

RFGAN: RF-Based Human Synthesis

Cong Yu, Zhi Wu, Dongheng Zhang, Zhi Lu, Yang Hu,
and Yan Chen, *Senior Member, IEEE*

Abstract—This paper demonstrates human synthesis based on the Radio Frequency (RF) signals, which leverages the fact that RF signals can record human movements with the signal reflections off the human body. Different from existing RF sensing works that can only perceive humans roughly, this paper aims to generate fine-grained optical human images by introducing a novel cross-modal RFGAN model. Specifically, we first build a radio system equipped with horizontal and vertical antenna arrays to transceive RF signals. Since the reflected RF signals are processed as obscure signal projection heatmaps on the horizontal and vertical planes, we design a RF-Extractor with RNN in RFGAN for RF heatmap encoding and combining to obtain the human activity information. Then we inject the information extracted by the RF-Extractor and RNN as the condition into GAN using the proposed RF-based adaptive normalizations. Finally, we train the whole model in an end-to-end manner. To evaluate our proposed model, we create two cross-modal datasets (*RF-Walk & RF-Activity*) that contain thousands of optical human activity frames and corresponding RF signals. Experimental results show that the RFGAN can generate target human activity frames using RF signals. To the best of our knowledge, this is the first work to generate optical images based on RF signals.

Index Terms—RF Sensing, Human Synthesis, GAN.

I. INTRODUCTION

VARIOUS recent works have built Radio Frequency (RF) sensing systems to perceive and understand the activities of humans. Compared with alternative sensing methods, RF sensing has improved usability due to the characteristics of RF signals, for example, the RF signals can work in all-day and all-weather scenarios, the sensing is contactless etc.. Existing RF-based human sensing works mainly include human position tracking [1]–[9], human speed estimation [10]–[12], human keypoint prediction [13]–[16], and gesture recognition [17]. However, these works can only perceive humans roughly and the sensing results usually lack fine details and are not as intuitive as optical sensing results.

In recent years, Generative Adversarial Networks (GAN) [18] have achieved promising results in modeling complex multimodal data and synthesizing realistic images. Furthermore, to generate meaningful images that meet actual requirements, many conditional GAN models have been proposed to control the generated results. Researchers have explored various kinds of conditions, e.g., category labels [19], text descriptions [20]–[22], and images [23]–[31]. From technology perspective, most existing GAN models require the conditions to be able to guide the GAN model

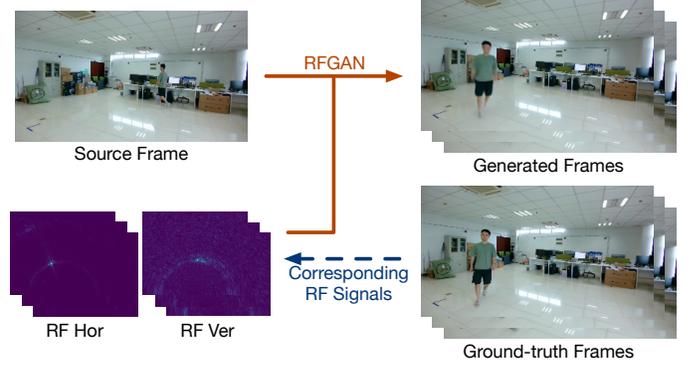


Fig. 1. With a source frame as reference, the RFGAN model can synthesize the target human activities based on the RF signals.

explicitly [19], [23], [25], [26], [31] or can be transformed to conditional variables for GAN using an existing pre-trained model [20]–[22].

In this paper, we propose a solution to overcome the limitation of existing RF-based human activity sensing by making the results more visually intuitive, which is valuable in practice. We leverage the power of GAN models to generate photo-realistic sensing results from RF signals. Specifically, a photograph of the people in the scene is provided so that the GAN model has sufficient information about the visual appearance of the people and the environment of the scene. We use millimeter wave (mmWave) radars to build our radio system, which is equipped with two antenna arrays, horizontal and vertical ones, to obtain the RF signals that reflect off the human body. We process the horizontal and vertical RF signal reflections to horizontal and vertical RF heatmaps, which record the activity information of the human. The RF signal is a new kind of conditional data for GAN models. Due to the characteristics of RF signals, the resolution of the horizontal and the vertical RF heatmaps are relatively low. Besides, their spatial structures are quite different from optical images. Therefore, to utilize the RF signals as the conditional data to guide the GAN model, some challenges need to be addressed: firstly, we need to train the RF conditioning encoding network without supervision labels to obtain the desirable human activity information; secondly, the information from the horizontal and vertical RF heatmaps need to be fused to characterize the overall human activity; thirdly, the fused information needs to be injected into the GAN model properly.

To tackle the above challenges, we design dual RF-Extractors and RNNs in the RFGAN model, one in the generative part and the other in the discriminative part, and we train them by adversarial learning. Two CNN encoders

C. Yu and Z. Lu are with University of Electronic Science and Technology of China. E-mail: congyu@std.uestc.edu.cn, zhilu@std.uestc.edu.cn

Z. Wu, D. Zhang, Y. Hu and Y. Chen are with University of Science and Technology of China, and Y. Chen is the corresponding author. E-mail: wzwyx@mail.ustc.edu.cn, dongheng@ustc.edu.cn, eeyhu@ustc.edu.cn, eecyan@ustc.edu.cn

in the RF-Extractor are used to extract features from the horizontal and vertical RF heatmaps, respectively. Then a novel fusion operation is designed to fuse the information by building relationships between the extracted features. To inject the fused information into the GAN model, inspired by [26], [32], [33], we propose to modify the distributions of the latent features in GAN by using a RF-based adaptive normalization. Furthermore, we create two cross-modal datasets (*RF-Walk* & *RF-Activity*) that consist of optical human activity frames and corresponding RF signals to train and test our proposed RFGAN model. The experimental results show that the RFGAN can generate better human results than alternative methods.

Since the RF signals do not rely on visible lights and can traverse occlusions, our proposed RFGAN model can also work when lights dim or the human is occluded by barriers. For example, when the environment is favorable, we capture a human frame as the source. Our radio system can record the RF signal reflections when the illumination becomes bad or the human is in occlusions. The proposed RFGAN model can synthesize human activities based on these collected multimodal data.

Therefore, the main contributions of this paper can be summarized as follows:

1. We propose a novel RFGAN model to enable RF-based human synthesis. To the best of our knowledge, this is the first work to generate human images from the mmWave radar signals. There are many potential applications that can be derived from this task, e.g., fine-grained human perception and all-day monitoring systems in the smart home.
2. Technically, for the new kind of conditional data, i.e., the RF signals, we propose to train the RF conditioning encoding network, i.e., the RF-Extractor and RNN, by adversarial learning. Then we design a novel fusion operation to fuse the horizontal and vertical RF information, which is an effective approach for overall human activity sensing from the two-dimensional RF heatmaps. Due to the spatial structure difference, we propose to use the RF-based adaptive normalizations to inject the fused information into the GAN model.
3. We create two cross-modal datasets, i.e., *RF-Walk* and *RF-Activity*, which contain thousands of optical human activity frames and corresponding RF signals. The datasets will be released to public.

II. RELATED WORK

Conditional GAN Many research works have shown that GAN [18] has the capability of generating realistic images based on the given conditional data. For example, [19] utilize category labels to generate target digit images. Some works [23]–[31] introduce the GAN-based image-to-image translation frameworks. For some more complex conditional data, e.g., text data, [20]–[22] use existing pre-trained models to transform the text into conditioning variables for GAN. To employ these conditions in the networks, some works [26], [32], [33] find that utilizing the conditional normalization in the hidden layers can contribute to generating target images.

