference in the Bi-LSTM layer versus the previously described LSTM layer comes from each cell having two hidden states, with two parallel pipelines feeding to both previous and next timesteps as illustrated in Fig. 5. Note that in this figure, the two different hidden states are denoted with capital  $H$  and forward and backward arrows, respectively, as they are in the equations. Differently from the LSTM, in the Bi-LSTM layer, the interconnections between the input, output, and hidden states through the relevant weights do not propagate through the forward and backward cells directly; instead, they interface separately by going through the forward cells ( $\rightarrow$ ) and backward cells ( $\leftarrow$ ) at the same timestep. The hidden states from these forward and backward cells are then combined to generate the output from the Bi-LSTM layer, denoted by  $y_t$ . The implication is that the Doppler information corresponding to specific body movements over a long duration in both forward-time and backward-time directions are characterized and captured by the Bi-LSTM layer. Essentially, this means that the network searches and memorizes recurring feature patterns in the past (previous actions) and any linked recurring feature patterns in the future (subsequent actions).

### C. Doppler LSTM and Doppler Bi-LSTM Performance Analysis

Fig. 7 shows the spectrogram of one of the sequences classified by the Doppler-LSTM. Furthermore, it shows the comparison of the classification and ground truth of the activities within this sequence. Initially at  $t = 0$ , we see that there is a sharp spike that detects A5: Drink while in truth the person was performing A4: Pick, since both of these activities have the central component of moving arms the classifier has a moment of indecision. It then correctly classifies A4: Pick, but it detects A5: Drink with a delay of 3 seconds, after which another “impulse-like” indecision, referring to the sharp spike at about 9 seconds, where A6: Fall is detected. In a fall detection system, the presence of these spikes for erroneous

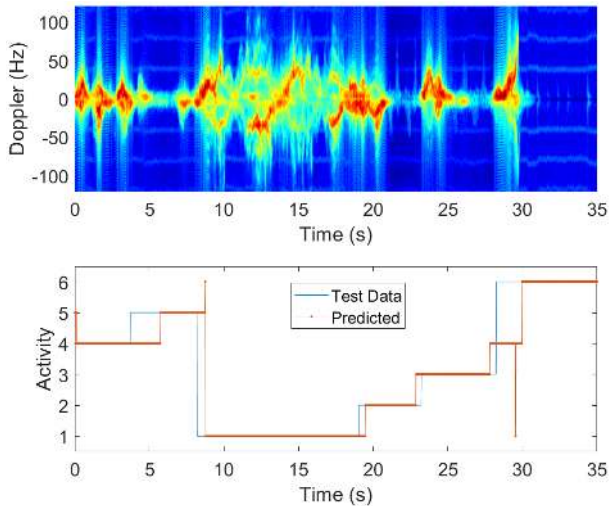


Fig. 7. Classifier input at the top sub-figure, ground truth in blue, and the test outcome in orange in the bottom sub-figure for a test sequence for the *Doppler-LSTM* network.

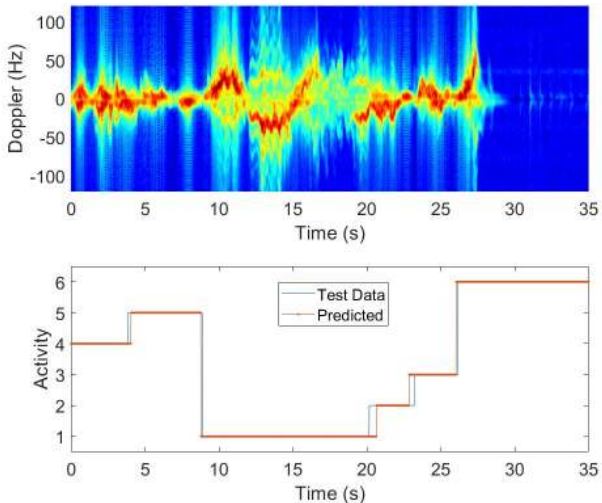


Fig. 8. Classifier input at the top sub-figure, ground truth in blue, and the test outcome in orange in the bottom sub-figure for a test sequence for the *Doppler-Bi-LSTM* network.

classifications could be undesirable and a potential source of false alarms.

Fig. 8 represents an example of results for the Doppler Bi-LSTM network. Note that the sequence of activities is the same as the one presented in the previous figure, but performed by a different subject. For this reason, the ground truth plots are identical, but the input spectrograms appear overall similar. Comparing the classifier output/test outcome in the orange line and the ground truth in the blue line, we can see that test outcome matches the classifier output to a very large extent. However, there are three noticeable segments at time points 4, 20, and 22 seconds where there is a slight mismatch between the test outcome and the observed ground truth. In the first case, at 4 seconds, there is a short delay in detecting the transition from A4: Pick up to A5: Drink. However, one can note that in the spectrogram input in Fig. 8, the signature is unclear at that time instance, with a difficult transition detectable by eye. This is typical of transitions where

TABLE III  
ACCURACY METRICS FROM THE TESTED PARTICIPANTS ACROSS  
DIFFERENT LSTM ARCHITECTURES WITH DOPPLER  
AND RANGE INPUT

| Classifier   | Mean | Standard deviation | Maximum | Minimum |
|--------------|------|--------------------|---------|---------|
| Bi-LSTM      | 91   | 5                  | 98      | 69      |
| LSTM         | 78   | 9                  | 92      | 54      |
| Range BiLSTM | 76   | 7                  | 87      | 54      |

the dynamic range of the macro movement of the body/torso and the micro-movements of the limbs change drastically. The network may respond to this by maintaining the classification from the previous time instances, so there is a short delay but no erroneous classification occurs. This is similar to the second case, at 20 seconds, where the classifier appears to detect A6: Fall with a short delay. In the third instance, at 22 seconds just before it happens. Reviewing the spectrogram, prior to the A3: Standing occurring, there is a precursory movement which the classifier notes and associates as part of the A6: Fall class, possibly due to the knowledge of the signature at future time instances in the spectrogram provided by the bidirectional capabilities of the network.

