# GraSens: A Gabor Residual Anti-aliasing Sensing Framework for Action Recognition using WiFi

Yanling Hao
Queen Mary University
of London
London, UK
Email: yanling.hao@qmul.ac.uk

Zhiyuan Shi
Onfido Research London
London, UK
Email: zhiyuan.shi@onfido.com

Xidong Mu
Queen Mary University
of London
London, UK
Email: x.mu@qmul.ac.uk

Yuanwei Liu
Queen Mary University
of London
London, UK
Email: yuanwei.liu@qmul.ac.uk

*Abstract*—WiFi-based human action recognition (HAR) has been regarded as a promising solution in applications such as smart living and remote monitoring due to the pervasive and unobtrusive nature of WiFi signals. However, the efficacy of WiFi signals is prone to be influenced by the change in the ambient environment and varies over different sub-carriers. To remedy this issue, we propose an end-to-end Gabor residual anti-aliasing sensing network (GraSens) to directly recognize the actions using the WiFi signals from the wireless devices in diverse scenarios. In particular, a new Gabor residual block is designed to address the impact of the changing surrounding environment with a focus on learning reliable and robust temporal-frequency representations of WiFi signals. In each block, the Gabor layer is integrated with the anti-aliasing layer in a residual manner to gain the shift-invariant features. Furthermore, fractal temporal and frequency self-attention are proposed in a joint effort to explicitly concentrate on the efficacy of WiFi signals and thus enhance the quality of output features scattered in different subcarriers. Experimental results throughout our wireless-vision action recognition dataset (WVAR) and three public datasets demonstrate that our proposed GraSens scheme outperforms state-of-the-art methods with respect to recognition accuracy.

## I. INTRODUCTION

Human action recognition (HAR) has attracted considerable attention in a range of applications, such as assisted living [1], behavior analysis [2], and health monitoring [3]. Many pioneering actions sensing attempts [4]–[6] have continuously emerged and developed in recent years to enhance measurement data and expand signal acquisition range [7]. These sensing techniques motivate the breakthrough of long-time monitoring in a non-intrusive way [6], [8]–[10].

The radio frequency (RF)-based technique is one of the most promising technologies among other action sensing technologies to localize people and track their motion [11], [12]. This attempt draws on the propagation of electromagnetic (EM) waves which are almost distributed at everyone's home. Benefit from the ubiquitous deployment, using WiFi signals for HAR in the indoor environment, is an economic solution [13], [14]. Furthermore, WiFi-based solutions have no requirements of line-of-sight (LOS) thereby enabling larger detection areas than vision-based techniques [8], [9]. Therefore, WiFi-based HAR methods have received increasing attention [7].

Extant researches have demonstrated the great potential of employing WiFi signals as a sensing approach [10]. Previously, most techniques for HAR are presented based on hand-crafted features from WiFi signals [15]. In essence, WiFi signals are susceptible to severe multipath and random noise in indoor surroundings. Hence, these manually designed features based mechanisms have certain limitations due to their heavy dependence on prior knowledge [11]. Furthermore, the efficacy of WiFi signals for HAR scatters over different sub-carriers since certain bands are sensitive to certain movements. Therefore, it is of vital importance to explore the problem of how to non-manually obtain robust and reliable representations from the WiFi signals. Deep learning is capable of automatic feature selection and has emerged as a new paradigm for mining the temporal-frequency information in the WiFi signals in diverse scenarios.

Deep learning has been evolving as a promising solution for HAR over the past few years [16], [17]. Past deep learning methods however are prone to cause distortions after downsampling operation [18]. In deep learning networks, the downsampling operation is broadly utilized to reduce parameters and computation cost [19]. After the sampling operation, high-frequency information signals degenerate into completely different ones, which further disturbs the feature information [20]. The standard solution of embedding a low-pass filter before sampling [21] is unsatisfying because it degrades performance.

To remedy the above limitations, in this paper, an end-to-end Gabor residual anti-aliasing sensing (GraSens) network is proposed for HAR in varied environments. The architecture exploiting the reliable temporal-frequency representations from wireless signals is in an end-to-end style. The main contributions are summarized as follows:

1) We propose a Gabor residual anti-aliasing sensing network to directly recognize the activities based on the WiFi signals from wireless devices such as smartphones and routers in diverse scenarios.

2) We design a Gabor residual block for exploiting reliable and robust WiFi signals representations to mitigate the influence of the change in the ambient environment. Specifically, the Gabor layer in this block is integrated with anti-aliasing operation in a residual manner to gain the shift-invariant features.

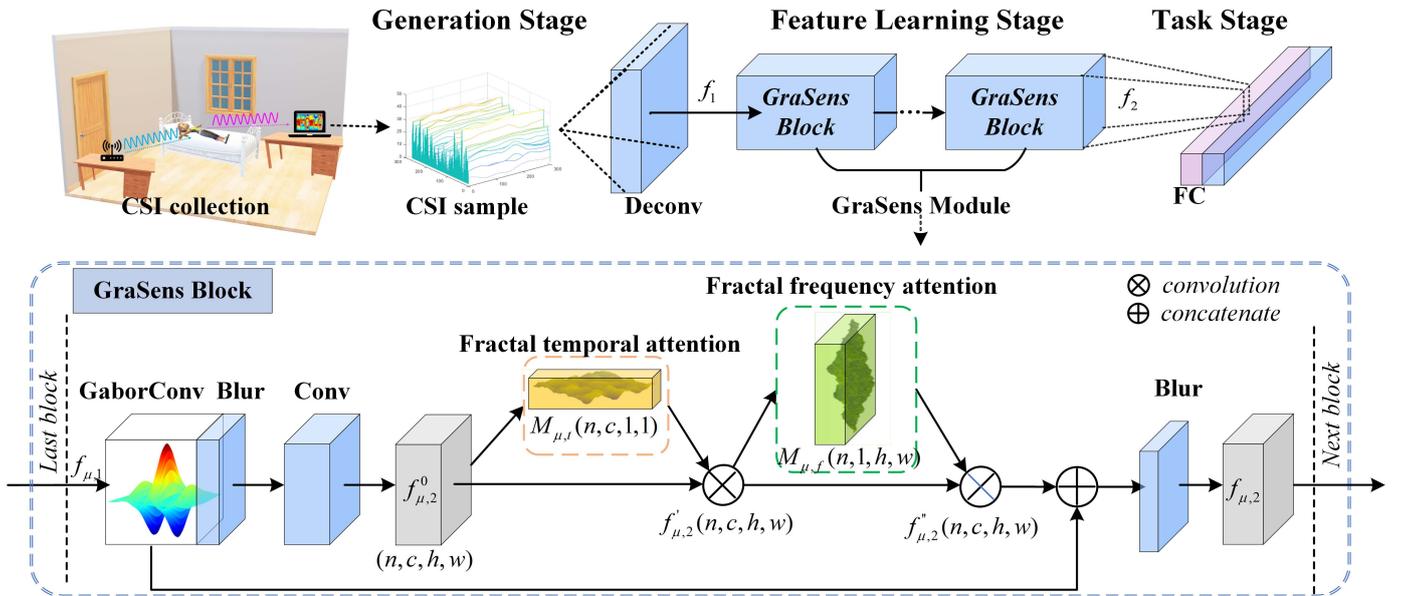3) We design a fractal temporal and frequency self-attention