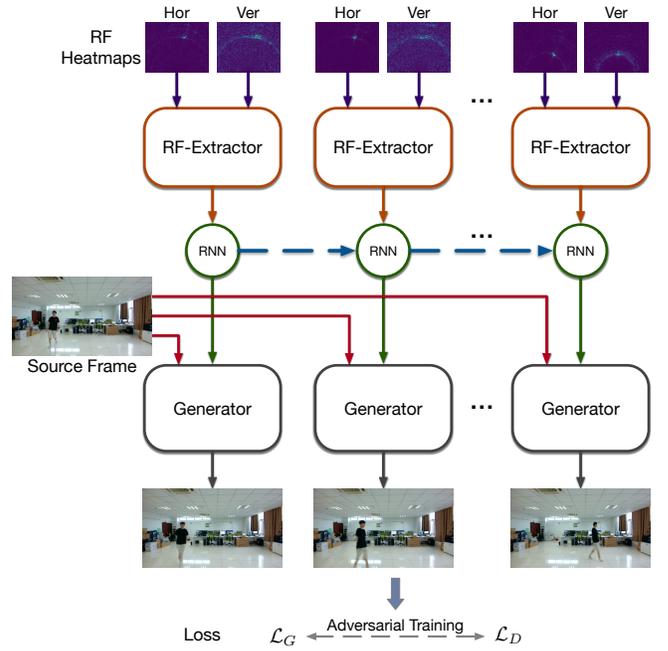


Fig. 2. The architecture of the RFGAN model for generating sequential human activity frames.

In our case, we take RF signals as the condition to guide the image synthesis, which is a new cross-modal conditional data that has obscure guidance for GAN and has no existing pre-trained model for conditioning encoding.

RF-Based Human Perception Recent years have witnessed much interest in using RF signals to enable various human perception tasks [34], including indoor localization and tracking [1]–[5], [7]–[9], human speed estimation or human movement detection [10]–[12], [17], human identification [35]–[37], and human vital signs inference [38]–[42]. Besides the above signal-processing-based methods, approaches based on deep learning are also utilized to handle radio human perception. For example, [43] combines convolutional and recurrent neural networks to learn sleep stages from radio signals. [13], [14] propose to predict the 2D/3D human keypoints based on RF signals by building a teacher-student network model. In this paper, we propose to use RF signals for human image synthesis by combining conditional GAN models.

Sequence Modeling Recurrent neural networks, such as vanilla RNN, GRU and LSTM, have been widely used for processing sequential data, such as text and speech. They have also been successfully applied to model the temporal dependencies in videos for various vision problems, such as video classification [44], action recognition [45]–[49], object segmentation [50], video prediction [51], etc. In this work, considering that the RF signals are sequential data and the RF heatmaps are the samples at different moments, we utilize recurrent neural networks as the backbone of our model to perceive human activities from RF signals and synthesize corresponding optical images.

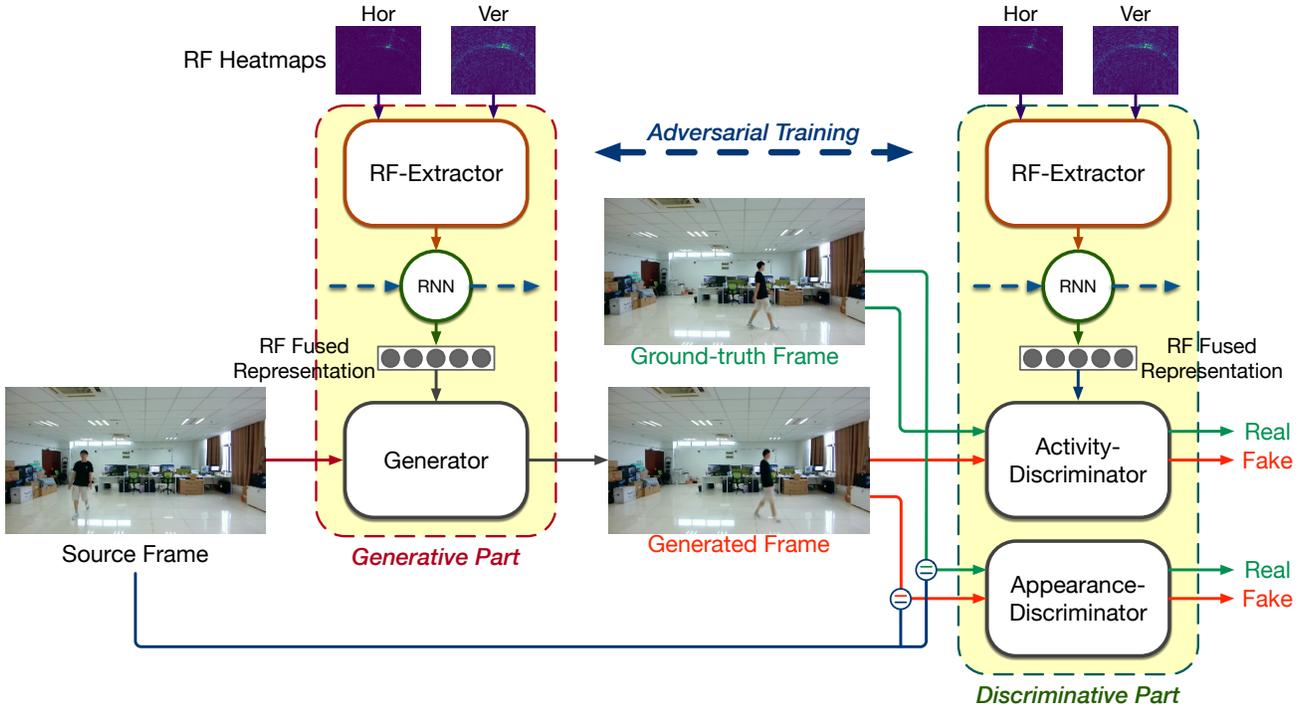


Fig. 3. The training framework of RFGAN at one moment. It consists of a generative part and a discriminative part. The whole model is trained by adversarial learning in an end-to-end manner.

III. PRELIMINARY

Our method relies on transmitting RF signals and receiving the reflections. We adopt Frequency Modulated Continuous Wave (FMCW) and linear antenna arrays for signal transceiving. Inspired by [13], our radio system is equipped with two antenna arrays: horizontal and vertical ones, which are utilized to acquire the signal projections on the plane parallel to the ground and the plane perpendicular to the ground, respectively. Hence, the RF data is composed of both horizontal and vertical heatmaps.

Compared with camera-based visual data, RF signals have some different characteristics. Firstly, RF signals have much lower resolution. The resolution is determined by the bandwidth of the signal and the aperture of the antenna array [52]. In our system, the depth resolution is about 7.5cm, and the angular resolution is about 1.3 degrees. Secondly, the RF signals suffer from severe multi-path propagation in an indoor environment [8], which introduces severe interference in the received signals. Thirdly, the RF signals have different representations of the scene compared with the camera, i.e., horizontal and vertical projections.

IV. RFGAN

The RFGAN model aims to generate sequential human activity frames using a sequence of RF heatmaps (horizontal & vertical) and a source frame. To extract and combine the human activity information from the horizontal and vertical RF heatmaps, we design a RF-Extractor, which is built with a sequence model, i.e., the Recurrent Neural Network (RNN), to process the RF sequence. To generate optical human activity

frames, we utilize the Generative Adversarial Network (GAN) as the main technological approach in our model, where the source frame is fed as the input layer and the information extracted from RF heatmaps is the condition of GAN.

The architecture of RFGAN model is shown in Figure 2. The RNN is the backbone of the model, which is designed for sequence data processing and generation. The RF-Extractor and the Generator are plugged into both sides of the RNN to process RF heatmaps and generate human frames. In the following subsections, we first introduce the training framework of the model and then discuss the network structures of the RF-Extractor, the RNN, and the RF-based Generator and Discriminators in detail. Finally, we describe the loss functions used to train the whole model.

