# Integrated Sensing and Communication from Learning Perspective: An SDP3 Approach

Guoliang Li, Shuai Wang, Jie Li, Rui Wang, Fan Liu,
Xiaohui Peng, Tony Xiao Han, and Chengzhong Xu, *Fellow, IEEE*

*Abstract*—Characterizing the sensing and communication performance tradeoff in integrated sensing and communication (ISAC) systems is challenging in the applications of learning-based human motion recognition. This is because of the large experimental datasets and the black-box nature of deep neural networks. This paper presents SDP3, a Simulation-Driven Performance Predictor and oPtimizer, which consists of SDP3 data simulator, SDP3 performance predictor and SDP3 performance optimizer. Specifically, the SDP3 data simulator generates vivid wireless sensing datasets in a virtual environment, the SDP3 performance predictor predicts the sensing performance based on the function regression method, and the SDP3 performance optimizer investigates the sensing and communication performance tradeoff analytically. It is shown that the simulated sensing dataset matches the experimental dataset very well in the motion recognition accuracy. By leveraging SDP3, it is found that the achievable region of recognition accuracy and communication throughput consists of a communication saturation zone, a sensing saturation zone, and a communication-sensing adversarial zone, of which the desired balanced performance for ISAC systems lies in the third one.

*Index Terms*—Integrated sensing and communication, resource allocation.

## I. INTRODUCTION

INTEGRATED sensing and communication (ISAC) is a promising technology for the next generation cellular system and wireless local area network (WLAN). It is expected to considerably improve the spectral, energy and hardware efficiencies of wireless systems [2]. Since wireless resource is shared between sensing and communication functionalities in ISAC systems, it is of significant interest to investigate their tradeoff relation. However, the analysis of sensing performance could be challenging.

Generally, most of the sensing tasks can be classified into three categories, including *detection*, *estimation* and *recognition*. The detection tasks aim to determine the state of a

Part of this paper has been presented at the IEEE International Workshop on Signal Processing Advances in Wireless Communications (SPAWC), Lucca, Italy, Sep. 2021 [1]. *(Corresponding Author: Shuai Wang and Rui Wang.)*

Guoliang Li, Jie Li, Rui Wang, and Fan Liu are with the Department of Electrical and Electronic Engineering, Southern University of Science and Technology (SUSTech), Shenzhen 518055, China (e-mail: {ligl2020, lij2019}@mail.sustech.edu.cn, {wang.r, liuf6}@sustech.edu.cn).

Shuai Wang is with the Guangdong-Hong Kong-Macao Joint Laboratory of Human-Machine Intelligence-Synergy Systems, Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, Shenzhen, 518055, China (e-mail: s.wang@siat.ac.cn).

Xiaohui Peng and Tony Xiao Han are with Huawei technologies, Co. Ltd., Shenzhen, China (e-mail: {pengxiaohui5, tony.hanxiao}@huawei.com).

Chengzhong Xu is with the State Key Laboratory of Internet of Things for Smart City (SKLIOTSC), Department of Computer Science, University of Macau, Macau, China (e-mail: czxu@um.edu.mo).

target, such as presence/absence. The detection probability is a typical metric to quantify the detector performance. In [3], an integrated communication and passive radar system was considered, where a generalized likelihood ratio test (GLRT) was proposed for target detection and the corresponding detection probability was approximated in a closed-form formula. Accordingly, the detection probability was maximized subject to a minimum communication rate constraint. Furthermore, the above method was extended to an integrated multi-static radar and communication system with the aim for power allocation optimization in [4].

The estimation tasks acquire the useful parameters, i.e., the distance, velocity and angle, from the sensed targets. For instance, a dual functional waveform was utilized to minimize the estimation mean square error (MSE) of the target response matrix while ensuring a worst communication performance in [5]. When the closed-form expression of MSE is not attainable, the Cramér-Rao bound (CRB) could be adopted for the sensing performance evaluation, which represents the lower bound of the variance of all the unbiased estimators. For example in [6], the CRB was minimized subject to the signal-to-interference-plus-noise ratio (SINR) constraint of communication receivers. Furthermore, in [7], [8], the estimation rate, quantifying the reduction of the uncertainties for the sensing parameters per second, was proposed, then the tradeoff analysis between the estimation rate and the communication rate was provided.