Fig. 1: Overview of the proposed GraSens method.

mechanism to jointly explore the frequency and temporal continuity inside WiFi signals to enhance the quality of output features scattered in different subcarriers.

4) We conduct experiments on our proposed wireless-vision action recognition dataset and the other three public datasets. The experimental results show that our method is robust over different scenes and outperforms competitive baselines with a good margin on the recognition accuracy.

## II. RELATED WORK

Current researches on HAR can be loosely classified into two types, namely, video-based methods [5], [22] and RF-based methods [12].

### A. Video-based human action recognition

Video-based sensing methods have been prevailing in human action recognition. These methods capture image sequences by exploiting the camera and realize human action recognition using classification algorithms. Generally, they can be categorized into three groups: part-based frameworks [22], two-step frameworks [5], multi-stream model frameworks. In the part-based HAR, body parts are firstly detected separately and further assembled for human pose estimations such as DeepCut [22]. However, the assembled pose is prone to be ambiguous when more than one person gathers together and causes occlusion. Moreover, the part-based scheme is unable to recognize human pose globally since it focuses only on the second-order dependence of human body parts. As for the two-step framework, human bounding boxes are first detected and the poses within each box are then estimated such as Faster RNN [23]. In this way, the quality of action recognition is highly attached to the accuracy of the detected human bounding boxes. In the presence of the multiple streams

framework like RGB flow and optical flow, it aims to improve the accuracy of action recognition by characterizing and integrating the patterns from various stream sources such as SlowFast [24]. However, most of the video-based methods are susceptible to ambient surroundings such as occlusion, lightning and privacy concerns, etc. To break the obstacles of the demand for line-of-sight (LOS), a time-series generative adversarial network (TS-GAN) [25] is proposed to generate inferences and hallucinations in recognizing videos related to unseen actions. In fact, such hallucinations tend to produce errors due to the deformable ability of the human body.

### B. WiFi based human action recognition

RF-based techniques include radars [8], LiDARs [26] and WiFi devices [12]. Radar and LiDARs sensors demand dedicated and specially designed hardware. In contrast, WiFi devices are ubiquitously deployed since they are cost-effective and power-efficient. Besides, WiFi devices are free from the influences of illumination and privacy concerns in comparison to video-based methods. Recently, an amount of WiFi-based sensing systems were developed for human action recognition, such as WifiU [27] and RT-Fall [28]. Yet, previous systems are fairly coarse. These systems either locate only one single limb or produce a rough and static representation of the human body [12]. Most of the methods often target the general perception, for example, the rough classification [12] and indoor localization [15]. To mitigate the situation, some researchers attempt to simulate 2D or 3D skeletons based on wireless signals for person perception [7]. Other researchers simulate the WiFi arrays to enhance the accuracy of recognition and localization [29]. These researches illuminate the optimizing applications of WiFi-based HAR in varied environmental conditions. Recently, Alazrai et al. proposed an end-to-end framework E2EDLF [30] to recognize human-to-human inter-

actions by sophisticated and careful construction of the input CSI image.

## III. ARCHITECTURE FOR WIFI SENSING VIA GABOR RESIDUAL ANTI-ALIASING

As seen in Fig. 1, the proposed GraSens is designed and conceived to fully exploit and explore the data collected from off-the-shelf commercial WiFi devices in an end-to-end style. Three stages can be generalized, namely generation stage, feature learning stage, and task stage. Specifically, the generation stage is aiming to enable the raw WiFi channel state information (CSI) data compatible with the input of the network while preserving the original frequency and temporal information. The feature learning stage as shown in the bottom part of Fig. 1 is defined as Gabor residual anti-aliasing attention module, which puts forward the up-sampled CSI samples for feature maps generation. This stage can greatly mitigate the influence of the ambient noises that are confused with the action signals, and improve the quality of output features from CSI information scattered in different subcarriers. These learned features are further fed to fully connected layers for a particular task in the last stage.

### A. The proposed GraSens network

*1) Generation Stage:* To preserve the temporal as well as frequency information within the CSI signals, the raw CSI signals are transformed into a set of CSI tensors with learnable parameters in the generation stage seen in Fig. 2(a). Firstly, the raw CSI signals of an action segment as shown in Fig. 1 are converted into a series of CSI tensors, aiming to interpret the action with multiple aspects. After this, all the CSI tensors are up-sampled by the deconvolution operation adapted to the network. The principle of WiFi-based sensing is to recognize the influence of perceived objects on the transmitted signals [28]. Generally, a WiFi system can be modeled and summarised as follows:

$$\boldsymbol{B}_s(i) = \boldsymbol{\gamma}_s(i)\boldsymbol{A}_s(i) + \boldsymbol{\theta}, \tag{1}$$

where $s \in [1, \cdots, N_s]$ depicts the index of the orthogonal frequency-division multiplexing (OFDM) subcarriers employed in the WiFi device, $N_s$ represents the total number of the OFDM subcarriers. $i$ defines the index of the transmitted and received packets. The $i^{th}$ transmitted and received packets pertinent to the OFDM subcarrier frequency $s$ are specified as $\boldsymbol{A}_s(i)$ and $\boldsymbol{B}_s(i)$, respectively. $\boldsymbol{\theta}$ represents the received noise, and a complex-valued matrix $\boldsymbol{\gamma}_s$ constitutes the CSI measurements for the OFDM subcarrier frequency $s$. $\boldsymbol{\gamma}_s$ is of dimensions $N_T \times N_R$ whose horizontal and vertical co-ordinates indicate the number of transmitting and receiving antennas, respectively.