A. Training Framework

The proposed human synthesis model aims to generate sequential human frames from a source frame and the corresponding sequential RF heatmaps. Figure 3 shows the adversarial training framework of the human synthesis model at one moment, which consists of a generative part and a discriminative part. The generative part contains a RF-Extractor, a RNN, and a Generator. The RF-Extractor and RNN extract the human position and posture information from the corresponding RF heatmaps and represent it as a RF fused representation. For the Generator, the source frame is fed as the input layer, and the extracted RF fused representation controls the network through normalization at the convolution layers. The output is the generated human frame. There are

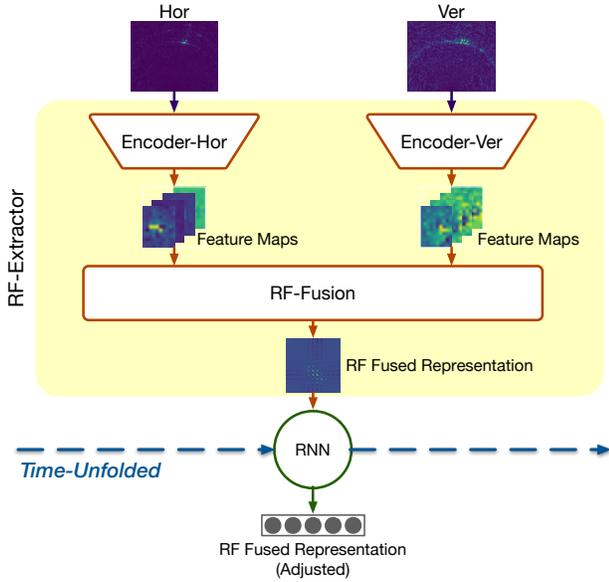


Fig. 4. The structure of RF-Extractor. It consists of two CNN encoders, a fusion operation, and a RNN.

two Discriminators in the discriminative part. The Activity-Discriminator is designed to ensure that the human position and posture in the generated frame are consistent with the RF signal. It takes the generated frame as the input layer. The RF fused representation extracted by the RF-Extractor and RNN in the discriminative part is used as condition of this Discriminator. The Appearance-Discriminator ensures that the generated frame maintains the same visual information, such as human appearance, with the source frame, thus the generated frame is concatenated with the source frame at the input layer.

Note that the RF-Extractors and RNNs in the generative part and the discriminative part have the same network structures and RF inputs, but they do not share the network parameters. In the previous GAN literature, the conditional variables that fed into the Generator and the Discriminator are obtained by the same network model, mainly due to the existence of the pre-trained model to enable the desirable conditioning encoding. However, in our task, there is no existing pre-trained model for RF encoding. Thus, we propose dual RF-Extractors and RNNs under an adversarial training framework to learn to transform the RF heatmaps, one for the generation task and the other for the discrimination task. Specifically, the RF-Extractor and RNN in the generative part update with the Generator, whereas the RF-Extractor and RNN in the discriminative part update with the Activity-Discriminator. The update process is adversarial training and the whole model is trained in an end-to-end manner.

B. RF-Extractor & RNN

The horizontal RF heatmaps and the vertical RF heatmaps record human activities from different viewpoints and each of them only contains partial human activity information, i.e., the horizontal RF heatmap is a projection of the signal reflections

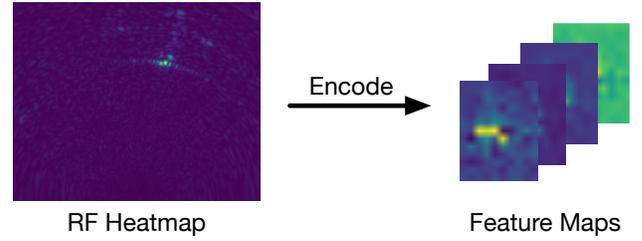


Fig. 5. The RF heatmap and the corresponding feature maps.

on a plane parallel to the ground, which leads to the loss of the human height information, whereas the vertical heatmap is a projection of the reflected signals on a plane perpendicular to the ground and the human width information is missed. Thus, it is a challenge that how to extract and combine the horizontal and vertical RF information to get the whole human activity information.

In our proposed network structure, as shown in Figure 4, we first use two standard CNN encoders to transform the horizontal and vertical RF heatmaps into feature maps, respectively. The original RF heatmaps record the reflected signals throughout the whole room. After the differential operation along the time, only the signals introduced by the moving human are retained. As shown in left part of Figure 5, we can find the signal reflections from the moving human (bright area) only occupy a very small area of the RF heatmap. Therefore, we use an encoder that consists of several convolution layers to reduce the RF heatmap size and focus on the bright area. Since the values of signal reflections from no human areas (dark areas) are very small and close to 0, the convolution results, which are denoted as feature maps (shown in right part of Figure 5), can capture the human posture information from the bright area and the human position information from the location of the bright area.

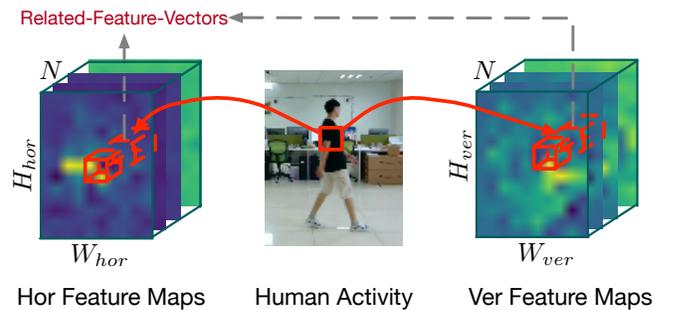


Fig. 6. The horizontal and vertical feature maps contain the human activity information on the horizontal and vertical plane, respectively. The red cuboids are the related-feature-vectors, which contain the activity information introduced by the same human body part.

After encoding RF heatmaps, a fusion operation (RF-Fusion) is proposed to combine the horizontal and vertical feature maps into a fused representation. As shown in Figure 6, the horizontal feature maps can be represented as a $H_{hor} \times W_{hor} \times N$ tensor, which uses $H_{hor} \times W_{hor}$ feature vectors on a horizontal plane to record the human activity

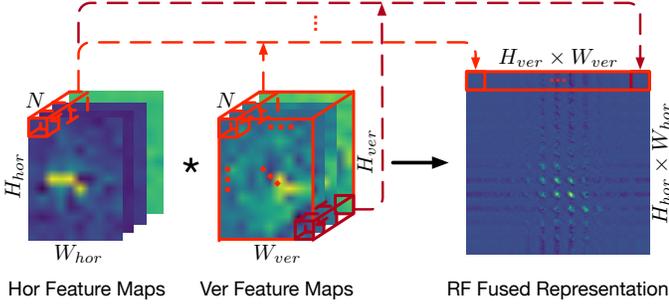


Fig. 7. The RF-Fusion operation.

information, and each feature vector is N dimensional. For the vertical feature maps, $H_{ver} \times W_{ver}$ feature vectors are used on a vertical plane to record the human activity information, and each feature vector is also N dimensional. We refer to the feature vectors in the horizontal and vertical feature maps as related-feature-vectors if they record the activity information introduced by the same human body part (see Figure 6). Combining these related-feature-vectors can help characterize the overall human activity. However, it is difficult to find the one-to-one correspondence between them directly.

To address this problem and bridge these related-feature-vectors, we define RF-Fusion as follows: for each feature vector in the horizontal feature maps, the dot product is applied with every feature vector in the vertical feature maps, and the results are denoted as a RF fused representation. For example, as shown in Figure 7, the dot products between the first horizontal feature vector and every vertical feature vector generate $H_{ver} \times W_{ver}$ values, which are the first row of the RF fused representation. In such a way, we can obtain the RF fused representation as follows:

$$R(i, j) = \frac{H(i)V(j)^T}{\sqrt{N}}, \quad (1)$$

$$i \in [0, H_{hor} \times W_{hor}), j \in [0, H_{ver} \times W_{ver}),$$

where $R(i, j)$ is the value at the point (i, j) in the RF fused representation, $H(i)$ and $V(j)$ refer to the feature vector in the horizontal feature maps and the feature vector in the vertical feature maps. The denominator \sqrt{N} is to scale the values.