The above two sensing tasks are usually processed in the physical layer, while the recognition tasks are usually accomplished in the application layer aided by the machine intelligence. They aim to acquire the semantic information of the sensed targets, e.g., human motion recognition (HMR). For example, in [11], the short time Fourier transform (STFT) was adopted to generate spectrograms of the sensing data. Then the support vector machine (SVM) was trained for motion classification, whose accuracy was higher than $90\%$. By replacing the SVM with a convolutional deep neural network (CNN), it was shown in [12] that a higher classification accuracy could be achievable. In fact, deep neural network has been adopted extensively in the application of activity recognition [13], [14], hand gesture recognition [15], [16], and gait recognition [17], [18] with radio wave. However, most of the HMR tasks are based on deep neural network, where the relation between sensing accuracy and sensing resource (e.g., sensing time and transmission power) can hardly be represented analytically. Moreover, the above works rely on the motion datasets generated from time-consuming and labor-intensive experiments. As a result, there are the following
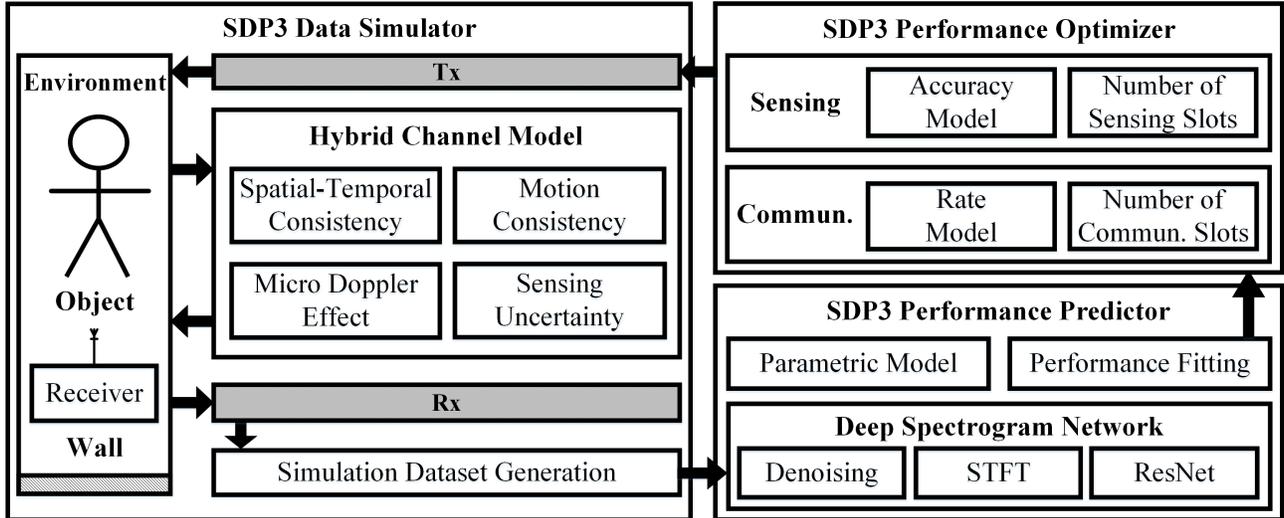
Fig. 1: The system diagram of the proposed SDP3 framework.

challenges in the sensing and communication tradeoff analysis with the recognition tasks: 1) there is no analytical model of recognition accuracy, which maps the sensing resource (e.g., sensing time, power and etc.) to the recognition accuracy; 2) even for numerical performance evaluation, there is no low-cost method to generate the motion sensing datasets for the training of recognition models.

In this paper, we would like to shed some light on the above open issues. In our preliminary work [1], a data-assisted hybrid channel model was proposed to simulate the received sensing signals. In this paper, we shall extend this work and propose a complete simulation-driven analysis framework, namely Simulation-Driven Performance Predictor and oPtimizer (SDP3), to quantify the tradeoff between recognition accuracy and communication throughput, which consists of the following three components as illustrated in Fig. 1.

1) The SDP3 data simulator mimicks experimental datasets in a virtual environment consisting of human target, communication receivers, and static objects (e.g., walls), such that the experimental costs can be saved.

2) The SDP3 performance predictor trains a deep spectrogram network (DSN) with the above dataset for motion recognition, tests the recognition accuracy, and approximates the relation between recognition accuracy and sensing duration with power function regression.

3) The SDP3 performance optimizer illustrates tradeoff performance between recognition accuracy and communication throughput. It is shown that the derived accuracy-throughput (A-T) region consists of a communication saturation zone, a sensing saturation zone, and a sensing-communication adversarial zone, where the last zone achieves the best balance between sensing and communication.

The rest of the paper is organized as follows. The system model considered in this paper is introduced in Section 2. The SDP3 data simulator is presented in Section 3. The procedure of the SDP3 performance predictor is described in Section 4. The proposed SDP3 performance optimizer and the sensing and communication tradeoff analysis are presented in Section 5. The simulation and experiment results are demonstrated in Section 6. Finally, the conclusion is drawn in Section 7.

## II. SYSTEM MODEL

An SDP3 framework to characterize the communication and human motion sensing tradeoff is proposed in this paper. In order to demonstrate the procedure of the framework, we consider an time-division-multiple-access-based (TDMA-based) ISAC system located in a conference room as an example, where the conference room follows the same specification as that in IEEE 802.11ay [35].