In each time serial sequence, the raw CSI signals are arranged in a 2D matrix of dimensions $\boldsymbol{\gamma} \times I$ with $\boldsymbol{\gamma} = N_T \times N_R \times N_S$ where $I$ indicates the index of packets recorded in a specific CSI time serial. A sliding window along the time axis divides the CSI signals into a bank of overlapped segments as CSI tensor $CSI(i)$ of size $\phi \times \boldsymbol{\gamma}$. $\phi$ defines the



(a)The CSI signal of throwing without occlusion  (b)The CSI signal of throwing with partial occlusion  (c)The CSI signal of throwing with full occlusion
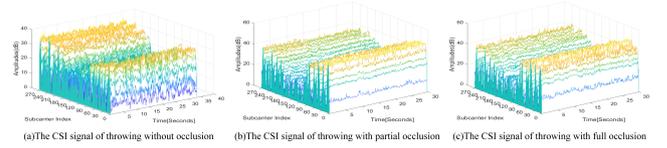
Fig. 2: The CSI signal of throwing. (a)-(c) are the CSI signals of throwing in scenes without occlusion, with partial occlusion, and with full occlusion, respectively.

number of packets and $\upsilon$ implies the overlap between every two adjacent segments, where $\upsilon \leq \phi$ and $i \leq I/\upsilon$.

The CSI samples are further put forward to the deconvolution layer. The deconvolution layer serves as an up-sampling layer to up-sample feature maps of the input CSI tensor and preserves the connectivity pattern. In the up-sampling process, the input CSI tensor is enlarged and densified by cross-channel convolutions with multiple filters. The spatial and frequency information in each channel is expanded and encoded into spatially-coded maps. In comparison with the extant resizing methods, the benefit of the deconvolution layers is that the parameters are trainable. During the training, the weights of deconvolution layers are constantly updated and refined. The CSI samples are up-sampled to be processed by feature learning modules as follows:

$$f_1 = Deconv(CSI). \tag{2}$$

where $Deconv(\cdot)$ is the deconvolution operation.

*2) Feature Learning Stage:* As depicted in Fig. 1, a Gabor residual anti-aliasing sensing module is proposed for shift-invariant feature learning. This GraSens module consists of several Gabor residual anti-aliasing blocks. In each block, a Gabor convolution layer filter replaces the first convolution layer in a traditional residual module and serves as initialization to gain more discriminative power. After this, an anti-aliasing layer is further added to keep the output feature maps shift-invariant. For block $\mu$, given the intermediate feature map $f_1 \in \mathcal{R}^{C \times H \times W}$ as the input, the output features can be generated as follows:

$$f_{\mu,2}^0 = Conv(Blur(GaborConv(f_1))). \tag{3}$$

where $GaborConv(\cdot)$ is the Gabor convolution operation and $Blur(\cdot)$ is the anti-aliasing operation. To explicitly concentrate on the efficacy of WiFi signals, GraSens sequentially infers a 1D fractal dimension based temporal attention map $M_{\mu,t} \in R^{C \times 1 \times 1}$ and a 2D fractal dimension based frequency attention map $M_{\mu,f} \in R^{C \times H \times W}$ as shown in Fig. 1. In short, the whole attention process can be generalized as follows:

$$\begin{aligned} f_{\mu,2}^{'} &= M_{\mu,t}(f_{\mu,2}) \otimes f_{\mu,2}, \\ f_{\mu,2}^{''} &= M_{\mu,f}(f_{\mu,2}^{'}) \otimes f_{\mu,2}^{'}, \end{aligned} \tag{4}$$

where $\otimes$ indicates the element-wise multiplication. The unique asset of multiplication locates in the way of duplication of attention values. Intuitively, temporal attention values replicated

along the frequency axis and vice versa. Herein, the refined output $f_{\mu,2}$ of stacked block $\mu$ can be formulated as follows:

$$f_{\mu,2} = Blur(f''_{\mu,2} \oplus f_1), \tag{5}$$

where $\oplus$ is the concatenate operation. Fig. 1 describes the calculation process of each attention map. After several blocks, $f_2$ is the final output temporal and frequency representation. The following section III-B describes the details of each attention module. The feature learning progress of GraSens module is as depicted in Algorithm 1.

---

**Algorithm 1** Feature Learning

**Input:** The up-sampled CSI sample $f_1$
**Output:** The output feature maps $f_2$ of GraSens module
1: Choose the number of stacked GraSens blocks as $\lambda$;
2: Initialize the block $\mu = 1$;
3: **repeat**
4:    **for** block $\mu$ **do**
5:       Update the Gabor anti-aliasing output $f^0_{\mu,2} \leftarrow f_1$ using Eqs. (3), (8) and (9);
6:       Update the fractal self-attention output $f''_{\mu,2} \leftarrow f^0_{\mu,2}$ using Eqs. (10)- (12);
7:       Update the anti-aliasing output $f_{\mu,2} \leftarrow f''_{\mu,2}$ using Eqs. (5) and (9);
8:    **end for**
9:    Renew the input for next block $f_1 = f_{\mu,2}$;
10:   Move to next block $\mu = \mu + 1$;
11: **until** $\mu = \lambda$;
12: Return $f_2 = f_{\lambda,2}$ and forward to the task stage.

---

*3) Task Stage:* During the task stage, the learned frequency and temporal features are fed to one fully connected layer to generate the outputs for a particular task. In the training of GraSens, the loss is computed by the activation function and loss function. In this way, the difference between the outputs of the GraSens network $f_3$ and the ground-truth G can be measured by the loss. The output $f_3$ is formulated as follows:

$$f_3 = Blur(FC(f_2)), \tag{6}$$

The cross-entropy loss is a basic option to be applied to optimize GraSens and given by:

$$\mathcal{L} = \sum_{j=1}^{J} f_3{}^j \log(G^j). \tag{7}$$

where $j$ is the snippet number of input training CSI samples. In addition, we utilize the Stochastic Gradient Descent with Momentum to learn the parameters.