Why does RF-Fusion work? The traditional feature map fusion approach is to concatenate feature maps along the channel directly, which is effective when feature maps have the same spatial structure, i.e., the feature vectors in the different feature maps are aligned and can be combined by concatenating. However, in the RF fusion step, for a given feature vector in the horizontal feature maps, we do not know which feature vector in the vertical feature maps is related to it. Thus our proposed RF-Fusion builds relationships between every feature vector in the horizontal feature maps and every feature vector in the vertical feature maps. For the learned RF-Extractor, the related-feature-vectors in the horizontal and vertical feature maps are supposed to be highly correlated and lead to generating large values in the RF fused representation, which are shown as bright points. Therefore, the distribution

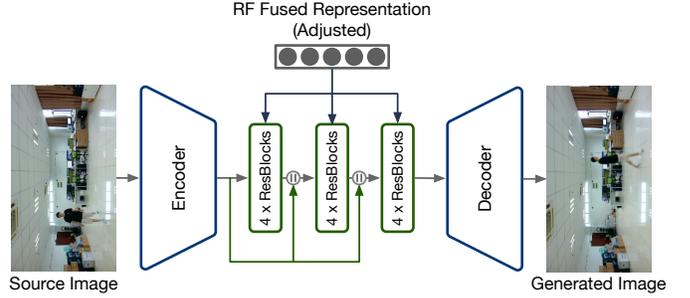


Fig. 8. The structure of Generator.

and values of these bright points can characterize the overall human activity.

Finally, the RF fused representations are fed into the RNN to get adjustments. We propose this procedure based on the following fact: human activities, such as arm swing, leg raising, etc., are generally continuous, thus the RF fused representation that contains the human position and posture information at a certain moment is interrelated with several preceding and subsequent RF fused representations. The RNN model can adjust the current representation by considering its neighbors, and the adjusted results would contain more smooth and more accurate human activity information. In our model, we use a three-layer BiLSTM as the proposed RNN, and the hidden states in the last layer are adopted as the adjusted RF fused representations.

C. RF-Based Generator & Discriminators

The Generator in our model consists of an encoder, several residual blocks, and a decoder (see Figure 8). The encoder and the decoder contain the same numbers of convolutional and deconvolutional layers. The residual blocks are divided into several groups and each group has the same numbers of blocks. The feature extracted from the source frame is concatenated with several feature maps in the residual blocks to maintain the appearance information. For the Activity-Discriminator and the Appearance-Discriminator, we use the network structures inspired by PatchGAN [23]. They both consist of convolutional layers, where the first layer does not use the normalization and the last layer is only a convolution to produce a 1-dimensional output.

Specifically, to enable the RF-based condition setting in RF-GAN, we propose a RF-based adaptive instance normalization (RF-InNorm) in the hidden layers of Generator and Activity-Discriminator, which injects the RF fused representation by modifying the feature distribution. The RF-InNorm is defined as

$$\text{RF-InNorm}(\mathbf{f}^n) = F_\gamma^n(\mathbf{h}) \cdot \frac{\mathbf{f}^n - \boldsymbol{\mu}^n}{\boldsymbol{\sigma}^n} + F_\beta^n(\mathbf{h}), \quad (2)$$

where \mathbf{f}^n is the feature map of the n -th layer in the Generator or Activity-Discriminator, $\boldsymbol{\mu}^n$ and $\boldsymbol{\sigma}^n$ are the mean and standard deviation of the feature map. \mathbf{h} refers to the RF fused representation, $F_\gamma^n(\cdot)$ and $F_\beta^n(\cdot)$ are the learned nonlinear functions, which specialize \mathbf{h} to RF-based modulation parameters. Therefore, the feature map \mathbf{f}^n is first normalized and

then scaled and biased by $F_\gamma^n(\mathbf{h})$ and $F_\beta^n(\mathbf{h})$ to incorporate the RF fused representation condition.

For the Appearance-Discriminator, the source frame condition is concatenated with the input and fed into the network.

D. Loss Functions

The training process of the RFGAN is a two-player minimax game between the generative part and the discriminative part. For the discriminative part, we set the loss for the Activity-Discriminator and the RF-Extractor and RNN as:

$$\mathcal{L}^{act} = \mathcal{L}_{LSD}^{act} + \lambda \mathcal{L}_{GP}^{act}, \quad (3)$$

where \mathcal{L}_{LSD}^{act} is the adversarial loss inspired by LSGAN [53], \mathcal{L}_{GP}^{act} is the gradient regularization term that penalizes the discriminator gradients only on the true data to stabilize the training process [54], which can be calculated by

$$\begin{aligned} \mathcal{L}_{LSD}^{act} = & \mathbb{E}_{\mathbf{x}_r \sim \mathbb{P}} [(D_{act}(\mathbf{x}_r | E_{dis}(\mathbf{r}_h, \mathbf{r}_v)) - 1)^2] + \\ & \mathbb{E}_{\mathbf{x}_f \sim \mathbb{Q}} [(D_{act}(\mathbf{x}_f | E_{dis}(\mathbf{r}_h, \mathbf{r}_v)) - 0)^2], \end{aligned} \quad (4)$$

and

$$\mathcal{L}_{GP}^{act} = \mathbb{E}_{\mathbf{x}_r \sim \mathbb{P}} [\|\nabla D_{act}(\mathbf{x}_r | E_{dis}(\mathbf{r}_h, \mathbf{r}_v))\|_2^2], \quad (5)$$

where D_{act} is the Activity-Discriminator, E_{dis} is the RF-Extractor and RNN in the discriminative part, \mathbf{x}_r and \mathbf{x}_f are the ground-truth and the generated human frame, respectively, \mathbf{r}_h and \mathbf{r}_v refer to the horizontal RF heatmap and the vertical RF heatmap, respectively.

For the Appearance-Discriminator, the loss function is similar to the loss for the Activity-Discriminator and the RF-Extractor and RNN:

$$\mathcal{L}^{app} = \mathcal{L}_{LSD}^{app} + \lambda \mathcal{L}_{GP}^{app}, \quad (6)$$

the \mathcal{L}_{LSD}^{app} and \mathcal{L}_{GP}^{app} are calculated by

$$\begin{aligned} \mathcal{L}_{LSD}^{app} = & \mathbb{E}_{\mathbf{x}_r \sim \mathbb{P}} [(D_{app}(\mathbf{x}_r | \mathbf{x}_s) - 1)^2] + \\ & \mathbb{E}_{\mathbf{x}_f \sim \mathbb{Q}} [(D_{app}(\mathbf{x}_f | \mathbf{x}_s) - 0)^2], \end{aligned} \quad (7)$$

and

$$\mathcal{L}_{GP}^{app} = \mathbb{E}_{\mathbf{x}_r \sim \mathbb{P}} [\|\nabla D_{app}(\mathbf{x}_r | \mathbf{x}_s)\|_2^2], \quad (8)$$

where \mathbf{x}_s is the source frame.

Therefore, the final loss function of the discriminative part is

$$\mathcal{L}_D = \mathcal{L}^{act} + \mathcal{L}^{app}. \quad (9)$$

For the generative part, the loss function is

$$\mathcal{L}_G = \mathcal{L}_{LSG} + \alpha \mathcal{L}_{IMG} + \beta \mathcal{L}_{FEA}, \quad (10)$$

where \mathcal{L}_{LSG} is the corresponding adversarial loss, \mathcal{L}_{IMG} and \mathcal{L}_{FEA} are designed for synthesizing images with better visual quality, which push the generated images towards the ground-truth images in the image space and the feature space. They are calculated by:

$$\begin{aligned} \mathcal{L}_{LSG} = & \mathbb{E}_{\mathbf{x}_f \sim \mathbb{Q}} [(D_{act}(\mathbf{x}_f | E_{gen}(\mathbf{r}_h, \mathbf{r}_v)) - 1)^2] + \\ & \mathbb{E}_{\mathbf{x}_f \sim \mathbb{Q}} [(D_{app}(\mathbf{x}_f | \mathbf{x}_s) - 1)^2], \end{aligned} \quad (11)$$

Algorithm 1 Training algorithm for RFGAN.