As illustrated in Fig. 2, the system consists of one ISAC-enabled base station (BS), $K$ communication receivers and one human target to be sensed. Both the radar and communication modulars are implemented at the BS, which are multiplexed in time domain. In order to facilitate transmission and sensing scheduling, the time is organized by slots, including the sensing slots and communication slots. The slot duration is $T_s$ and the channel impulse response is assumed to be quasi-static in each slot. The communication modular is enabled in the communication slots for downlink data transmission, and the radar modular is enabled in the sensing slots for the motion recognition of the human target. The sensing and communication tradeoff is studied by adapting the number of sensing and communication slots in every N slots, which are referred to as a scheduling period. The numbers of communication and sensing slots in a scheduling period are denoted as $N_c$ and $N_s$, respectively. Thus, $N_c + N_s = N$. Moreover, in order to better capture the micro-Doppler effect, the sensing slots are not successive. Instead, there are $N_m$ communication slots between two neighboring sensing slots. Large $N_s$ could lead to better resolution of micro-Doppler effect and capture more motion details of human target, at the price of lower communication throughput. The communication and sensing models are elaborated below, respectively.
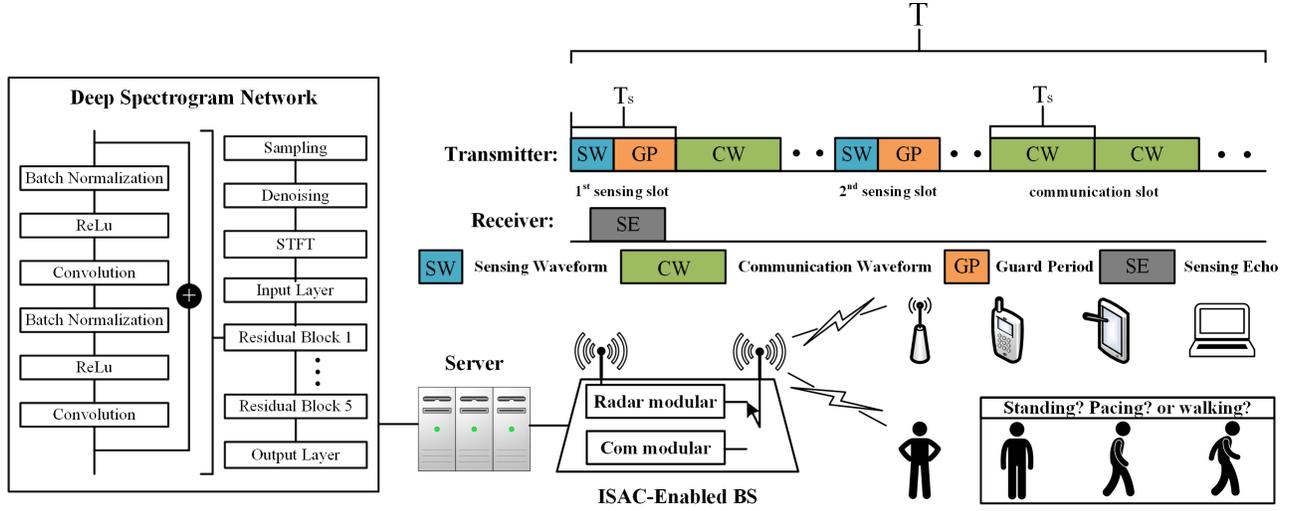
Fig. 2: Illustration of the ISAC system model.

## A. Communication Model

One receiver is selected in one communication slot for downlink transmission. The transmission signal is modulated via orthogonal frequency-division multiplexing (OFDM) with M subcarriers. Let $h_{k,j}(t)$ be the channel impulse response from the BS to the k-th communication receiver in the j-th slot. [1] The channel gain of the m-th subcarrier from the BS to the k-th receiver can be written as

$$H_{k,j,m} = \sum_{l=0}^{L-1} h_{k,j}(t - lT) e^{\frac{-j2\pi ml}{M}}, \ 0 \le m \le M-1, \quad (1)$$

where T is the sampling period of OFDM transceiver [36]. Then, the received signal of the m-th subcarrier is given by

$$Y_{k,j,m} = H_{k,j,m} X_{k,j,m} + Z_{k,j,m}, \quad (2)$$

where $X_{k,j,m}$ and $Z_{k,j,m} \sim \mathcal{CN}(0, \sigma_z^2)$ are the transmitted signal and the white Gaussian noise, respectively, and $\sigma_z^2$ is the noise power. As a result, the throughput of the k-th receiver in the j-th slot (if this receiver is selected) can be expressed as

$$R_{k,j} = \frac{T_s}{T_o} \sum_{m=1}^{M} \log_2 (1 + \gamma_{k,j,m}), \quad (3)$$

where $T_o$ is the duration of one complete OFDM symbol (including the cyclic prefix),

$$\gamma_{k,j,m} = \frac{|H_{k,j,m}|^2 P_{k,j,m}}{\sigma_z^2}, \quad (4)$$

and $P_{k,j,m}$ is the power allocated to the m-th subcarrier. Notice that the total power constraint $\sum_{m=1}^{M} P_{k,j,m} \le P$ should be satisfied in the power allocation.

## B. Sensing Model

In each sensing slot, the frequency-modulated continuous wave (FMCW) is broadcasted for human motion detection, followed by a guarding interval. Let $s(t)$ denote the FMCW and $h_{0,j}(t)$ be the channel impulse response between the radar modular transmitter and receiver of the j-th slot (if the j-th slot is sensing slot), the received signal at the radar modular is

$$r_j(t) = h_{0,j}(t) * s(t) + n_j(t), \quad (5)$$

where $n_j(t)$ is the Gaussian noise with average power $\sigma_z^2$, and $s(t)$ follows the power constraint $\frac{1}{T} \int_{-\frac{T}{2}}^{\frac{T}{2}} |s(t)|^2 dt = P$.

Based on the received signals of the sensing slots $\{r_j(t)|j \bmod (N_m + 1) = 1, \ 1 \le j \le N_s(N_m + 1)\}$, the human motion can be recognized via a neural network which should have been trained in advance via a dataset collected in the same environment. Note it is inefficient to collect dataset of $\{r_j(t)|j \bmod (N_m + 1) = 1, \ 1 \le j \le N_s(N_m + 1)\}$ via experiments for various environments, and it is also difficult to investigate the recognition accuracy of neural network with respect to $N_s$ analytically. The study of communication and sensing tradeoff is challenging. In order to address the above issues, in the following Section 3, the SDP3 data simulator based on a data-assisted hybrid channel (DAHC) model is proposed to generate the datasets of received sensing signals via simulation, which has been adopted as a part of the channel model for WLAN (wireless local area network) sensing in IEEE 802.11bf [37]. Based on it, a Deep Spectrogram Network (DSN) for motion recognition is proposed, and the approximated expression of recognition accuracy versus the number of sensing slots, denoted as $A = \Theta(N_s)$, is obtained in Section 4, which is referred to as the SDP3 performance predictor. Finally, the tradeoff between the sensing and communication performance is studied via SDP3 performance optimizer in Section 5.