Fig. 9 shows the classification accuracy for the 45 sequences collected where each one was the test sequence in turn, as discussed in section II.C. Note that the hyperparameters and training/testing approach were kept consistent between the two network architectures. This provides a more effective way to compare the performances of the proposed Doppler LSTM and Bi-LSTM networks across the whole dataset of continuous signatures, rather than observing individual sequences.

The range of classification accuracy is, on average higher for the Bidirectional LSTM network compared with the unidirectional LSTM, and there is less variability across different subjects and different sequences performed by the same subject. This can be described by the mean and standard deviation across the 45 classification tests, that are recorded in Table III for the Doppler Bi-LSTM and LSTM networks. The mean increases to approximately 91% from 78%, whereas the standard deviation is reduced; the maximum (best case) and minimum (worst case) are also increased when using a bidirectional architecture of approximately +6% and +15%, respectively.

The metrics in Table III and the detailed results in Fig. 9 for each test sequence show how the bidirectional capability of the proposed Bi-LSTM provides superior capabilities to classify human activities in a continuous sequence with respect to a conventional unidirectional LSTM (for example the increase in accuracy from LSTM to Bi-LSTM is +30% in the best case of sequence #31). The robustness and good generalization of the proposed approach across the diverse set of 45 recorded sequences and 15 subjects are demonstrated. Furthermore, for subsequent sequences where there is a drop in the accuracy, the Bi-LSTM appears to be more robust than the unidirectional LSTM, for example for sequences 5 and 6, where for the Bi-LSTM the accuracy drops from 94% to 88% and for the LSTM it drops from 78% to 57%. As a note, the accuracy in Fig. 9 and Table III is calculated as the number of correct classification of the activity (1 out of the 6 performed) in each time bin of the spectrogram, over the total number of time bins

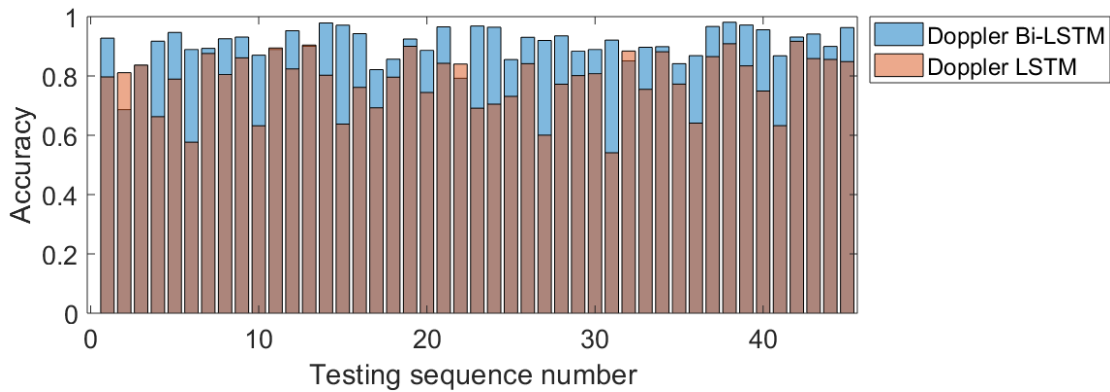


Fig. 9. Comparison between *Doppler-LSTM* and *Doppler Bi-LSTM* architectures output as classification accuracy over the 45 test sequences. Although the layers between these classifiers are different, the hyperparameters and training and testing methodologies are consistent between both network architectures.

in the 35s total duration of each testing sequence. We think that this is a conservative approach, which labels as mistakes even the very short spikes with erroneous prediction lasting for only a few time bins, as seen in Fig.7. As part of future work, smoothing filters approaches could be applied to the predictions of the LSTM or Bi-LSTM networks to disregard labels for activities that would last only for few time bins, and therefore be unrealistic as the subject could not perform a given activity in such short physical time.

#### D. Range Bi-LSTM

In the previous sections, we have analyzed the results of using Doppler spectrograms input to the LSTM and Bi-LSTM networks, but spectrograms need an additional level of processing after the generation of the range profiles to be calculated. This prompts the question of whether sufficient information can be inferred from the data in the range-time domain, leaving to the networks the task of extracting the Doppler information, i.e. the changes between subsequent range profiles implicitly.

Range profiles do not show the different activities in the signature in an easily perceivable manner compared to spectrograms since only the location relative to the radar is given, and in the specific case of our radar, the range resolution is limited to approximately 40cm with the 400 MHz bandwidth. Hence, to the human eye, the different activities in the range-time plots appear much less distinguishable than in the spectrograms, but this may not necessarily be a limitation for neural networks.