Set: The batch size m is 2, the hyperparameters

$\lambda = \alpha = \beta = 10.0$, the learning rate η is 0.0002.

Initialize: Initial $\Phi_{E_{gen}}$ for the RF-Extractor and RNN in generative part, initial $\Phi_{E_{dis}}$ for the RF-Extractor and RNN in discriminative part, initial Φ_G for the Generator, initial $\Phi_{D_{act}}$ for the Activity-Discriminator, and initial $\Phi_{D_{app}}$ for the Appearance-Discriminator.

```

1: while  $\Phi_{E_{gen}}, \Phi_G$  has not converged do
2:   Sample a batch of  $\{\mathbf{r}_h, \mathbf{r}_v, \mathbf{x}_s, \mathbf{x}_r\}$  from the dataset
3:   Update  $\Phi_{E_{dis}}, \Phi_{D_{act}}, \Phi_{D_{app}}$  using Adam with:
4:      $\Phi_{E_{dis}} \leftarrow \Phi_{E_{dis}} - \eta \frac{1}{m} \nabla \Phi_{E_{dis}} \sum_{i=1}^m \mathcal{L}^{act}$ 
5:      $\Phi_{D_{act}} \leftarrow \Phi_{D_{act}} - \eta \frac{1}{m} \nabla \Phi_{D_{act}} \sum_{i=1}^m \mathcal{L}^{act}$ 
6:      $\Phi_{D_{app}} \leftarrow \Phi_{D_{app}} - \eta \frac{1}{m} \nabla \Phi_{D_{app}} \sum_{i=1}^m \mathcal{L}^{app}$ 
7:   Update  $\Phi_{E_{gen}}, \Phi_G$  using Adam with:
8:      $\Phi_{E_{gen}} \leftarrow \Phi_{E_{gen}} - \eta \frac{1}{m} \nabla \Phi_{E_{gen}} \sum_{i=1}^m \mathcal{L}_G$ 
9:      $\Phi_G \leftarrow \Phi_G - \eta \frac{1}{m} \nabla \Phi_G \sum_{i=1}^m \mathcal{L}_G$ 
10: end while

```

and

$$\mathcal{L}_{IMG} = \mathbb{E}_{\mathbf{x}_f \sim \mathbb{Q}, \mathbf{x}_r \sim \mathbb{P}} \|\mathbf{x}_f - \mathbf{x}_r\|_1, \quad (12)$$

$$\begin{aligned} \mathcal{L}_{FEA} = & \sum_i^K \mathbb{E}_{\mathbf{x}_f \sim \mathbb{Q}, \mathbf{x}_r \sim \mathbb{P}} \|\mathbf{f}_{\mathbf{x}_f}^{i,act} - \mathbf{f}_{\mathbf{x}_r}^{i,act}\|_1 + \\ & \sum_i^K \mathbb{E}_{\mathbf{x}_f \sim \mathbb{Q}, \mathbf{x}_r \sim \mathbb{P}} \|\mathbf{f}_{\mathbf{x}_f}^{i,app} - \mathbf{f}_{\mathbf{x}_r}^{i,app}\|_1, \end{aligned} \quad (13)$$

where E_{gen} is the RF-Extractor and RNN in the generative part, $\mathbf{f}_{\mathbf{x}}^{i,act}$ refers to the feature map of \mathbf{x} at layer i in the Activity-Discriminator, $\mathbf{f}_{\mathbf{x}}^{i,app}$ refers to the feature map of \mathbf{x} at layer i in the Appearance-Discriminator, and K is the total number of layers.

The whole training procedure is described in Algorithm 1.

V. EXPERIMENTS

A. Implementation

Data We collected the RF signal reflections at 20Hz from our mmWave radar system, i.e., the horizontal and vertical antenna arrays generate 20 pairs of heatmaps per second. To obtain the optical human images, we attach an RGB camera with the mmWave radar system to record videos at 10 FPS. In order to reduce the coupling between the human and the environment, we collected the data under 9 indoor scenes. There were 6 volunteers involved in the data collection and each volunteer wears multiple dresses.

In total, we create two types of RF-Vision datasets, i.e., *RF-Walk* and *RF-Activity*. For *RF-Walk*, it contains 67,860 human random walking frames and 135,720 pairs of corresponding RF heatmaps. We use 54,525 frames of human walking images and 109,050 pairs of RF heatmaps for training and the rest for testing. For *RF-Activity*, it contains 68,680 human daily activity (e.g., stand, walk, squat, sit, etc.) frames and 137,360 pairs of corresponding RF heatmaps. We use 55,225 frames of human activity images and 110,450 pairs of RF heatmaps for training and the rest for testing. Each human activity frame is

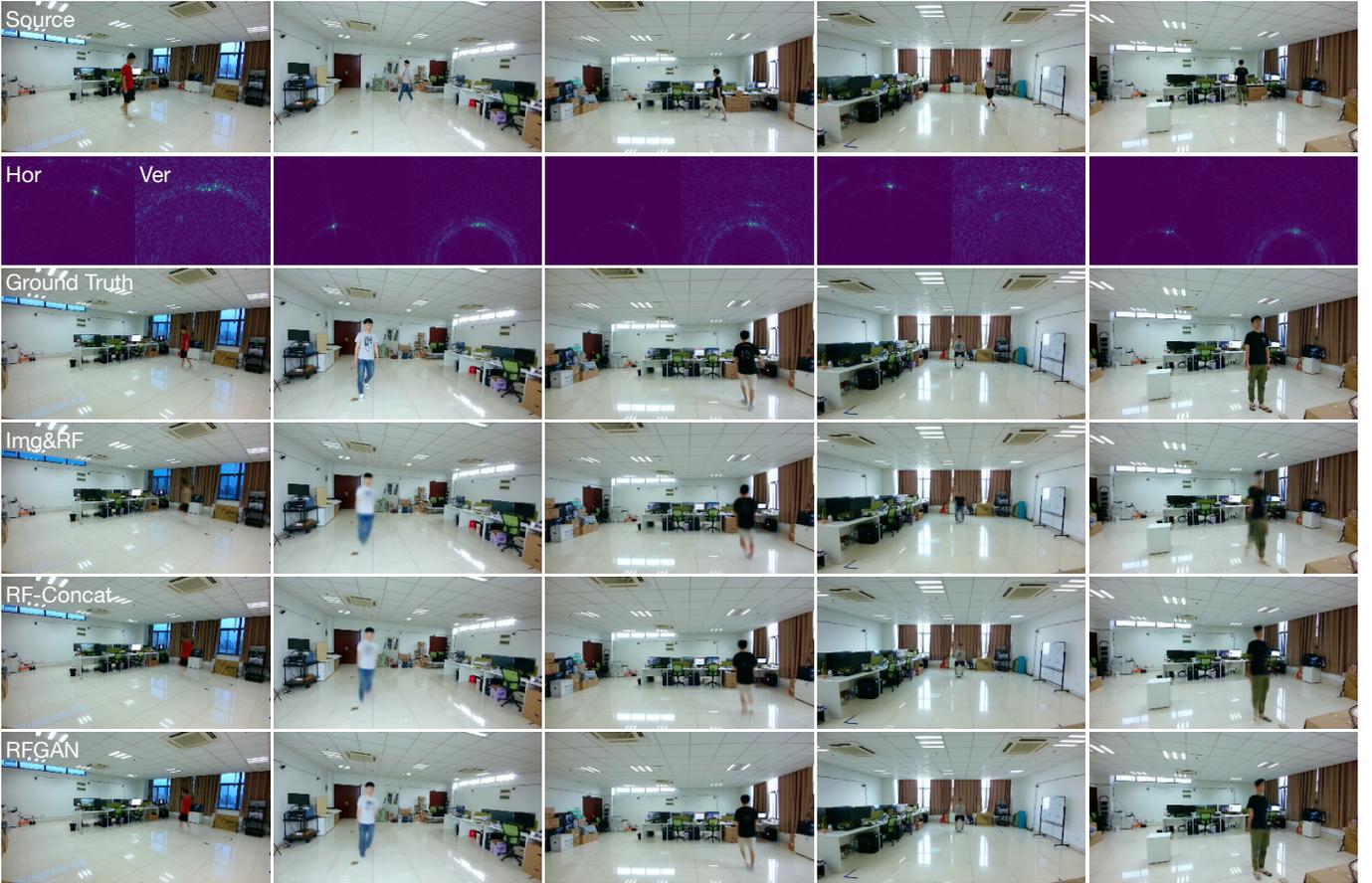


Fig. 9. Qualitative comparison of different methods. The 1st row shows the source frames. The 2nd row shows the horizontal and vertical RF heatmaps. The 3rd row shows the ground-truth human activity frames captured by the optical camera. The 4th to the 6th rows show the generated results by Img&RF, RF-Concat, and RFGAN.