---

[1]Although there is Doppler effect in the channel due to the motion of sensing target, the phase shift raised by the Doppler effect is negligible in one slot. Moreover, the propagation paths without Doppler shift is dominant in the overall channel impulse response. Hence, it is assumed that the CSI is quasi-static in one slot.

TABLE I: A Comparison of Existing and Proposed Channel Models

| Type | Literature | Methodology | Application Scenario | Spatial-Temporal Consistency | Motion Consistency | Micro Doppler | Sensing Uncertainty |
|------|-----------|-------------|---------------------|------------------------------|--------------------|---------------|---------------------|
| S | [28] | cluster random process | commun. | L1 | ✗ | ✗ | ✗ |
| D | [29] | ray tracing | both | L2 | ✓ | ✓ | ✗ |
|  | [30] | primitive based | sensing | L2 | ✓ | ✓ | ✗ |
| Q–D | 3GPP TR 38.901 [33] | GBSM based | commun. | L2 | ✗ | ✗ | ✗ |
|  | QuaDRiGa [34] | GBSM based | commun. | L2 | ✗ | ✗ | ✗ |
|  | IEEE 802.11 ay [35] | Q-D based | commun. | L2 | ✗ | ✗ | ✗ |
|  | METIS [26] | GBSM, map based hybrid | commun. | L2 | ✗ | ✗ | ✗ |
|  | **Proposed** | primitive based hybrid | both | L2 | ✓ | ✓ | ✓ |

GBSM means geometry-based stochastic channel model.
"S" means statistical, "D" means deterministic, "Q–D" means quasi-deterministic.
"L1" means large-scale spatial consistency, "L2" means large-scale and small-scale spatial consistency.
"✓" means functionality supported, "✗" means functionality not supported.

## III. SDP3 DATA SIMULATOR

In this section, we first summarize the drawbacks of the existing channel models in sensing dataset generation, then propose a novel data-assisted hybrid channel (DAHC) model for the puropose of efficient sensing dataset generation. Moreover, the human kinematic model for the motion simulation is also introduced.

### A. Preliminaries

A wireless channel model for sensing performance evaluation should generate consistent channel impulse response spatially and temporally so that the receiver can capture the micro-Doppler effect via the received signals in a time interval. In fact, the following two kinds of consistency in both spatial and temporal domains have been proposed for communication channel model design [26]: 1) large-scale spatial-temporal consistency refers to consistent power fading, delay spreads and angular spreads at two close locations and time instances; 2) small-scale spatial-temporal consistency refers to consistent delays and angles of rays at two close locations and time instances. Besides them, the channel model for wireless sensing should be composed of the rays consistent with the environment and the motions of sensing target. Finally, random interference should also be included to model the unpredictable motions except the sensing target.

As summarize in Table I, current channel models can be categorized into statistical models, deterministic models, and quasi-deterministic models. As an example of the statistical model proposed in [28], the Non-line-of-sight (NLoS) rays between the transmitter and the receiver were generated via scattering clusters. The phases, amplitudes and delays of rays were generated via independent distributions. Lack of the spatial-temporal consistency in small scale, this model cannot simulate the micro-Doppler effects due to the non-rigid human motions.

On the other hand, the quasi-deterministic models could maintain the spatial-temporal consistency in both large and small scales, which are adopted in many existing industrial standards for communication systems [26], [33]–[35]. In [26], [33], [34], the large-scale channel parameters, including delay spread, angular spreads, Ricean K-factor and shadow fading, were computed via ray tracing and the small-scale channel parameters, inculding time delays, cluster powers and arrival/departure angles, were obtained using statistical method, where the small-scale consistency is maintained via the distribution correlation. For example, the delays and angles of rays were generated based on uniform distributions in [33] where the distribution parameters linearly depend on the correlation distance. In [38], a method to calculate the correlation distance according to the large-scale parameters was proposed. In [35], a channel model consisting of both deterministic and random rays was proposed, where the deterministic rays were generated via ray tracing to maintain the small-scale consistency. However, the modeling of human motions is not considered in the above channel models. For example, the human body was simply treated as a blocker of rays in [33] and [35]. Thus, the micro-Doppler effect of human motion are not characterized in the above models.