TABLE I: Classification accuracy of the dataset WVAR.

| Methods | fall_down | throw | push | kick | punch | jump | phone_talk | seat | drink | OA |
|---|---|---|---|---|---|---|---|---|---|---|
| SVM | **1.00** | 0.92 | 0.90 | 0.94 | 0.93 | 0.94 | 0.91 | 0.88 | **1.00** | 0.94 |
| WNN | **1.00** | **1.00** | **1.00** | 0.86 | 0.88 | **1.00** | **1.00** | 0.81 | **1.00** | 0.94 |
| GraSens | **1.00** | **1.00** | 0.95 | **0.97** | **0.99** | **1.00** | 0.88 | **0.90** | 0.92 | **0.95** |

TABLE II: Classification accuracy of the dataset WAR.

| Methods | lie_down | fall | run | sit_down | stand_up | walk | OA |
|---|---|---|---|---|---|---|---|
| RF [31] | 0.53 | 0.60 | 0.81 | 0.88 | 0.49 | 0.57 | 0.65 |
| HMM [32] | 0.52 | 0.72 | 0.92 | 0.96 | 0.76 | 0.52 | 0.73 |
| LSTM [33] | **0.95** | 0.94 | **0.97** | 0.81 | 0.83 | **0.93** | 0.91 |
| SVM | 0.91 | 0.96 | 0.93 | 0.96 | 0.71 | 0.87 | 0.93 |
| WNN | 0.93 | 0.93 | 0.93 | **0.98** | 0.90 | 0.86 | 0.95 |
| GraSens | 0.94 | **0.97** | 0.95 | **0.98** | **0.91** | 0.85 | **0.96** |

*B. GraSens Module*

*1) Gabor Filtering based Anti-aliasing:* As for each GraSens block, the Gabor layer builds a convolution kernel library for feature extraction. To obtain the strong auxiliary feature information, the Gabor convolution kernel group is optimized by the network training and further convolved with the CSI samples. Generally, the Gabor function describes a complex sinusoid modulated by Gaussian in accordance with monotonicity and differentiability, i.e.,

$$\begin{aligned} GaborConv &= g(x, y, \varpi, \theta, \psi, \sigma) \\ &= \exp(-\tfrac{x'^2 + y'^2}{2\sigma^2}) \cos(\varpi x' + \psi), \end{aligned} \tag{8}$$

where $x' = x \cos\theta + y \sin\theta$, and $y' = -x \cos\theta + y \cos\theta$. Gabor layers prove to be efficient for spatially localized features extracting [34]. To extract the features from the WiFi signals, a set of Gabor filters are used as ref [35]. Frequencies $\varpi_n$ of the Gabor filters is obtained by $\varpi_n = \frac{\pi}{2}\sqrt{2}^{-(n-1)}$, $n = 1, 2, \ldots, 5$. The orientations $\theta_m$ is set as $\theta_m = \frac{\pi}{8}(m-1)$, where $m = 1, 2, \ldots, 8$. In addition, the $\sigma$ is defined by the relationship between $\sigma$ and $\varpi$ where $\sigma \approx \frac{\pi}{\varpi}$. $\psi$ follows the uniform distribution U$(0,\pi)$. Accordingly, the Gabor Layer weights in this paper are initialized similarly.

Subsequently, the anti-aliasing layer is leveraged to enable the extracted feature shift-invariant. The anti-aliasing layer serves as two steps. To begin with, a set of low-pass filters $\Psi$ are arranged and generated in terms of varied spatial locations and channel groups within each GraSens block. After than, the predicted filters are adopted and applied back onto the input feature maps on account of anti-aliasing. We assume an input feature $X$. To be specific, a low-pass filter $\Psi_{i,j}^{p,q}$, for example, a 3×3 convolution filter, is generated to down-sample the input feature $X$ over each spatial location $(i, j)$ as follows:

$$Blur = \sum_{p,q \in \Omega} \Psi_{i,j}^{p,q} \cdot X_{i+p,j+q}. \tag{9}$$

*2) Fractal Dimension based Self-Attention:* Fractal describes unusual objects of irregular shapes which have a high degree of complex properties. Fractal dimension can indicate the degree of the complexity of objects, such as the irregular WiFi signals. For the convenience, a general expression has been defined to measure the fractal dimension as follows:

$$FD = -\lim_{\varepsilon \to 0} \frac{log(\eta(\varepsilon))}{log(\varepsilon)}, \tag{10}$$

where $\eta$ measures self-similarity and $\varepsilon$ denotes the scale. In our work, $FD$ is employed to calculate the fractal dimension
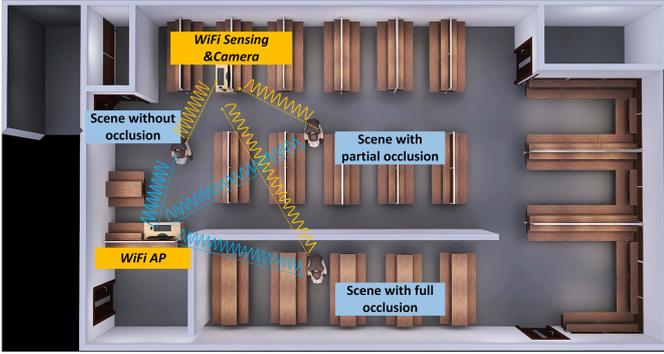
Fig. 3: Three experiment scenes of WVAR dataset.

of feature maps along with the frequency and temporal domain.

***Fractal temporal attention module.*** Each channel within a feature map can reflect the diverse temporal characteristics of the input CSI samples. Inspired by the CBAM [36], we calculate the fractal dimensions for all the frequencies in feature maps input as the temporal attention as follows:

$$M_{\mu,t}(f_{\mu,2}) = \xi(MLP(FD(f_{\mu,2}))), \quad (11)$$

where $\xi$ implies the sigmoid function. $MLP$ specifies a multi-layer perceptron operation.

***Fractal frequency attention module.*** Cross-channels within a feature map can capture the frequency characteristics. For this purpose, a frequency attention map is generated to exploit the cross-channel relationship of features. Fractal dimensions across the channel are utilized to generate one feature map as the fractal feature maps. Those fractal feature maps are further fed to a standard convolution layer and thus generate the frequency attention map. In brief, the fractal frequency attention is calculated as follows:

$$M_{\mu,f}(f'_{\mu,2}) = \xi(Conv(FD(f'_{\mu,2}))), \quad (12)$$

where $Conv$ represents a convolution operation.

## IV. EXPERIMENTS

### A. Datasets

**Our WVAR dataset**. WVAR collection was implemented in one spacious office apartment by 2 volunteers who performed 9 activities with five repeated trials in different simulating occlusion occasions as seen in Fig. 3. The experimental hardware as seen in Fig. 1 constitutes two desktop computers as transmitter and receiver, both of which are carried out in IEEE 802.11n monitor mode operating at 5.4 GHz with a sampling rate of 100 Hz. WVAR also contains the synchronized video data recorded at 20 FPS, i.e. every frame is corresponding to five CSI packets.