Methods	<i>RF-Walk</i>			<i>RF-Activity</i>		
	FID ↓	SSIM ↑	User study ↑	FID ↓	SSIM ↑	User study ↑
Img&RF	27.84	0.9622	42.11%	22.03	0.9643	35.89%
RF-Concat	21.08	0.9689	69.23%	19.19	0.9707	68.42%
RFGAN	15.75	0.9695	80.76%	15.05	0.9708	78.12%

TABLE I
QUANTITATIVE COMPARISON OF DIFFERENT METHODS.

resized to 320×180 , and the shape of each RF heatmap is 201×160 .

Training details The proposed model is trained using Adam solver. The learning rate is set to 0.0002 for both the generative part and the discriminative part. The number of epochs is 80 and the batch size is 2. The hyperparameters λ , α , and β are equal and set to 10.0. We implement our method using PyTorch and all experiments can be run on a commodity workstation with a single GTX-1080 graphics card.

B. Evaluation Metric

We evaluate our proposed model from the following aspects:

- **Image Quality (FID):** We use the most popular metric FID [55] to evaluate the quality of the generated images.

It computes the Fréchet Inception Distance between the sets of generated images and the real images. The smaller the distance, the better the quality.

- **Image Similarity (SSIM):** For each test sample, we calculate the visual structural similarity (SSIM) [56] to measure the similarity between the generated and the ground-truth human frames. A higher value means that the model can generate a human frame more similar to the ground truth.

- **User Study:** We conduct user surveys to evaluate whether our model can synthesize human frames with correct positions and postures. We first show our subjects some generated human frames and the corresponding ground-truth human frames, then each subject is asked to assess (yes/no) the generated results based on human positions and postures. There are 10

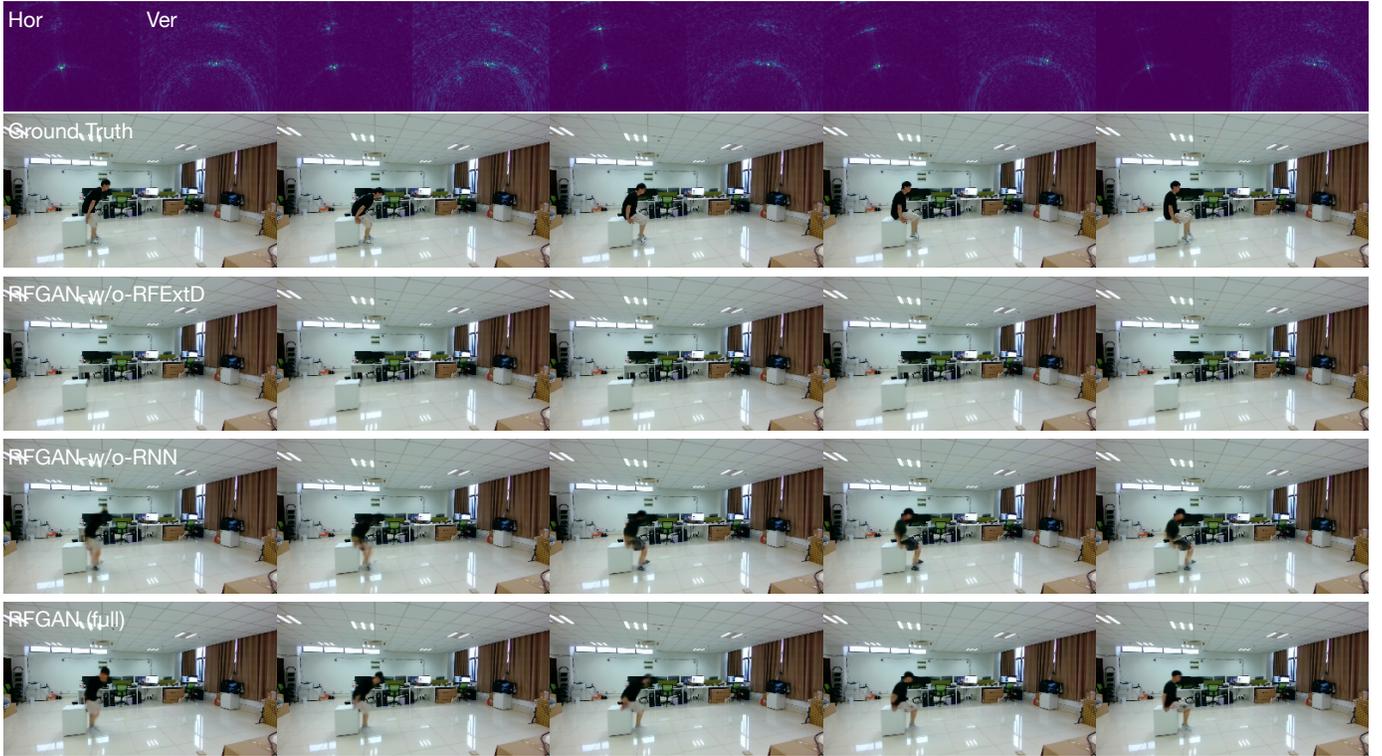


Fig. 10. Qualitative comparison of the ablation study. The 1st row shows the horizontal and vertical RF heatmaps. The 2nd row shows the ground-truth human activity frames captured by the optical camera. The 3rd to 5th rows show the generated results by RFGAN-w/o-RFExtD, RFGAN-w/o-RNN, and RFGAN (full).

subjects involved in the study.

C. Baselines

To our knowledge, this work is the first attempt that utilizes the RF signals to generate realistic human activity frames and there is no existing and suitable baseline method to be compared with. Therefore, we modify our model with some classic approaches that are widely used in GANs or related works, and the modified models are set as the baselines:

- **Img&RF:** To enable the RF-based condition setting, we propose a RF-Extractor with RNN to encode RF heatmaps and use RF-InNorm to inject the extracted information. Another alternative approach is to concatenate the RF condition with the input image directly, which is effective when the conditions have explicit guidance for GAN, e.g., pose-guided human synthesis [25]. However, the RF conditions are obscure data and have totally different spatial structures with optical images.
- **RF-Concat:** In our model, we propose a novel RF-Fusion operation to combine the horizontal and the vertical RF information, whereas the state-of-the-art approach for fusing RF information is to concatenate the features from RF signals along the channel directly, as in [13], [57]. We can find most existing learning-based RF sensing works just follow the common approach in computer vision literature to combine the two-dimensional RF information. In this paper, we design a specialized operation for RF signal data.

- **Ground Truth:** Another baseline is the ground-truth human activity frames captured by the optical camera.