The deterministic channel models such as ray-tracing [29] have potential to capture information of both environment and human target. However, the computation complexity could be huge if the ray tracing method is directly applied on non-rigid human motions. On the other hand, primitive-based
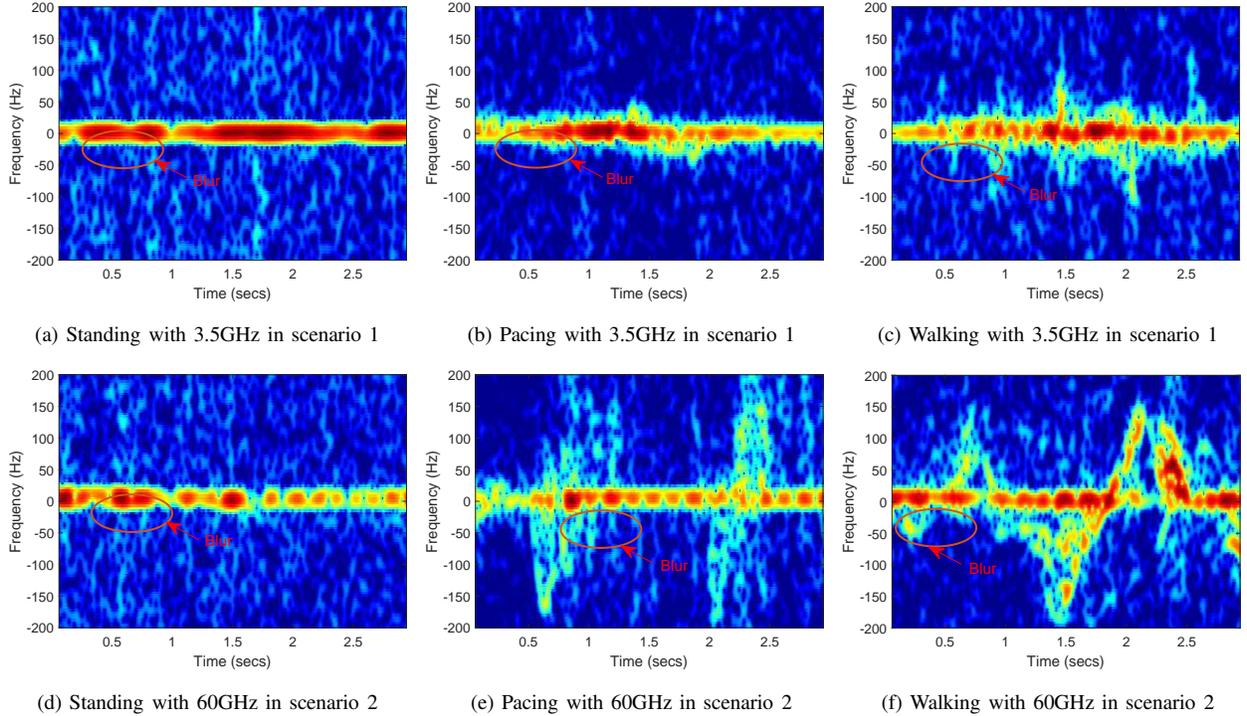
Fig. 3: The experimental results of wireless sensing in scenario 1 and 2.

(a) Standing with 3.5GHz in scenario 1

(b) Pacing with 3.5GHz in scenario 1

(c) Walking with 3.5GHz in scenario 1

(d) Standing with 60GHz in scenario 2

(e) Pacing with 60GHz in scenario 2

(f) Walking with 60GHz in scenario 2

channel model [30] is a computationally efficient approximation method for ray tracing, where each dynamic object is modeled as an extended target made of multiple point scatters (primitives) distributed along its body (e.g., a human body usually consists of 16 primitives). The received signal from each primitive is computed using the electromagnetic field method in [31, Sec. 5.8], where the radar cross section (RCS) for simple shapes could be obtained from [32].

As a summary, the quasi-deterministic channel models are able to simulate the rays' propagation in the sensing environment, and the primitive-based model is able to efficiently simulate the rays from human target. The integration of both models could be used for sensing dataset generation. Furthermore, due to the complexity and uncertainty of the real sensing scenario, it is necessary to keep randomness in the channel impulse response, which is referred to as the *sensing uncertainty* in this paper. Specifically, sensing uncertainty could be raised at least by the following three factors: 1) unpredictable reflections from the wall and random scatters, 2) undesired movements of the non-target objects and 3) noise. Examples of spectrograms generated from the real experiment are illustrated in Fig. 3. It can be observed that there are significant blurs aroused by the sensing uncertainty. To our best knowledge, the sensing uncertainty has not been not explicitly considered in the existing channel models.

### B. Channel Modeling

In order to keep the consistency with the sensing scenario and model the sensing uncertainty, the autoregressive model and quasi-deterministic model are integrated in the proposed data-assisted hybrid channel (DAHC) model as illustrated in

Fig. 4. Specifically, the channel impulse response from the BS to the k-th receiver (the 0-th receiver refers to the co-located BS radar receiver) in the j-th time slot can be represented by

$$h_{k,j}(t) = u_{k,j}(t) + v_{k,j}(t), \quad (6)$$

where

- $u_{k,j}(t)$ is the target-related channel component consisting of the rays reflecting from the human target.
- $v_{k,j}(t)$ is the target-unrelated channel component, consisting of the rays via the LoS path and the other reflection paths. For example, the reflection paths via walls. Due to the self-interference cancellation at the BS radar receiver, we neglect the LoS path in $v_{0,j}(t)$.

Similar to [30], the target related channel $u_{k,j}(t)$ is modeled using the following primitive-based method:

$$u_{k,j}(t) = \frac{A}{\sqrt{4\pi}} \sum_{b=1}^{B} \frac{\sqrt{G_{b,j}}}{D_{b,j}^2} \exp\left(-j\frac{2\pi f_c}{c} 2D_{b,j}\right)$$
$$\times \exp\left(j\varphi_b\right) \delta\left(t - \frac{2D_{b,j}}{c}\right). \quad (7)$$

In the above equation, $A$ is the constant related to wave length $\lambda$ and antenna gain $P_t$, which is $\lambda^2\sqrt{P_t}$ in the LoS case; B is the number of primitives; $G_{b,j}$ is the amplitude accounting for radar cross section of the $b$-th primitive in the j-th time slot; $D_{b,j}$ is the distance from the $b$-th primitive to the radar in the j-th time slot, $f_c$ is the carrier frequency; $c$ is the speed of light; $\varphi_b$ is the initial phase of the $b$-th ray, which follows uniform distribution in $[-\pi, \pi]$. $\delta(a)$ is the indicator function, whose value is 1 for $a = 0$ and 0 otherwise.