Table IV shows the classification accuracy of the dataset CSNLOS. We test two LOS scenarios' data E1 and E2. The results of GraSens rank first compared to all other two methods in two LOS scenes E1 and E2. As for E1, GraSens achieves the best results by 3% average accuracy higher than SVM [40].

With regard to E2, the performance of GraSens is better except for no_movement and walking which still are comparable with those of SVM [40]. In other words, GraSens has good robustness in comparison to the other two models.

**WAR, HHI, and CSLOS**. The public available dataset WAR [33] consists of 6 persons, 6 activities with 20 trials for each in an indoor office. The sampling rate is 1 kHz.

The publicly available CSI dataset of HHIs [41] is composed of 12 different human-to-human interactions (HHI) which performed by 40 distinct pairs of subjects in an indoor environment inside an office with 10 different trials, e.g. approaching, departing, hand_shaking, etc.

Another public available cross-scene dataset (CSLOS) [42] is provided by the same group as the HHI. LOS contains five experiments in three different indoor environments, where two are of LOS nature and the third environment is of a non-line-of-sight (NLOS) nature. 30 different subjects were included with 20 repeated trials for each of the experiments in terms of the variations of human movements.

*1) Evaluation Metrics:* Accuracy and precision are utilized in the sort of performance evaluation. Accuracy defines the percentage of total actions classified correctly. Precision reflects the correct percentage of classified actions from all predicted ones. It should be underlined that false positives are also included in precision. Both metrics are denoted as follows: $Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$ and $Precision = \frac{TP}{TP+FP}$, where $TP$, $FP$, $TN$ and $FN$ represent the true and false positives and negatives, respectively.

### B. Comparison with state-of-the-art methods

*1) Quantitative Results:* We compare GraSens with several state-of-the-art approaches on all four datasets, namely WVAR, WAR, HHI, and CSLOS. Apart from SVM and WNN, we used the reported accuracy of their original paper unless otherwise stated for comparison.

Table I illustrates the classification accuracy of the dataset WVAR. GraSens surpassed all other methods in most of the actions with an OA of 95%, which is slightly higher than these of SVM and WNN 1%. The reason behind this may be due to the fact that the dataset WVAR is relatively too small to reflect the advantages of GraSens. In addition, it can be observed that some action classes (i.e. push, phone talk, and drink) of GraSens obtained a slightly lower accuracy than WNN. The possible reason for this can be that all are simple activities whose changes in waveform characteristics over time were similar. Compared with WNN, GraSens has fewer advantages in this case.

Table II shows the results on the dataset WAR. GraSens outperforms all the baselines with a large margin of 5% than LSTM and 1% than our baseline WNN. Notably, WNN has the same network structure as GraSens. This confirms the effectiveness of the design of our network. Compared with the results of RF, HMM, and SVM, the results of GraSens had obvious improvements in all the six activities. This reason behind this is due to the fact that GraSens can extract more robust and shift-invariant features than machine

TABLE III: Classification accuracy of the dataset HHI.

| Methods | approaching | departing | hand_shaking | high five | hugging | kicking_left_leg | kicking_right_leg | pointing_left_hand | pointing_right_hand | punching_left_hand | punching_right_hand | pushing | OA |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GoogleNet [37] | 0.93 | 0.93 | 0.79 | 0.76 | 0.64 | 0.54 | 0.50 | 0.78 | 0.77 | 0.59 | 0.59 | 0.68 | 0.71 |
| ResNet-18 [38] | 0.92 | 0.90 | 0.85 | 0.79 | 0.77 | 0.68 | 0.60 | 0.82 | 0.80 | 0.60 | 0.65 | 0.76 | 0.76 |
| Squeeze-Net [39] | 0.95 | 0.93 | 0.83 | 0.76 | 0.70 | 0.66 | 0.62 | 0.78 | 0.79 | 0.60 | 0.72 | 0.74 | 0.76 |
| E2EDLF [30] | 0.96 | 0.92 | 0.89 | 0.84 | 0.86 | **0.78** | **0.82** | **0.85** | **0.90** | 0.73 | **0.80** | 0.86 | 0.85 |
| SVM | **0.99** | 0.96 | 0.90 | 0.83 | 0.82 | 0.73 | 0.79 | 0.69 | 0.62 | **0.74** | 0.77 | 0.74 | 0.78 |
| WNN | 0.97 | 0.96 | 0.83 | 0.84 | 0.72 | 0.52 | 0.65 | 0.76 | 0.81 | 0.63 | 0.69 | 0.78 | 0.79 |
| GraSens | **0.99** | **0.97** | **0.91** | **0.89** | **0.89** | 0.58 | 0.68 | 0.83 | 0.79 | 0.55 | 0.75 | **0.93** | **0.86** |



(a)Falling-down without occlusion  (b)Seating without occlusion  (c)Falling-down with partial occlusion  (d)Seating with partial occlusion

(e)Falling-down without occlusion  (f)Seating without occlusion  (g)Falling-down with partial occlusion  (h)Seating with partial occlusion
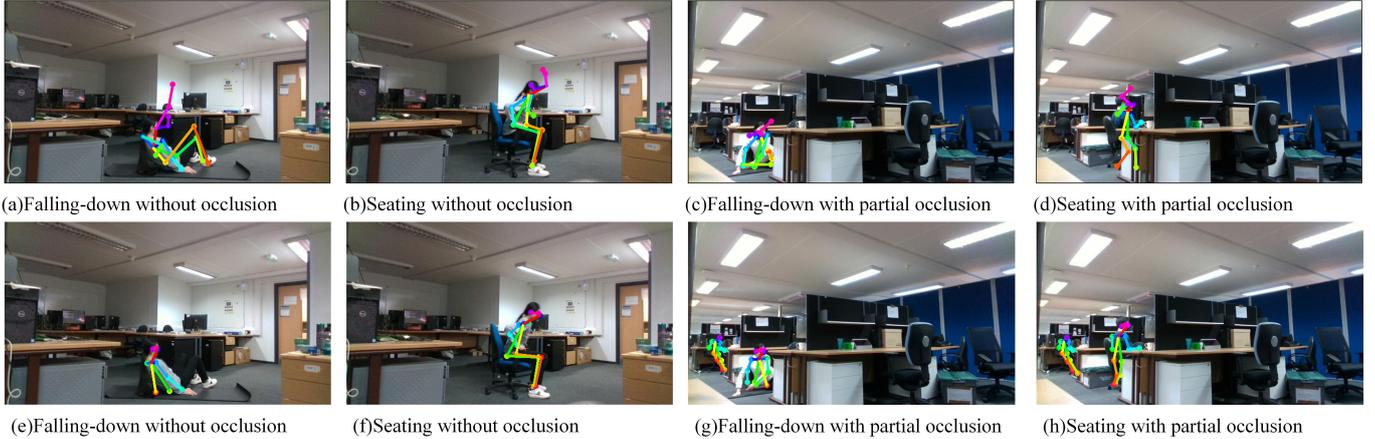
Fig. 4: The skeleton results by WiFi (a)-(d) and by video (e)-(h). In the scene without occlusion as the first two columns show, the skeleton results by WiFi are comparable in seating, and better in self-occlusion cases like falling down than those by video. As for the scene without occlusion in the last two columns, the skeleton results by WiFi are more precise seen in the legs in (d) compared to (h) and have less false detection like the chairs than those by video.