The qualitative and quantitative comparisons are shown in Figure 9 and Table I. From the visual results, we can see that our proposed RFGAN model can capture the human position and posture information from RF signals and generate desirable activity frames. Although the baselines, i.e., RF-Concat and Img&RF, can capture the human position information, people in the generated frames are quite blurred. The user survey results also confirm the higher position and posture accuracy of the generated human activity frames by RFGAN. According to the FID and SSIM measurements, we find that the human activity frames generated by the proposed RFGAN have better quality and are more similar to the ground truth. The experimental results demonstrate the effectiveness of our proposed RF-Fusion and the RF conditioning encoding network.

D. Ablation Study

In this subsection, we conduct ablation studies to evaluate some important components in our proposed RFGAN model:

- **RFGAN-w/o-RFExtD:** In our full RFGAN model, there are two RF-Extractors and RNNs, one in the generative part and the other in the discriminative part. In RFGAN-w/o-RFExtD, We remove the RF-Extractor and RNN in the discriminative part and use the RF fused representation extracted in the generative part as the condition for Activity-Discriminator.
- **RFGAN-w/o-RNN:** We remove the RNN module from our full RFGAN model in this setting, which means the RFGAN-

Methods	<i>RF-Walk</i>			<i>RF-Activity</i>		
	FID ↓	SSIM ↑	User study ↑	FID ↓	SSIM ↑	User study ↑
RFGAN-w/o-RFExtD	58.36	0.9618	0.00%	45.71	0.9630	0.00%
RFGAN-w/o-RNN	16.41	0.9691	71.79%	18.11	0.9705	72.00%
RFGAN (full)	15.75	0.9695	80.76%	15.05	0.9708	78.12%

TABLE II
QUANTITATIVE COMPARISON FOR THE ABLATION STUDY.

w/o-RNN generates human activity frames only based on current RF signal inputs.

The qualitative and quantitative results are shown in Figure 10 and Table II. We find that the dual RF-Extractors and RNNs setting under the adversarial learning framework, i.e., RFGAN (full), can synthesize the target human activity frames, which is mainly due to the fact that the RF-Extractors and RNNs in the generative and the discriminative parts have different focuses. The RF-Extractor and RNN in the generative part pay more attention to guide the Generator for better synthesis, whereas the RF-Extractor and RNN in the discriminative part aim to help the Activity-Discriminator to distinguish different human poses. They are trained by adversarial learning. For RFGAN-w/o-RFExtD, only one RF-Extractor and RNN are used for RF conditioning encoding. Due to the lack of supervision labels for training guidance or another RF-Extractor and RNN for adversarial learning, one RF-Extractor and RNN cannot get the desirable human activity information from the RF signals and lead to ignoring the human object in the generated frames (see 3rd row in Figure 10). For RFGAN-w/o-RNN, due to the lack of information from the RF input neighbors, which can be used to adjust the current RF fused representation, it performs worse than the full RFGAN model. From the visual results, we can see the people in the generated frames by RFGAN-w/o-RNN are full of artifacts (4th row in Figure 10).

E. Deployment in New Scene

To deploy our model in a new scene, we do not need to retrain the whole model from the start. We can fine-tune the pre-trained RFGAN using very little data (about 40s data) to get similar results (see Table III). Specifically, the learning rate during the fine-tuning process is 0.0002 for both the generative part and the discriminative part. The epochs and batch size are set to 40 and 2, respectively. The loss functions and hyperparameters are the same with the training stage. From the quantitative results, we find that the pre-trained RFGAN model can generate desirable human activity frames in the new scene after fine-tuning with only a little data, which means our proposed model has the potential for being widely used.

	FID ↓	SSIM ↑	User study ↑
New Scene	20.64	0.9739	73.33%

TABLE III
QUANTITATIVE RESULTS IN THE NEW SCENE.

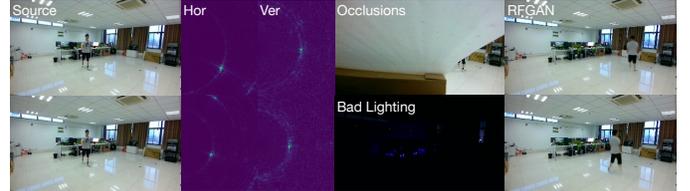


Fig. 11. Performance under occlusions and bad lighting.

F. Occlusions and Bad Lighting

RF signals can traverse occlusions and do not rely on lights, thus our model can work in the occluded or bad lighting environments (see Figure 11).

VI. LIMITATIONS

Since our method relies on the natural characteristics of RF signals, the solution that we present in this paper has some limitations. Firstly, in our mmWave radar system, the depth resolution of the RF signals is about 7.5cm, and the angular resolution is about 1.3 degrees. Thus, some micro-motion behaviors that are smaller than the resolution thresholds may be missed by our model. Secondly, the operating distance of our radar system depends on the transmission power, which works up to 20m. Finally, the datasets we use in this paper mainly contain the data of human daily activities under indoor scenes. Exploring more RF-based sensing models and synthesizing people in the wild is left for future work.

VII. CONCLUSION

In this paper, we aim to use RF signals to guide human synthesis. To tackle the challenge of using this new kind of driving signal, we propose a novel RFGAN model, which introduces a RF-Extractor with RNN to obtain the human activity information from the horizontal and vertical RF heatmaps and utilize the RF-InNorm to inject the information into the GAN networks. Furthermore, we propose to train the RF-Extractor and RNN under an adversarial learning framework to enable the encoding of the new kind of conditional data. To evaluate our proposed model, we create two cross-modal datasets and the experimental results show that the RFGAN can achieve a promising performance. We believe this work opens up a research opportunity to use a new form of conditional data, i.e., RF signals, to guide the GAN model, and the performance of the RFGAN confirms that RF signals have great potential in the imaging applications.