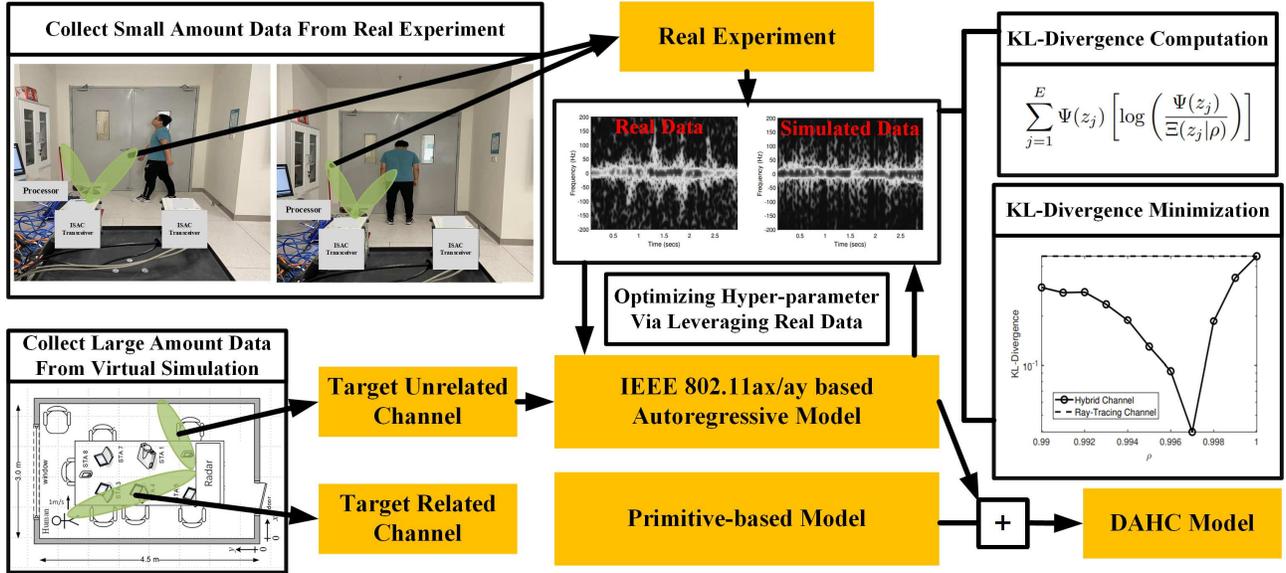On the other hand, it is not necessary to capture the mo-

Fig. 4: The framework of the proposed DAHC model.

tions via the primitive-based approach in the target-unrelated channel $v_{k,j}(t)$ as that in target-related channel. Hence, the following autoregressive method is used to model the sensing uncertainty statistically

$$v_{k,j}(t) = \begin{cases} \Upsilon_1(t), & \text{if } j = 1 \\ \rho v_{k,j-1}(t-T_0) + (1-\rho)\Upsilon_j(t), & \text{if } j > 1 \end{cases},$$

(8)

where $\Upsilon_1(t)$ is the target-unrelated channel impulse response of the first slot, and $\rho$ is a hyper-parameter controlling the intensity of sensing uncertainty. A larger (or smaller) $\rho$ leads to a weaker (or stronger) sensing uncertainty. In this paper, we adopt the quasi-deterministic channel model in [35] to generate $\Upsilon_j(t)$, $\forall j$, thus,

$$\Upsilon(t)_j = \sum_{n=1}^{N} \sqrt{H_n} \frac{\lambda}{4\pi(D_0 + \tau_n^{\text{cluster}}c)}$$
$$\times \left[ \sum_{m=1}^{M} a_{n,m}\exp\left(j\phi_{n,m}\right)\delta(t-\tau_{n,m}^{\text{ray}}) \right]. \quad (9)$$

In the above expression, N is the number of scattering clusters, $H_n$ is the reflection factor for both first-order and second-order reflections; $\lambda$ is the wave length; $\tau_n^{\text{cluster}}$ is the $n$-th cluster's time delay (in seconds) obtained from ray-tracing; $\tau_{n,m}^{\text{ray}}$, $a_{n,m}$ and $\phi_{n,m}$ are the time delay, amplitude and initial phase of $m$-th ray via the $n$-th cluster, which are obtained from Poisson distribution, Rayleigh distribution and uniform distribution respectively.

### C. Data-Assisted Model Calibration

In order to fit the hyper-parameter $\rho$ of the DAHC to the real uncertainty level in a particular sensing scenario, it is necessary to make a measurement in the target environment. Denote the received signal in the real scenario and DAHC model simulation as $\{r_j(t)|j \bmod (N_m+1) = 1, 1 \leq j \leq N_s(N_m+1)\}$ and $\{\hat{r}_j(t)|j \bmod (N_m+1) = 1, 1 \leq j \leq$

$N_s(N_m+1)\}$, respectively. The adaptation of hyper-parameter to real uncertainty level is elaborated below.