TABLE IV: Classification accuracy of the dataset CSLOS

| Scenes | Methods | no_move | falling | walking | sitting/standing | turning | picking_up | Average |
|---|---|---|---|---|---|---|---|---|
| E1 | SVM [40] | 0.98 | 0.86 | **1.00** | 0.91 | 0.90 | 0.92 | 0.94 |
| | WNN | 0.89 | 0.80 | 0.73 | 0.86 | 0.67 | 0.94 | 0.81 |
| | GraSens | **0.97** | **0.97** | 0.95 | **0.98** | **0.96** | **0.99** | **0.97** |
| E2 | SVM [40] | **0.95** | 0.82 | **0.99** | 0.82 | 0.81 | 0.82 | 0.89 |
| | WNN | 0.84 | 0.78 | 0.75 | 0.83 | 0.69 | 0.84 | 0.79 |
| | GraSens | 0.93 | **0.94** | 0.98 | **0.91** | **0.92** | 0.91 | **0.93** |

learning methods. Compared to WNN and LSTM, GraSens achieved the best performance on fall, sit-down, and stand-up, which means that GraSens can capture the characteristics of rapidly changing motion in time and space. These results demonstrated that the GraSens is able to explore the frequency and temporal continuity inside WiFi signals to enhance the quality of output features scattered in different subcarriers. As for lie-dow, GraSens obtained slightly lower but similar performance with 1% than LSTM. The reason is due to that the signals change fast at the beginning but keep similar after in space. With regard to the action walk which behaved similarly in time and space, the accuracy of GraSens was 8% lower than LSTM. The possible reason is that the spectrum of the signals behaves similarly in time. The results indicated that GraSens is good at sophisticated action recognition but slightly poor at simple actions.

Table III shows the classification accuracy of the dataset HHI. GraSens obtains the most satisfying results by obvi-

ous margins and surpassed the original method E2EDLF. GraSens outperforms the WNN with 7% which confirms the effectiveness of fractal dimension-based self-attention as well as Gabor filtering-based anti-aliasing. Specifically, for the actions of approaching and departing, all of these methods achieved satisfied accuracy over 90%. On the basis of the results of hand-shaking, high five, hugging, and pushing, the proposed GraSens outperformed other algorithms. However, the evaluation of GraSens on kicking, pointing and punching lacked effectiveness. The possible reason is that these actions were single limb linear movements and last shortly in time series sequences thus the input CSI samples contained an amount of the noises included in the ambient environment. GraSens augmented the characteristics of WiFi signals and was inevitably affected by these noises. Overall, the performance of GraSens was moderate, but it was still more convenient to realize action recognition with no requirements for the sophisticated preprocessing than the state-of-art E2EDLF, especially on complex actions in the temporal and frequency domains.

*2) Qualitative Results:* We also show the effectiveness of WiFi and Video data on WVAR. Fig. 2(b) and (c) illustrate that CSI signals are not affected by the occlusion and exhibit similar patterns in the same actions.

Skeleton visualization is further to show the effectiveness of WVAR. Inspired by the work [8], the skeletons derived from Alphapose [43] are used to train the GraSens in LOS

TABLE V: Ablation study of the number of GraSens blocks

| Blocks | approaching | departing | hand_shaking | high five | hugging | kicking_left_leg | kicking_right_leg | pointing_left_hand | pointing_right_hand | punching_left_hand | punching_right_hand | pushing | OA |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\lambda = 4$ | 0.96 | **0.98** | 0.84 | 0.80 | 0.70 | 0.50 | 0.49 | 0.83 | **0.84** | **0.65** | **0.81** | **0.95** | 0.84 |
| $\lambda = 8$ | **0.99** | 0.97 | **0.91** | **0.89** | **0.89** | **0.58** | **0.68** | 0.83 | 0.79 | 0.55 | 0.75 | 0.93 | **0.86** |
| $\lambda = 16$ | 0.96 | 0.96 | 0.84 | 0.83 | 0.77 | 0.52 | 0.64 | 0.81 | 0.80 | 0.53 | 0.59 | 0.91 | 0.82 |

TABLE VI: Ablation study of Gabor filtering-based anti-aliasing mechanism and fractal dimension-based self-attention distilling

| Ablation Study | Methods | approaching | departing | hand_shaking | high five | hugging | kicking_left_leg | kicking_right_leg | pointing_left_hand | pointing_right_hand | punching_left_hand | punching_right_hand | pushing | OA |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Gabor filtering based anti-aliasing mechanism** | **Baseline1** | 0.96 | 0.96 | 0.79 | 0.85 | 0.69 | 0.55 | 0.65 | 0.66 | 0.65 | 0.54 | 0.58 | **0.93** | 0.78 |
| | **Baseline1+Anti-aliasing** | 0.97 | **0.98** | 0.83 | 0.91 | **0.91** | **0.61** | 0.63 | 0.74 | 0.78 | 0.46 | 0.69 | **0.93** | 0.84 |
| | **Baseline1+Gabor** | **1.00** | 0.95 | 0.83 | **0.92** | 0.67 | 0.53 | 0.63 | **0.85** | **0.92** | 0.42 | 0.69 | 0.90 | 0.85 |
| | **GANet** | 0.99 | 0.97 | **0.91** | 0.89 | 0.89 | 0.58 | **0.68** | 0.83 | 0.79 | **0.55** | 0.75 | **0.93** | **0.86** |
| **fractal dimension based self-attention distilling** | **Baseline2** | 0.91 | **0.98** | 0.84 | 0.85 | 0.74 | 0.57 | 0.54 | 0.71 | 0.67 | 0.51 | 0.62 | 0.90 | 0.79 |
| | **Baseline2+FrequencyAttention** | 0.97 | 0.94 | 0.75 | 0.80 | 0.71 | 0.57 | 0.56 | 0.84 | 0.81 | 0.45 | 0.61 | 0.90 | 0.82 |
| | **Baseline2+TemporalAttention** | 0.79 | 1.00 | 0.95 | **0.90** | 0.86 | 0.50 | **0.91** | **0.89** | **0.91** | **0.64** | 0.50 | 0.62 | 0.84 |
| | **GANet** | **0.99** | 0.97 | **0.91** | 0.89 | **0.89** | **0.58** | 0.68 | 0.83 | 0.79 | 0.55 | **0.75** | **0.93** | **0.86** |