REFERENCES

- [1] F. Adib and D. Katabi, "See through walls with wifi!" in *ACM SIGCOMM*, 2013.
- [2] M. Kotaru, K. Joshi, D. Bharadia, and S. Katti, "Spotfi: Decimeter level localization using wifi," in *ACM SIGCOMM*, 2015.
- [3] X. Li, D. Zhang, Q. Lv, J. Xiong, S. Li, Y. Zhang, and H. Mei, "Indotrack: Device-free indoor human tracking with commodity wi-fi," in *ACM IMWUT*, 2017.
- [4] R. Ghazalian, A. Aghagolzadeh, and S. M. H. Andargoli, "Energy optimization and qoe satisfaction for wireless visual sensor networks in multi target tracking scenario," *IEEE Transactions on Multimedia*, 2020.
- [5] K. Qian, C. Wu, Z. Yang, Y. Liu, and K. Jamieson, "Widar: Decimeter-level passive tracking via velocity monitoring with commodity wi-fi," in *ACM MOBIHOC*, 2017.
- [6] Y. Chen, X. Su, Y. Hu, and B. Zeng, "Residual carrier frequency offset estimation and compensation for commodity wifi," *IEEE Transactions on Mobile Computing*, 2019.
- [7] D. Zhang, Y. He, X. Gong, Y. Hu, Y. Chen, and B. Zeng, "Multitarget aoa estimation using wideband lfmw signal and two receiver antennas," *IEEE Trans. Veh. Technol.*, 2018.
- [8] D. Zhang, Y. Hu, and Y. Chen, "Mtrack: Tracking multi-person moving trajectories and vital signs with radio signals," *IEEE Internet Things J.*, 2020.
- [9] Y. Liu, W. Zhou, M. Xi, S. Shen, and H. Li, "Multi-modal context propagation for person re-identification with wireless positioning," *IEEE Transactions on Multimedia*, 2021.
- [10] Y. Chen, H. Deng, D. Zhang, and Y. Hu, "Speednet: Indoor speed estimation with radio signals," *IEEE Internet Things J.*, 2020.
- [11] K. Qian, C. Wu, Z. Yang, Y. Liu, F. He, and T. Xing, "Enabling contactless detection of moving humans with dynamic speeds using csi," *ACM Trans. Embed. Comput. Syst.*, 2018.
- [12] C. Wu, Z. Yang, Z. Zhou, X. Liu, Y. Liu, and J. Cao, "Non-invasive detection of moving and stationary human with wifi," *IEEE J. Sel. Areas Commun.*, 2015.
- [13] 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.
- [14] M. Zhao, Y. Tian, H. Zhao, M. A. Alsheikh, T. Li, R. Hristov, Z. Kabelac, D. Katabi, and A. Torralba, "Rf-based 3d skeletons," in *ACM SIGCOMM*, 2018.
- [15] F. Wang, S. Zhou, S. Panev, J. Han, and D. Huang, "Person-in-wifi: Fine-grained person perception using wifi," in *ICCV*, 2019.
- [16] W. Jiang, H. Xue, C. Miao, S. Wang, S. Lin, C. Tian, S. Murali, H. Hu, Z. Sun, and L. Su, "Towards 3d human pose construction using wifi," in *MOBICOM*, 2020.
- [17] K. Niu, F. Zhang, X. Wang, Q. Lv, H. Luo, and D. Zhang, "Understanding wifi signal frequency features for position-independent gesture sensing," *IEEE Transactions on Mobile Computing*, 2021.
- [18] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *NeurIPS*, 2014.
- [19] M. Mirza and S. Osindero, "Conditional generative adversarial nets," in *arXiv preprint arXiv:1411.1784*, 2014.
- [20] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee, "Generative adversarial text to image synthesis," in *ICML*, 2016.
- [21] H. Zhang, T. Xu, H. Li, S. Zhang, X. Wang, X. Huang, and D. Metaxas, "Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks," in *ICCV*, 2017.
- [22] B. Li, X. Qi, T. Lukasiewicz, and P. H. Torr, "Manigan: Text-guided image manipulation," in *CVPR*, 2020.
- [23] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *CVPR*, 2017.
- [24] Y. Chen, Y.-K. Lai, and Y.-J. Liu, "Cartoogan: Generative adversarial networks for photo cartoonization," in *CVPR*, 2018.
- [25] H. Tang, D. Xu, G. Liu, W. Wang, N. Sebe, and Y. Yan, "Cycle in cycle generative adversarial networks for keypoint-guided image generation," in *ACM MM*, 2019.
- [26] T. Park, M.-Y. Liu, T.-C. Wang, and J.-Y. Zhu, "Semantic image synthesis with spatially-adaptive normalization," in *CVPR*, 2019.
- [27] C. Yu, Y. Hu, Y. Chen, and B. Zeng, "Personalized fashion design," in *ICCV*, 2019.
- [28] L. Chen, L. Wu, Z. Hu, and M. Wang, "Quality-aware unpaired image-to-image translation," *IEEE Transactions on Multimedia*, 2019.
- [29] M. Zhang and Q. Ling, "Supervised pixel-wise gan for face super-resolution," *IEEE Transactions on Multimedia*, 2020.
- [30] H. Emami, M. M. Aliabadi, M. Dong, and R. B. Chinnam, "Spa-gan: Spatial attention gan for image-to-image translation," *IEEE Transactions on Multimedia*, 2020.
- [31] E. Richardson, Y. Alaluf, O. Patashnik, Y. Nitzan, Y. Azar, S. Shapiro, and D. Cohen-Or, "Encoding in style: a stylegan encoder for image-to-image translation," in *CVPR*, 2021.
- [32] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *ICCV*, 2017.
- [33] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *CVPR*, 2019.
- [34] Y. He, Y. Chen, Y. Hu, and B. Zeng, "Wifi vision: Sensing, recognition, and detection with commodity mimo-ofdm wifi," *IEEE Internet of Things Journal*, 2020.
- [35] Y. Zeng, P. H. Pathak, and P. Mohapatra, "Wiwho: wifi-based person identification in smart spaces," in *IPSN*, 2016.
- [36] L. Fan, T. Li, R. Fang, R. Hristov, Y. Yuan, and D. Katabi, "Learning longterm representations for person re-identification using radio signals," in *CVPR*, 2020.
- [37] C.-Y. Hsu, R. Hristov, G.-H. Lee, M. Zhao, and D. Katabi, "Enabling identification and behavioral sensing in homes using radio reflections," in *CHI*, 2019.
- [38] D. Zhang, Y. Hu, Y. Chen, and B. Zeng, "Breathtrack: Tracking indoor human breath status via commodity wifi," *IEEE Internet Things J.*, 2019.
- [39] S. Yue, H. He, H. Wang, H. Rahul, and D. Katabi, "Extracting multi-person respiration from entangled rf signals," in *ACM IMWUT*, 2018.
- [40] T. Rahman, A. T. Adams, R. V. Ravichandran, M. Zhang, S. N. Patel, J. A. Kientz, and T. Choudhury, "Dopplesleep: A contactless unobtrusive sleep sensing system using short-range doppler radar," in *ACM UBICOMP*, 2015.
- [41] D. Zhang, Y. Hu, Y. Chen, and B. Zeng, "Calibrating phase offsets for commodity wifi," *IEEE Syst. J.*, 2019.
- [42] C.-Y. Hsu, A. Ahuja, S. Yue, R. Hristov, Z. Kabelac, and D. Katabi, "Zero-effort in-home sleep and insomnia monitoring using radio signals," in *ACM IMWUT*, 2017.
- [43] M. Zhao, S. Yue, D. Katabi, T. S. Jaakkola, and M. T. Bianchi, "Learning sleep stages from radio signals: A conditional adversarial architecture," in *ICML*, 2017.
- [44] J. Yue-Hei Ng, M. Hausknecht, S. Vijayanarasimhan, O. Vinyals, R. Monga, and G. Toderici, "Beyond short snippets: Deep networks for video classification," in *CVPR*, 2015.
- [45] S. Song, C. Lan, J. Xing, W. Zeng, and J. Liu, "An end-to-end spatio-temporal attention model for human action recognition from skeleton data," in *AAAI*, 2017.
- [46] S. Zhang, Y. Yang, J. Xiao, X. Liu, Y. Yang, D. Xie, and Y. Zhuang, "Fusing geometric features for skeleton-based action recognition using multilayer lstm networks," *IEEE Transactions on Multimedia*, 2018.
- [47] C. Si, Y. Jing, W. Wang, L. Wang, and T. Tan, "Skeleton-based action recognition with spatial reasoning and temporal stack learning," in *ECCV*, 2018.
- [48] S. Agethen and W. H. Hsu, "Deep multi-kernel convolutional lstm networks and an attention-based mechanism for videos," *IEEE Transactions on Multimedia*, 2019.
- [49] C. Si, W. Chen, W. Wang, L. Wang, and T. Tan, "An attention enhanced graph convolutional lstm network for skeleton-based action recognition," in *CVPR*, 2019.
- [50] C. Ventura, M. Bellver, A. Girbau, A. Salvador, F. Marques, and X. Giro-i Nieto, "Rvos: End-to-end recurrent network for video object segmentation," in *CVPR*, 2019.
- [51] C. Lu, M. Hirsch, and B. Scholkopf, "Flexible spatio-temporal networks for video prediction," in *CVPR*, 2017.
- [52] M. A. Richards, *Fundamentals of radar signal processing*. McGraw-Hill Education, 2014.
- [53] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, and S. Paul Smolley, "Least squares generative adversarial networks," in *ICCV*, 2017.
- [54] L. Mescheder, A. Geiger, and S. Nowozin, "Which training methods for gans do actually converge?" *arXiv preprint arXiv:1801.04406*, 2018.
- [55] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, G. Klambauer, and S. Hochreiter, "Gans trained by a two time-scale update rule converge to a nash equilibrium," in *arXiv preprint arXiv:1706.08500*, 2017.
- [56] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: From error visibility to structural similarity," *IEEE Trans. Image Process.*, 2004.
- [57] A. Sengupta, F. Jin, R. Zhang, and S. Cao, "mm-pose: Real-time human skeletal posture estimation using mmwave radars and cnns," *IEEE Sensors Journal*, 2020.