The spectrograms of both measured and simulated signals are first obtained according to Section IV. The Doppler frequency strength versus time and frequency of both spectrograms are quantized into F levels. The probability mass function (PMF) of Doppler frequency strength for both spectrograms can then be obtained, denoted as $\Psi$ and $\Xi$ respectively. Finally, the proper hyper-parameter $\rho$ is the one minimizing the Kullback-Leibler (KL) divergence as follows.

$$\min_{\rho} \quad \sum_{f=1}^{F} \Psi(z_f) \left[ \log\left( \frac{\Psi(z_f)}{\Xi(z_f|\rho)} \right) \right],$$
$$\text{s.t.} \quad 0 \leq \rho \leq 1, \quad (10)$$

where $z_f$ denotes the signal strength of the f-th level. The above problem can be solved by one-dimensional search. In practice, $\rho$ varies in different scenarios and we can store the values of $\rho$ for typical scenarios in a look-up table.

### D. Human Kinematic Modeling

The human target is represented by 16 primitives in the proposed DAHC model. In order to model the motion of the 16 primitives, a global human model, namely the Thalmann model [40], is adopted. In the Thalmann model, the human body is represented as a series of 16 segments, which corresponds to the primitives of the proposed channel model. Based on the biomechanical experimental data, the motions of all the segments, including the their positions and orientations of all time slots, are obtained for different categories human motions.

However, this global human model averages out the personification in the same category of motion, losing the diveristy in motion's dataset generation. In the future work, the motion capture methods from the the areas of Graphics [41] could be exploited to simulate the primitives' motion with personification.

## IV. SDP3 PERFORMANCE PREDICTOR

The SDP3 performance predictor is trained to predict the sensing performance versus the sensing resource. Specifically, with the motion sensing dataset generated by the SDP3 data simulator, a deep spectrogram network (DSN) is trained for motion recognition, then the approximate expression of the motion detection accuracy versus the number of sensing slots $N_s$ is derived.

### A. Deep Spectrogram Network

Let $C$ and $\mathcal{C} = \{1, \ldots, C\}$ be the number of human motion categories and the set of human motion categories respectively, the input of DSN is the signals $\{r_j(t) | j \mod (N_m + 1) = 1, \ 1 \le j \le N_s(N_m+1)\}$ and the output is the index of human motion category $\widehat{c} \in \mathcal{C}$. In DSN, the spectrogram of input signals is first generated via data cleaning and transformation, followed by model training (in the training phase) or inference (in the inference phase), as shown in Fig. 2.

**Data Cleaning**. In this step, the received signal $r_j(t)$ within the sweep duration $t \in [0, T_{sw}]$ is sampled with a frequency $f_s$. Let $\mathbf{x}_j = [r_j(\frac{1}{f_s}) \ r_j(\frac{2}{f_s}) \ldots r_j(T_{sw})]^T \in \mathbb{C}^L$ be the vector of samples in the j-th slot, where $L = T_{sw} f_s$, and $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_{N_s}] \in \mathbb{C}^{L \times N_s}$ be the aggregation of sample vectors of all sensing slots. The matrix $\mathbf{X}$ is a superposed signal consisting of the desired signals reflected from the target and undesired signals reflected from walls (or other objects). To extract the desired information, we first dechirp the sampled signal $\mathbf{X}$ as follows

$$\widetilde{\mathbf{X}} = \left( \left[ \widetilde{\mathbf{s}} \odot \mathbf{x}_1^* \ \widetilde{\mathbf{s}} \odot \mathbf{x}_2^* \ \cdots \ \widetilde{\mathbf{s}} \odot \mathbf{x}_{N_s}^* \right] \right)^*, \quad (11)$$

where $\widetilde{\mathbf{s}} = [s(\frac{1}{f_s}) \ s(\frac{2}{f_s}) \ldots s(T_{sw})]^T$, $\mathbf{x}_1^*$ denotes the conjugate of $\mathbf{x}_1$, and $\odot$ denotes the Hadamard product. Then the singular value decomposition (SVD) is applied to $\widetilde{\mathbf{X}}$, yielding $\widetilde{\mathbf{X}} = \sum_{i=1}^{D} a_i \mathbf{u}_i \mathbf{v}_i^H$, where $D$ is the rank of $\widetilde{\mathbf{X}}$, $a_i$ is the i-th largest singular value, $\mathbf{u}_i$ and $\mathbf{v}_i$ are the $i$-th left-singular vector and right-singular vector respectively. Removing the first $d - 1$ components, which represent the undesired signal paths, the denoised signal matrix is

$$\mathbf{Y} = \sum_{i=d}^{D} a_i \mathbf{u}_i \mathbf{v}_i^H. \quad (12)$$

**Data Transformation**. In this step, the Short Time Fourier transform (STFT) [19] is applied on $\mathbf{Y}$. We first define the sliding window function as

$$w_\beta[n] = \frac{I_0 \left( \beta \sqrt{1 - \left( \frac{n - (N_w - 1)/2}{((N_w-1)/2)} \right)^2} \right)}{I_0(\beta)}, 0 \le n \le N_w - 1,$$

where $I_0$ is the zeroth-order modified Bessel function of the first kind and $N_w$ is the length of sliding window. Let $\mathbf{y} = \mathbf{1}^T \mathbf{Y}$, the STFT of $\mathbf{y}$ at time $k$ and the frequency $f$ with the sliding window function $w$ can be expressed as

$$z[k, f] = \sum_{k'=-\infty}^{+\infty} y\left[k'\right] w_\beta \left[k' - k\right] \exp(-\mathrm{j}2\pi k' f/N),$$

$$k \in \left\{ 0, \ldots, \left\lfloor \frac{N_s - N_w}{N_w - N_{overlap}} \right\rfloor \right\} * (N_w - N_{overlap}),$$