conditions. On the basis of the skeletons, the trained GraSens can further generate skeletons in non-line-of-light scenes. Skeleton visualization is further to show the effectiveness of WVAR. As seen in Fig. 4(a)-(d), in the scene without any occlusions, our GraSens yielded robust skeletons in good agreement with the truth images which were close to these of Alphapose. In partially covered situations, GraSens provided the most convincing skeleton results such as seating in Fig. 4(d) compared to Alphapose in Fig. 4(h), with the skeleton boundary being visually close to the raw truth image. This clearly demonstrates that our CSI data on WVAR has a good efficiency in these scenarios.

### C. Ablation Study

In this subsection, we have implemented the experiments to reveal how the different number of GraSens blocks influence the classification accuracy. In addition, we also conducted additional experiments on GraSens with ablation consideration. In this study, we use HHI as the benchmark to test the additional effects of the different number of GraSens blocks as well as self-attention and anti-aliasing mechanisms.

*1) The performance of number of GraSens blocks:* The number of stacked blocks $\lambda$ has a trade-off between the accuracy and efficiency of the proposed GraSens method. To further verify the influence of the number of stacked blocks on performance, we have added an experiment as illustrated in Table V. As shown in Table V, the GraSens achieves the better performance with a growth of 2% when $\lambda = 8$ compared with when $\lambda = 4$. In contrast, when we add the number of blocks to $\lambda = 16$, the classification accuracy decreases by 2%. It is noted that the 16 GraSens blocks network architecture is over-fitting for the training data and generalizes poorly on new testing data. As a result, the classification accuracy decreases on the contrary. According to the results, we choose $\lambda = 8$ as the number of blocks used in our experiments empirically.

*2) The performance of Gabor filtering-based anti-aliasing mechanism:* In this study, we testify to the potential accuracy of our Gabor filtering, anti-aliasing, and Gabor filtering-based anti-aliasing in acquiring "generative" results illustrated in Table VI. Firstly, WNN with the fractal dimension-based self-attention is set as the main pipeline 'baseline1'. For the second, we replace the pooling with an anti-aliasing operation. For the third, the Gabor filtering replaces the first layer of baseline as the Gabor convolution layer. Surprisingly, both anti-aliasing operation and Gabor filtering largely improve the classification accuracy by 8% and 9%, respectively. In addition, the fusion of two operations continues to enhance the performance by 9%. This confirms both the correlation between Gabor filtering and anti-aliasing operation and the importance of the fusion of each other. Thereafter, Gabor filtering-based anti-aliasing further improves the performance, widening the gap with the existing methods.

*3) The performance of fractal dimension-based self-attention distilling:* In the overall results Table VI, we distill frequency and temporal attention separately for self-attention. Firstly, WNN with Gabor filtering-based anti-aliasing is used as the 'baseline2'. Firstly, we add the baseline2 with fractal dimension frequency attention only. As for the second, we add the baseline2 with fractal dimension temporal attention. The fractal dimension-based self-attention determines how the network distributes the contribution of the features. We notice that both the frequency attention and the temporal attention contribute to the improvements of accuracy by 3% and 4%. The integration of both can further refine the accuracy by 7%.

### V. CONCLUSION

In this paper, we identified the inherent limitation of the WiFi signal-based convolution neural networks, with observations that the efficacy of WiFi signals is prone to be influenced by the change in the ambient environment and varies over different sub-carriers. Thereafter, based on their characteristics, we proposed to formulate reliable and robust temporal and frequency shift-invariant representations. We first designed the Gabor filtering based on anti-aliasing to obtain the shift-invariant feature information of actions with the strong auxiliary function. Furthermore, fractal dimension-based frequency and temporal self-attention are proposed to focus on the dominant features scattered in different subcarriers. In addition, we collected synchronous video and WiFi datasets WVAR to simulate the complex visual conditions like

the occlusions scenarios. The ablation study verified that both our Gabor filtering-based anti-aliasing and fractal dimension-based frequency and temporal self-attention are beneficial for the improvement of classification accuracy. Through the experiments on the four most popular datasets, our GraSens achieved a new state-of-the-art with a large margin. We believe it would be a promising future direction to adopt the Gabor filtering-based anti-aliasing and fractal dimension-based attention to the HAR or other related tasks.

## REFERENCES

[1] B. Wu, Z. Ma, S. Poslad, and Y. Li, "WiFi fingerprint based, indoor, location-driven activities of daily living recognition," in *BESC*. IEEE, 2018, pp. 148–151.

[2] F. Wang, J. Han, S. Zhang, X. He, and D. Huang, "CSI-net: Unified human body characterization and action recognition," *arXiv:1810.03064*, 2018.

[3] B. Tan, Q. Chen, K. Chetty, K. Woodbridge, W. Li, and R. Piechocki, "Exploiting WiFi channel state information for residential healthcare informatics," *IEEE Commun. Mag.*, vol. 56, no. 5, pp. 130–137, 2018.

[4] D. B. Lindell, G. Wetzstein, and V. Koltun, "Acoustic non-line-of-sight imaging," in *CVPR*, 2019, pp. 6780–6789.

[5] M. Isogawa, Y. Yuan, M. O'Toole, and K. M. Kitani, "Optical non-line-of-sight physics-based 3D human pose estimation," in *CVPR*, 2020, pp. 7013–7022.

[6] T. Li, L. Fan, M. Zhao, Y. Liu, and D. Katabi, "Making the invisible visible: Action recognition through walls and occlusions," in *ICCV*, 2019, pp. 872–881.

[7] Y. Luo, Y. Li, M. Foshey, W. Shou, P. Sharma, T. Palacios, A. Torralba, and W. Matusik, "Intelligent Carpet: Inferring 3D human pose from tactile signals," in *CVPR*, 2021, pp. 11 255–11 265.