Fig. 10 shows an example of such a range-time plot in its top part; it is evident how this image is less clear than the corresponding spectrograms in Fig. 8 and Fig. 7 for the same sequence of actions. A significant difference between spectrograms and range-time plots, is the number of time bins, of temporal units that the LSTM or Bi-LSTM network will need to process. For 35s of data in each sequence of activities, spectrograms consisted of 1750 time bins with the selected STFT parameters, whereas range-time plots had 35,000 observations or range profiles, as the data were sampled at 1 kHz PRF. This increased size of the data led to a modification of the network with a different number of inputs reflecting the number of range bins in the range profile. This network is referred to as *Range Bi-LSTM*. In terms of its layers

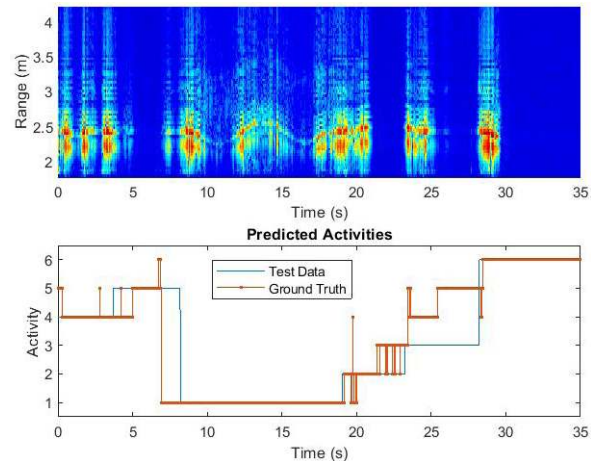


Fig. 10. Classifier input, ground truth, and the test outcome for a test sequence for the *Range Bi-LSTM* network.

and architecture, the *Range Bi-LSTM* is similar to the one shown in Fig. 6.

The bottom part of Fig.10 shows an example of representative results from this alternative network using range-time data as input. The performance is reduced compared to the Doppler based networks.

At 0, 4, and 5 seconds, we see transient detection of A5: Drink at multiple instances until an apparent misdetection of activity A1: Walking as A2: Sitting as the target comes to a halt which is visible in the range-time plot. This is followed by an early detection of the A2: Sitting. At 28 seconds, multiple instances of A5: Drinking is detected before A4: picking up item is correctly identified, which is reminiscent of the spike transients observed with the Doppler-LSTM. In general, more spikes and instability in providing a steady prediction are shown at other transitions, and there are misdetections of all activities throughout the sequence.

#### E. Range-Time Bi-LSTM Performance Analysis

The results in Fig.10 for the usage of range-time data as inputs to the Bi-LSTM show a degradation in performance compared to the usage of micro-Doppler information. To view the performance of the network on a sequence by sequence basis across the whole dataset, Fig. 11 shows the results



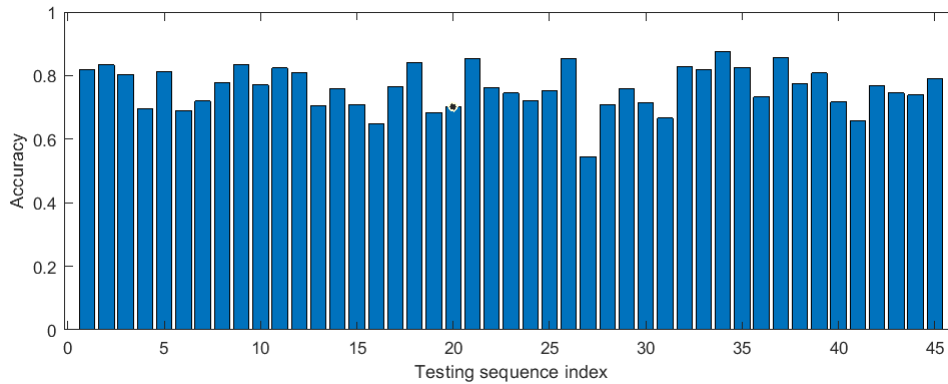


Fig. 11. Range Bi-LSTM output as classification accuracy over the 45 test sequences.

for the 45 test sequences, and Table III shows the overall performance metrics.

There are cases where the Range BiLSTM performs well. For example, in the sequences 31-35, an accuracy of approximately 80% is attained. However, it does not maintain this rate for all of the test sequences as the classification challenge of detecting complex activities designed for this set, and also delayed and transient detection of classes occur as demonstrated in Fig. 10. Viewing Table III while comparing Fig.11 and Fig. 9, show the performance loss of using the range-time profiles as inputs to the proposed LSTM networks, despite the potential advantage of avoiding the calculation of spectrograms at the pre-processing stage before the network. To put it into perspective, the best classification accuracy, or the maximum in Table III (87%), for any range input is 4% less than the mean accuracy for the Doppler Bi-LSTM (91%). In other words, the best case with range input cannot match the average case with Doppler input with a similar or even a range focused network architecture. Directly comparing the mean accuracy shows an improvement of 15% through the use of the Doppler Bi-LSTM (91%) instead of Range Bi-LSTM (76%) and improvements in the maximum by 5% and minimum rate by 15% when the former architecture and its corresponding input is used.

#### IV. FURTHER EXPERIMENTAL VALIDATION

In this section we present further tests to validate the proposed methods. We discuss the influence of the classifier having prior information from participants, compared to when no such information is provided. Additionally, the effect of static clutter on the classification is analysed and a comparison with a simpler support vector machine classifier (SVM) is made, with a comparative analysis is provided.

##### A. Known Prior vs Unknown

Table V shows the results from the networks and input domains discussed in this paper with the “Known Prior” and “New” training and testing methodologies. For each of the participants, one sequence of activities was taken as a test sample with the best performing classifier and input domain combined. In the case of *Known Prior*, the training set had

TABLE IV

SIZE AND PROPERTY OF LAYERS USED IN *Range Bi-LSTM* NETWORK. LAYERS WERE RESIZED TO FIT INPUT DOMAIN AND MEMORY LIMIT

| Layer   | Size | Properties   |
|---------|------|--|
| Input   | 64   | based on the range bins of the input spectrogram   |
| LSTM    | 240  | number selected to store large sequences in memory |
| Bi-LSTM | 240  | number selected to store large sequences in memory |
| FC      | 6    | Based on the number of possible output classes     |
| Output  | 1    | Single output                                      |

TABLE V

ACCURACY METRICS FOR RANGE VS DOPPLER DOMAIN NETWORKS WITH “KNOWN PRIOR” AND “NEW” TRAINING AND TESTING APPROACH

| Classifier      | Subject     | Mean | Standard deviation |
|-----------------|-------------|------|--------------------|
| Range Bi-LSTM   | Known Prior | 74   | 7                  |
| Range Bi-LSTM   | New         | 77   | 7                  |
| Doppler LSTM    | Known prior | 78   | 10                 |
| Doppler LSTM    | New         | 78   | 9                  |
| Doppler Bi-LSTM | Known prior | 91   | 4                  |
| Doppler Bi-LSTM | New         | 90   | 6                  |

included the other two sequences performed by the same subject (but with a different order of activities), whereas for the *New* set, the classifier did not have information on that specific subject from the other sequences. The Table shows there appears to be a marginal difference in the prediction between a classifier that has prior information from the test subject and one which does not. The only factor inducing a significant change is the selection of the network and corresponding input domain, where it follows the trends discussed at the end of section III-E. Substantially, the Doppler Bi-LSTM outperforms the other architectures/input domains. As the prior knowledge of the test subject induces no significant change in the classification accuracy, we assume that during the data collection the same activities were not reproduced in the exact form in all sequences, despite the test subject remained the same. This was due to the fact that the duration of each activity was unconstrained and that as the order was different, the transitions between the activities happening before or after a given one created diversity in the data. Therefore, each sequence is distinct, and there is no much difference in providing to the network knowledge of other sequences at the training stage for a given test subject. Conversely, it can be seen that, when the best performing combination of network and radar data format is used, the proposed approach based on recurrent LSTM networks is robust enough to generalize across the cohort of 15 subjects and 45 sequences. The analysis

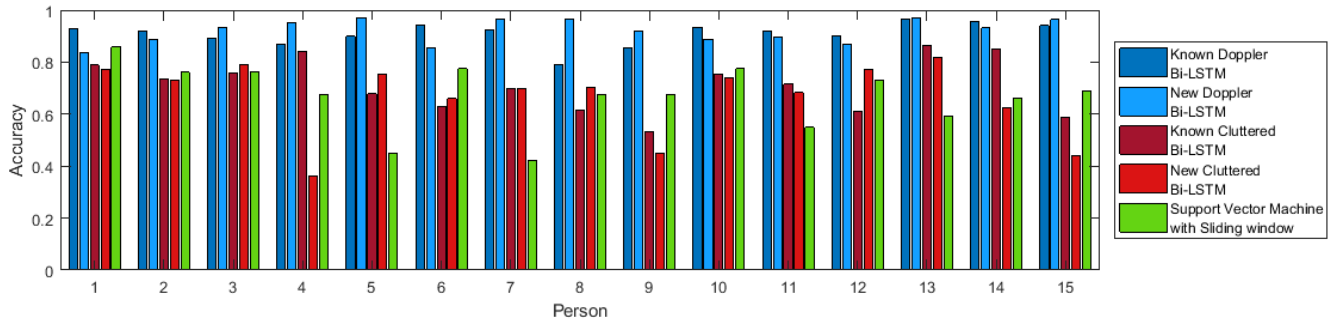


Fig. 12. Comparison between testing methodologies: where the network had samples of the test subject but with a separate sequence called Known prior; and another where the network had no prior information from the test subject called new. Repeated for scenario with and without clutter and SVM results presented as a classical classifier for comparison. Results generated with the *Doppler Bi-LSTM* network.

of some of the individual test sequences for the “Known Prior” and “New” approaches in Fig. 12 shows cases where prior knowledge appears to help with improving accuracy. For example, with person 1 and 6, there is an improvement of more than 10%. However, there are also cases where the opposite occurs, for example, with person 8. Therefore, prior knowledge of the test subject with the classifier does not appear to be a major factor in the overall accuracy.

### B. Influence of Static Clutter

The role of clutter is another aspect which is questioned about in the research area of using radar for monitoring human activities as indoor environments consistently have objects which generate static clutter and possible multipath. This is usually mitigated by using MTI filters, as it has been done in this paper since it removes the effect of static clutter. To demonstrate the impact of clutter, the signatures without MTI filtering have been used as the input to the classifiers. Fig. 12 shows the results when the spectrogram signature has the filter removed, therefore the effect of background clutter on the Doppler signature is included. Both cases of “Known Prior” and “New” approaches for training the Doppler Bi-LSTM network are reported for the cluttered data.

There are certain cases, e.g. Person 1 and person 13, where the presence of clutter results in a decrease of approximately 12% from the regular case where the MTI filter is present. Similarly, there are cases where an extreme decrease is present, e.g. person 15 where for the “Known Prior” cases the presence of clutter has a 50% decrease in accuracy compared to the filtered/regular counterpart. Incidentally, there appears to be a marginal benefit in this case where prior knowledge is useful in classification as in average, there is a 3% difference between the “Known Prior” and “New” cases with clutter considered as shown in Table VI. With Fig. 12 and Table VI, we see that there is a decrease in performance for all the participants in both “Known Prior” and “New” cases. This suggests that filtering static clutter is essential to ensure accurate recognition of sequences with the proposed method.