[8] 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 *CVPR*, 2018.

[9] F. Wang, S. Panev, Z. Dai, J. Han, and D. Huang, "Can WiFi estimate person pose?" *arXiv:1904.00277*, 2019.

[10] F. Wang, S. Zhou, S. Panev, J. Han, and D. Huang, "Person-in-WiFi: Fine-grained person perception using WiFi," *arXiv:1904.00276*, 2019.

[11] M. Li, Y. Meng, J. Liu, H. Zhu, X. Liang, Y. Liu, and N. Ruan, "When CSI meets public WiFi: Inferring your mobile phone password via WiFi signals," in *CCS*. ACM, 2016, pp. 1068–1079.

[12] K. Qian, C. Wu, Z. Zhou, Y. Zheng, Z. Yang, and Y. Liu, "Inferring motion direction using commodity Wi-Fi for interactive exergames," in *CHI*. ACM, 2017, pp. 1961–1972.

[13] D. Zhang and L. M. Ni, "Dynamic clustering for tracking multiple transceiver-free objects," in *PerCom*. IEEE, 2009, pp. 1–8.

[14] F. Li, M. Valero, H. Shahriar, R. A. Khan, and S. I. Ahamed, "Wi-COVID: A COVID-19 symptom detection and patient monitoring framework using WiFi," *Smart Health*, vol. 19, p. 100147, 2021.

[15] F. Adib and D. Katabi, "See through walls with WiFi!" in *ACM SIGCOMM Conf. SIGCOMM.*, 2013, pp. 75–86.

[16] V.-M. Khong and T.-H. Tran, "Improving human action recognition with two-stream 3D convolutional neural network," in *MAPR*. IEEE, 2018, pp. 1–6.

[17] A. Sharif Razavian, H. Azizpour, J. Sullivan, and S. Carlsson, "CNN features off-the-shelf: an astounding baseline for recognition," in *CVPR*, 2014, pp. 806–813.

[18] G. W. Taylor, R. Fergus, Y. LeCun, and C. Bregler, "Convolutional learning of spatio-temporal features." Springer, 2010, pp. 140–153.

[19] I. Pitas, *Digital image processing algorithms and applications.* John Wiley & Sons, 2000.

[20] X. Zou, F. Xiao, Z. Yu, and Y. J. Lee, "Delving deeper into anti-aliasing in convnets," *arXiv:2008.09604*, 2020.

[21] R. Zhang, "Making convolutional networks shift-invariant again," in *ICML*. PMLR, 2019, pp. 7324–7334.

[22] L. Pishchulin, E. Insafutdinov, S. Tang, B. Andres, M. Andriluka, P. V. Gehler, and B. Schiele, "Deepcut: Joint subset partition and labeling for multi person pose estimation," in *CVPR*, 2016, pp. 4929–4937.

[23] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *NIPS*, vol. 28, 2015, pp. 91–99.

[24] C. Feichtenhofer, H. Fan, J. Malik, and K. He, "SlowFast networks for video recognition," in *ICCV*, 2019, pp. 6202–6211.

[25] J. Gao, T. Zhang, and C. Xu, "I know the relationships: Zero-shot action recognition via two-stream graph convolutional networks and knowledge graphs," in *AAAI*, vol. 33, no. 01, 2019, pp. 8303–8311.

[26] A. Garcia-Garcia, F. Gomez-Donoso, J. Garcia-Rodriguez, S. Orts-Escolano, M. Cazorla, and J. Azorin-Lopez, "PointNet: A 3D convolutional neural network for real-time object class recognition," in *IJCNN*. IEEE, 2016, pp. 1578–1584.

[27] W. Wang, A. X. Liu, and M. Shahzad, "Gait recognition using WiFi signals," in *ACM Int. Joint Conf. Pervasive Ubiquitous Comput.*, 2016, pp. 363–373.

[28] 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, 2016.

[29] P. M. Holl and F. Reinhard, "Holography of Wi-Fi radiation," *Phys. Rev. Lett.*, vol. 118, no. 18, p. 183901, 2017.

[30] R. Alazrai, M. Hababeh, A. Baha'A, M. Z. Ali, and M. I. Daoud, "An end-to-end deep learning framework for recognizing human-to-human interactions using Wi-Fi signals," *IEEE Access*, vol. 8, pp. 197 695–197 710, 2020.

[31] T. K. Ho, "Random decision forests," in *ICDAR*, vol. 1. IEEE, 1995, pp. 278–282.

[32] S. R. Eddy, "What is a hidden Markov model?" *Nat. Biotechnol.*, vol. 22, no. 10, pp. 1315–1316, 2004.

[33] S. Yousefi, H. Narui, S. Dayal, S. Ermon, and S. Valaee, "A survey on behavior recognition using WiFi channel state information," *IEEE Commun. Mag.*, vol. 55, no. 10, pp. 98–104, 2017.

[34] S. Luan, C. Chen, B. Zhang, J. Han, and J. Liu, "Gabor convolutional networks," *IEEE Trans. Image Process.*, vol. 27, no. 9, pp. 4357–4366, 2018.

[35] A. Alekseev and A. Bobe, "GaborNet: Gabor filters with learnable parameters in deep convolutional neural network," in *EnT*. IEEE, 2019, pp. 1–4.

[36] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "Cbam: Convolutional block attention module," in *ECCV*, 2018, pp. 3–19.

[37] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *CVPR*, 2015, pp. 1–9.

[38] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.

[39] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, "SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and¡ 0.5 mb model size," *arXiv:1602.07360*, 2016.

[40] A. Baha'A, M. M. Almazari, R. Alazrai, and M. I. Daoud, "Exploiting Wi-Fi signals for human activity recognition," in *ICICS*. IEEE, 2021, pp. 245–250.

[41] R. Alazrai, A. Awad, A. Baha'A, M. Hababeh, and M. I. Daoud, "A dataset for Wi-Fi-based human-to-human interaction recognition," *Data Brief*, vol. 31, p. 105668, 2020.

[42] A. Baha'A, M. M. Almazari, R. Alazrai, and M. I. Daoud, "A dataset for wi-fi-based human activity recognition in line-of-sight and non-line-of-sight indoor environments," *Data in Brief*, vol. 33, p. 106534, 2020.

[43] H.-S. Fang, S. Xie, Y.-W. Tai, and C. Lu, "RMPE: Regional multi-person pose estimation," in *ICCV*, 2017, pp. 2334–2343.