### C. Comparison With Conventional Support Vector Machine

A simpler classifier, a support vector machine, is used with features derived from segmented windows of the whole

TABLE VI  
ACCURACY METRICS FOR THE DOPPLER BI-LSTM WITH VS WITHOUT STATIC CLUTTER FILTERED, FOR “KNOWN PRIOR” AND “NEW” TRAINING AND TESTING APPROACH. SVM RESULTS ALSO SHOWN FOR COMPARISON

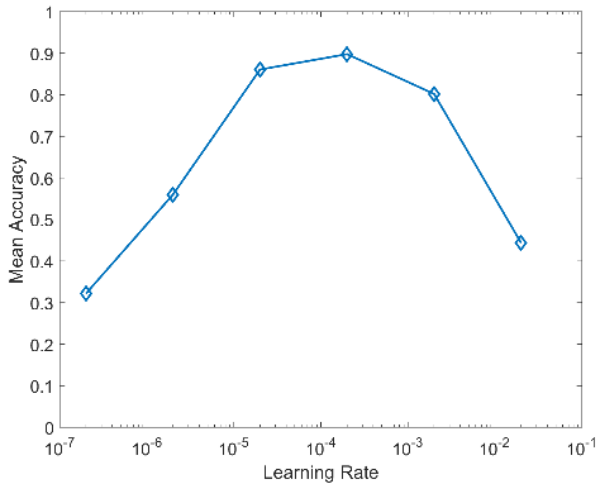
| Classifier                | Subject     | Mean | Standard deviation |
|---------------------------|-------------|------|--------------------|
| Doppler Bi-LSTM           | Known prior | 91   | 4                  |
| Doppler Bi-LSTM           | New         | 90   | 6                  |
| Cluttered Doppler Bi-LSTM | Known prior | 71   | 14                 |
| Cluttered Doppler Bi-LSTM | New         | 68   | 12                 |
| SVM                       | Both        | 66   | 11                 |

sequence to detect the activities, and the results can be then compared to those generated by the proposed LSTM networks. This is done to establish a benchmark for the Bi-LSTM architecture and to validate their use for this classification problem. SVM utilises as input a selection of features extracted from the centroid, bandwidth, and singular value decomposition of the spectrogram signature. This has been previously used in literature to identify discrete activities [6], [7], together with a sliding segmenting window. The window length was 4.5s and overlap was 90% and the kernel of the chosen SVM was linear. These parameters for the sliding segmenting window were selected as they resulted in the best performance in previous work [43].

Table VI shows a summary of the results for the SVM to compare them with those of the Bi-LSTM networks; the results for individual sequences were shown in Fig. 12. Note that there was only negligible difference in the SVM case between the “Known Prior” and “New” approach for training the classifier, hence they are reported together in the table. In general, the SVM results (green bars in Fig. 12 are lower than using Doppler Bi-LSTM networks (blue bars in Fig. 12. While the performance can be close in rare instances, e.g. for Person 1, in general the SVM is between 20 to 40 % lower when compared to the regular Doppler Bi-LSTM. This result emphasises the value of the proposed approach with a temporal-aware classifier such as the proposed Doppler Bi-LSTM for recognising continuous human activities.

### D. Optimising Learning Rate

A final point of note in the analysis is the selection of the learning rate. Fig 13 shows the sweep of the learning rate value with increments of a factor of ten for the Doppler Bi-LSTM architecture. It shows that even with an optimized architecture



**Fig. 13.** Parameter sweep of the learning rate with the best performing architecture and radar data domain: Doppler Bi-LSTM. A suboptimal initial learning rate can be as harmful to the classification accuracy as using a less suitable architecture or input domain.

and input domain if the initial learning rate selected is not optimal; the classification accuracy can degrade significantly to the levels where sub-optimal networks and input domains were used. Note that the mean accuracy presented in this figure refers to the average across the 45 diverse test sequences, as discussed in previous tables.

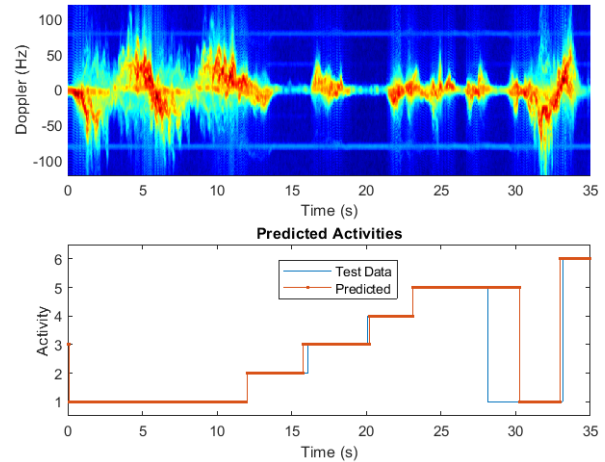
In suboptimal cases, be it with architecture, input domain or parameter selection, the presence of delayed classification with respect to the ground truth, transient states (spikes in the predicted labels), and complete misclassifications exhibited in [fig 10](#) increase significantly, thus reducing the performances.

### E. Line of Sight and Future Direction

One of the recurring questions about using radar for human activities monitoring is its performance when the target is not in a direct line of sight (LOS). In this case, the signal-to-noise ratio of the returned signal will vary, especially if the activity is being conducted at the edge of the antenna beamwidth. This is due to the combination of the effects from the attenuation of the antenna radiation pattern, the RCS fluctuation of the target, and the relation of the micro-Doppler shift with the aspect angle.

To assess this effect in a preliminary test, we recorded a sequence similar to [Fig.2](#) with an added A1: walking. The participant started the activities outside of the direct LOS before walking across the LOS to perform the A6: fall event at the edge of the beam. This meant that the aspect angle between the target and the radar varied from 0 degrees to approximately 30 degrees. This sequence was then tested with the Bi-LSTM trained with prior single LOS data, as discussed in the previous sections.

[Figure 14](#) shows the input, output, and ground truth of this test. The classifier makes two errors. The first is at 0 seconds where it detects A3: Standing for a small duration, and at 30 s, where it does not immediately detect the transition to A1: Walking. This is due to a non-movement gap visible in the



**Fig. 14.** Classifier input at the top sub-figure, ground truth in blue and the test outcome in orange in the bottom sub-figure for a test sequence with varying aspect angles and signal to noise ratio for the *Doppler-LSTM* network.

Doppler-time map between 28-30 seconds. Other than this, it appears to track the activities and the sequence well, with an average of 91.56% overlap between the predicted and ground truth data. This is consistent with other findings in literature where an angle up to 30 degrees gives an acceptable performance [15]. When this angle is larger (i.e. more than 30 degrees) or the target is out of the beam, performances may degrade further and a different radar deployment would necessary, such as for instance multistatic, interferometric, or multi-platform. Improving the robustness of classification for irrespective of the aspect angles is the subject of further research, but beyond the scope of this article.

## V. CONCLUSIONS AND FUTURE WORK

This paper analysed continuous sequences of experimental radar data to classify human activities and movements. Unlike the majority of current work in the literature, the data were not collected as separate recordings for each specific activity, but as a continuous stream where transitions between each activity can happen at any time and have unconstrained duration. These sequences were processed using novel recurrent Bidirectional LSTM networks that interpret the data as a temporal series, rather than as 2D images (i.e. matrices of pixels), as more conventional classification approaches based on convolutional networks do.

The proposed approach was validated with experimental data collected using a C-band FMCW radar with 15 participants performing 6 activities. Sequences with 3 different combinations of these activities were recorded to capture and classify diverse transitions between them. Different architectures for the recurrent networks were investigated, namely conventional LSTM and Bi-LSTM layers, as well as the effect of key hyperparameters such as the learning rate and of different formats of the input radar data, namely spectrograms and range-time sequences.

The results show that the proposed Bi-LSTM architecture outperforms unidirectional LSTM, as the former can capture

connections within the data in both the backward (past) and forward (future) temporal directions. This is particularly important for the classification of continuous sequences of human movement data, as the activity/movement performed at a specific time has a strong dependence on what was performed previously and influences what the subject can perform afterward. Classification accuracy over 90% was achieved for the optimized Bi-LSTM architecture across 45 different sequences of activities tested with a leave-one-subject-out cross-validation approach, demonstrating promising robustness and generalization capabilities for the proposed approach. Micro-Doppler data yielded higher accuracy than using range-time profiles as inputs to the networks. It is anticipated that the range information can be more relevant for classification when the subjects perform activities at an unfavorable aspect angle for Doppler-based measurements, or where a radar with wider bandwidth and finer range resolution is available (for example those operating in the mm-wave spectral region). It was also shown the benefit in carefully designing the network architecture and select its hyperparameters to fit the selected radar input data domain, as optimal architectures for the micro-Doppler domain did only provide sub-optimal performances when fed with range profiles.

An open problem faced by the radar research community for human monitoring is when multiple people are in the radar field of view and the recognition of activities while subjects are occluded by other subjects or objects. Techniques to separate the signatures of multiple subjects have been proposed using the fine range information of UWB radar [44] or the separation of the scatterers points of multiple subjects in the 3D radar cube [45]. These techniques could help separate and decompose the total signature into individual signatures that can then be subsequently processed by the proposed classification approach. The thorough investigation of this challenging scenario is left to future work, with different subjects, activities, environments, and trajectories or aspect angles.

#### ACKNOWLEDGMENT

The authors are grateful to the volunteers who supported the data collection, in particular, Dr. Lun Ma, Chang'an University, Xian, China.

#### REFERENCES

- [1] J. Le Kerneec *et al.*, "Radar signal processing for sensing in assisted living: The challenges associated with real-time implementation of emerging algorithms," *IEEE Signal Process. Mag.*, vol. 36, no. 4, pp. 29–41, Jul. 2019.
- [2] S. Z. Gurbuz and M. G. Amin, "Radar-based human-motion recognition with deep learning: Promising applications for indoor monitoring," *IEEE Signal Process. Mag.*, vol. 36, no. 4, pp. 16–28, Jul. 2019.
- [3] K. Chaccour, R. Darazi, A. H. El Hassani, and E. Andres, "From fall detection to fall prevention: A generic classification of fall-related systems," *IEEE Sensors J.*, vol. 17, no. 3, pp. 812–822, Feb. 2017.
- [4] H. Wang, D. Zhang, Y. Wang, J. Ma, Y. Wang, and S. Li, "RT-fall: A real-time and contactless fall detection system with commodity WiFi devices," *IEEE Trans. Mobile Comput.*, vol. 16, no. 2, pp. 511–526, Feb. 2017.
- [5] M. G. Amin, Y. D. Zhang, F. Ahmad, and K. C. D. Ho, "Radar signal processing for elderly fall detection: The future for in-home monitoring," *IEEE Signal Process. Mag.*, vol. 33, no. 2, pp. 71–80, Mar. 2016.
- [6] H. Li, A. Shrestha, H. Heidari, J. L. Kerneec, and F. Fioranelli, "A multisensory approach for remote health monitoring of older people," *IEEE J. Electromagn., RF Microw. Med. Biol.*, vol. 2, no. 2, pp. 102–108, Jun. 2018.
- [7] H. Li, A. Shrestha, H. Heidari, J. Le Kerneec, and F. Fioranelli, "Magnetic and radar sensing for multimodal remote health monitoring," *IEEE Sensors J.*, vol. 19, no. 20, pp. 8979–8989, Oct. 2019.
- [8] M. Ehatisham-Ul-Haq *et al.*, "Robust human activity recognition using multimodal feature-level fusion," *IEEE Access*, vol. 7, pp. 60736–60751, 2019.
- [9] C. Chen, R. Jafari, and N. Kehtarnavaz, "UTD-MHAD: A multimodal dataset for human action recognition utilizing a depth camera and a wearable inertial sensor," in *Proc. IEEE Int. Conf. Image Process. (ICIP)*, Quebec City, QC, Canada, Sep. 2015, pp. 168–172.
- [10] Z. Ahmad and N. Khan, "Human action recognition using deep multilevel multimodal ( $M^2$ ) fusion of depth and inertial sensors," *IEEE Sensors J.*, vol. 20, no. 3, pp. 1445–1455, Feb. 2020.
- [11] B. Vandersmissen *et al.*, "Indoor person identification using a low-power FMCW radar," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 7, pp. 3941–3952, Jul. 2018.
- [12] H. Sadreazami, M. Bolic, and S. Rajan, "CapsFall: Fall detection using ultra-wideband radar and capsule network," *IEEE Access*, vol. 7, pp. 55336–55343, 2019.
- [13] F. Luo, S. Poslad, and E. Bodanese, "Human activity detection and coarse localization outdoors using micro-Doppler signatures," *IEEE Sensors J.*, vol. 19, no. 18, pp. 8079–8094, Sep. 2019.
- [14] X. Bai, Y. Hui, L. Wang, and F. Zhou, "Radar-based human gait recognition using dual-channel deep convolutional neural network," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 12, pp. 9767–9778, Dec. 2019.
- [15] Y. Kim and H. Ling, "Human activity classification based on micro-Doppler signatures using a support vector machine," *IEEE Trans. Geosci. Remote Sens.*, vol. 47, no. 5, pp. 1328–1337, May 2009.
- [16] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, pp. 436–444, May 2015.
- [17] P. Cao, W. Xia, M. Ye, J. Zhang, and J. Zhou, "Radar-ID: Human identification based on radar micro-Doppler signatures using deep convolutional neural networks," *IET Radar, Sonar Navigat.*, vol. 12, no. 7, pp. 729–734, Jul. 2018.
- [18] Y. Kim and T. Moon, "Human detection and activity classification based on micro-Doppler signatures using deep convolutional neural networks," *IEEE Geosci. Remote Sens. Lett.*, vol. 13, no. 1, pp. 8–12, Jan. 2016.
- [19] Y. Yang, C. Hou, Y. Lang, T. Sakamoto, Y. He, and W. Xiang, "Omnidirectional motion classification with monostatic radar system using micro-Doppler signatures," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 5, pp. 3574–3587, May 2020.
- [20] M. S. Seyfioglu, B. Erol, S. Z. Gurbuz, and M. G. Amin, "DNN transfer learning from diversified micro-Doppler for motion classification," *IEEE Trans. Aerosp. Electron. Syst.*, vol. 55, no. 5, pp. 2164–2180, Oct. 2019.
- [21] M. S. Seyfioglu, A. M. Özbayoglu, and S. Z. Gürbüz, "Deep convolutional autoencoder for radar-based classification of similar aided and unaided human activities," *IEEE Trans. Aerosp. Electron. Syst.*, vol. 54, no. 4, pp. 1709–1723, Aug. 2018.
- [22] X. Li, Y. He, and X. Jing, "A deep multi-task network for activity classification and person identification with micro-Doppler signatures," in *Proc. Int. Radar Conf. (RADAR)*, Toulon, France, Sep. 2019, pp. 1–5.
- [23] X. Shi, Y. Li, F. Zhou, and L. Liu, "Human activity recognition based on deep learning method," in *Proc. Int. Conf. Radar (RADAR)*, Brisbane, QLD, Australia, Aug. 2018, pp. 1–5.
- [24] I. Alnujaim, D. Oh, and Y. Kim, "Generative adversarial networks for classification of micro-Doppler signatures of human activity," *IEEE Geosci. Remote Sens. Lett.*, vol. 17, no. 3, pp. 396–400, Mar. 2020.
- [25] B. Erol, S. Z. Gurbuz, and M. G. Amin, "Motion classification using kinematically sifted ACGAN-synthesized radar micro-Doppler signatures," *IEEE Trans. Aerosp. Electron. Syst.*, early access, Jan. 27, 2020, doi: [10.1109/TAES.2020.2969579](https://doi.org/10.1109/TAES.2020.2969579).
- [26] H. G. Doherty, L. Cifola, R. I. A. Harmanny, and F. Fioranelli, "Unsupervised learning using generative adversarial networks on micro-Doppler spectrograms," in *Proc. 16th Eur. Radar Conf. (EuRAD)*, Paris, France, 2019, pp. 197–200.
- [27] G. Duggal, S. Vishwakarma, K. V. Mishra, and S. S. Ram, "Doppler-resilient 802.11ad-based ultra-short range automotive joint radar-communications system," *IEEE Trans. Aerosp. Electron. Syst.*, early access, May 13, 2020, doi: [10.1109/TAES.2020.2990393](https://doi.org/10.1109/TAES.2020.2990393).

- [28] B. Jokanović and M. Amin, "Fall detection using deep learning in range-Doppler radars," *IEEE Trans. Aerosp. Electron. Syst.*, vol. 54, no. 1, pp. 180–189, Feb. 2018.
- [29] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [30] A. Graves, A.-R. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process.*, May 2013, pp. 6645–6649.
- [31] A. Graves and J. Schmidhuber, "Framewise phoneme classification with bidirectional LSTM and other neural network architectures," *Neural Netw.*, vol. 18, nos. 5–6, pp. 602–610, Jul. 2005.
- [32] Z. Yu *et al.*, "Using bidirectional lstm recurrent neural networks to learn high-level abstractions of sequential features for automated scoring of non-native spontaneous speech," in *Proc. IEEE Workshop Autom. Speech Recognit. Understand. (ASRU)*, Scottsdale, AZ, USA, Dec. 2015, pp. 338–345.
- [33] L. Gao, Z. Guo, H. Zhang, X. Xu, and H. T. Shen, "Video captioning with attention-based LSTM and semantic consistency," *IEEE Trans. Multimedia*, vol. 19, no. 9, pp. 2045–2055, Sep. 2017.
- [34] G. Klarenbeek, R. I. A. Harmanny, and L. Cifola, "Multi-target human gait classification using LSTM recurrent neural networks applied to micro-Doppler," in *Proc. Eur. Radar Conf. (EURAD)*, Nuremberg, Germany, Oct. 2017, pp. 167–170.
- [35] M. Wang, G. Cui, X. Yang, and L. Kong, "Human body and limb motion recognition via stacked gated recurrent units network," *IET Radar, Sonar Navigat.*, vol. 12, no. 9, pp. 1046–1051, Sep. 2018.
- [36] M. Wang, Y. D. Zhang, and G. Cui, "Human motion recognition exploiting radar with stacked recurrent neural network," *Digit. Signal Process.*, vol. 87, pp. 125–131, Apr. 2019.
- [37] S. Yang, J. L. Kerneç, F. Fioranelli, and O. Romain, "Human activities classification in a complex space using raw radar data," in *Proc. Int. Radar Conf. (RADAR)*, Toulon, France, Sep. 2019, pp. 1–4.
- [38] X. Li, Y. He, Y. Yang, Y. Hong, and X. Jing, "LSTM based human activity classification on radar range profile," in *Proc. IEEE Int. Conf. Comput. Electromagn. (ICCEM)*, Shanghai, China, Mar. 2019, pp. 1–2.
- [39] C. Ding *et al.*, "Continuous human motion recognition with a dynamic range-Doppler trajectory method based on FMCW radar," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 9, pp. 6821–6831, Sep. 2019.
- [40] V. C. Chen, F. Li, S.-S. Ho, and H. Wechsler, "Micro-Doppler effect in radar: Phenomenon, model, and simulation study," *IEEE Trans. Aerosp. Electron. Syst.*, vol. 42, no. 1, pp. 2–21, Jan. 2006.
- [41] M. Schuster and K. K. Paliwal, "Bidirectional recurrent neural networks," *IEEE Trans. Signal Process.*, vol. 45, no. 11, pp. 2673–2681, Nov. 1997, doi: [10.1109/78.650093](https://doi.org/10.1109/78.650093).
- [42] C. M. Bishop, *Neural Networks for Pattern Recognition*. London, U.K.: Oxford Univ. Press, 1995.
- [43] H. Li, A. Shrestha, H. Heidari, J. L. Kerneç, and F. Fioranelli, "Activities recognition and fall detection in continuous data streams using radar sensor," in *Proc. IEEE MTT-S Int. Microw. Biomed. Conf. (IMBioC)*, Nanjing, China, May 2019, pp. 1–4.
- [44] T. Sakamoto, T. Sato, P. J. Aubry, and A. G. Yarvoy, "Texture-based automatic separation of echoes from distributed moving targets in UWB radar signals," *IEEE Trans. Geosci. Remote Sens.*, vol. 53, no. 1, pp. 352–361, Jan. 2015.
- [45] H. Du, T. Jin, M. Li, Y. Song, and Y. Dai, "Detection of multi-people micro-motions based on range-velocity-time points," *Electron. Lett.*, vol. 55, no. 23, pp. 1247–1249, Nov. 2019.



**Aman Shrestha** (Student Member, IEEE) received the B.Eng. degree in electrical and electronic engineering from the University of Birmingham. He is currently pursuing the Ph.D. degree with the School of Engineering, University of Glasgow, working on application of radar systems and radar signal processing for assisted living with deep learning.



**Haobo Li** (Student Member, IEEE) received the B.Eng. degree in electrical and electronic engineering from Northumbria University, Newcastle upon Tyne, and the M.S. degree in communication and signal processing from the University of Newcastle in 2015 and 2016, respectively. He is currently pursuing the Ph.D. degree with the School of Engineering, University of Glasgow. He is also working on information fusion of multiple sensing technologies for assisted living applications and gesture recognition.



**Julien Le Kerneç** (Senior Member, IEEE) received the B.Eng. and M.Eng. degrees in electronic engineering from the Cork Institute of Technology, Ireland, in 2004 and 2006, respectively, and the Ph.D. degree in electronic engineering from Pierre and Marie Curie University, France, in 2011. He is currently a Lecturer with the School of Engineering, University of Glasgow. He is also seconded as a Lecturer with the University of Electronic Science and Technology of China and an Adjunct Associate Professor with the

ETIS Laboratory, University Cergy-Pontoise, France. His research interests include radar system design, software-defined radio/radar, signal processing, and health applications.



**Francesco Fioranelli** (Senior Member, IEEE) received the Ph.D. degree from Durham University, U.K., in 2014. He was a Postdoctoral Research Associate with University College London from 2014 to 2016, and a Lecturer with the School of Engineering, University of Glasgow from 2016 to 2019. He is an Assistant Professor of Microwave Sensing Signals and Systems Section with the Department of Microelectronics, TU Delft, The Netherlands. His research interests include development of radar systems and

radar-based classification and machine learning methods for applications, such as human signatures analysis for healthcare and security, drones and UAVs' detection and classification, automotive radar, wind farm, and sea clutter characterization.