TABLE II: Candidate Parameteric Learning Curve Models

| Name | Experssion | Tuning Parameters |
|---|---|---|
| vapor pressure | $\exp(\alpha + \beta/\mathrm{N_s})$ | $\alpha, \beta$ |
| pow$_3$ | $\gamma - \alpha \mathrm{N_s}^{-\beta}$ | $\alpha, \beta, \gamma$ |
| log power | $\alpha/(1 + (\mathrm{N_s}/e^\beta)^\gamma)$ | $\alpha, \beta, \gamma$ |
| exp$_4$ | $\gamma - e^{-\alpha \mathrm{N_s}^\epsilon + \beta}$ | $\alpha, \beta, \gamma, \epsilon$ |
| log log linear | $\log(\alpha \log(\mathrm{N_s}) + \beta)$ | $\alpha, \beta$ |
| ilog2 | $\beta - \alpha/\log(\mathrm{N_s})$ | $\alpha, \beta$ |
| pow$_4$ | $\gamma - (\alpha \mathrm{N_s} + \beta)^\epsilon$ | $\alpha, \beta, \gamma, \epsilon$ |

$$f \in \left\{ -\frac{N_{fft}}{2} + 1, \ldots, \frac{N_{fft}}{2} \right\}, \quad (13)$$

where $N_{overlap}$ and $N_{fft}$ specify the number of overlap samples between adjoining STFT windows and the number of frequency points respectively. Then the spectrogram of sensing signal $\{\mathbf{x}_j | j \mod (N_m + 1) = 1, \ 1 \le j \le N\}$ could be illustrated via $z[k, f]$.

**Model Training and Inference**. To classify the motions, the ResNet-32 [20] is adopted as the backbone for the feature extraction. It consists of one input layer, five identical residual blocks and one output layer as illustrated in Fig. 2. The input layer consists of a convolution layer and a pooling layer, and the output layer consists of a global average pooling layer and a fully-connected layer with softmax as the activation function. Each residual block consists of six layers: a batch normalization layer, a ReLu activation layer, a convolution layer, a batch normalization layer, a ReLu activation layer, and a convolution layer. The input of the ResNet is the spectrogram, and the output is the index of the estimated human motion category $\widehat{c}$.

### B. Motion Recognition Accuracy Model

It is difficult to directly derive the analytical relationship between the motion recognition accuracy and the number of sensing slots, which is denoted as $A = \Theta(N_s)$. This is because there is no analytical expression to quantify the learning performance of ResNet–32.

To address the above challenge, a promising solution is the performance regression approach proposed in [21]–[24]. Specifically, it is observed that the recognition accuracy $A = \Theta(N_s)$ is a nonlinear function of $N_s$ satisfying the following properties:

(i) $\Theta(N_s)$ is a monotonically increasing function of $N_s$;
(ii) As $N_s$ increases, the magnitude of the partial derivative $|\partial \Theta / \partial N_s|$ would gradually decrease and become zero when $N_s$ is sufficiently large, meaning that increasing the number of sensing time slots will not help wireless sensing performance at large $N_s$.

Hence the candidate parametric learning curve models to approximate the $\Theta(N_s)$ are shown in Table II (as proposed in [23]), where $\alpha, \beta, \gamma, \epsilon$ are tuning parameters.

In order to determine the tuning parameters, we first generate the dataset of sensing signals $\{r_j(t)|j \bmod (N_m + 1) = 1, \ 1 \leq j \leq N_s(N_m + 1)\}$ for $Q$ different numbers of sensing slots, denoted as $\{N_s^{(i)}| \ 1 \leq i \leq Q\}$. After the training and inferencing with DSN, the corresponding recognition accuracies are denoted as $\{A^{(i)}| \ 1 \leq i \leq Q\}$. Then the parameters $\alpha, \beta, \gamma, \epsilon$ can be calculated via the following least squares fitting,

$$\min_{\alpha, \beta, \gamma, \epsilon} \ \frac{1}{Q} \sum_{i=1}^{Q} \left| \Theta(N_s^{(i)}) - A^{(i)} \right|^2. \tag{14}$$

The above problem can be solved by brute-force search, or gradient descent method.

## V. SDP3 PERFORMANCE OPTIMIZER AND TRADEOFF ANALYSIS

The SDP3 performance optimizer, which investigates the accuracy-throughput (A-T) region of the ISAC system, is presented in this section. Denote $N_{c,k}$ as the number of time slots assigned for the k-th receiver, the throughput of k-th receiver in the whole scheduling period $R_k$ can be approximated as

$$R_k \approx \frac{N_{c,k} T_s}{T_o} \sum_{m=1}^{M} \log_2 (1 + \gamma_{k,m}), \tag{15}$$

where $\gamma_{k,m} = \mathbb{E}[\gamma_{k,j,m}] = \frac{\mathbb{E}[|H_{k,m}|^2 P_{k,j,m}]}{\sigma_z^2}$. In (15), we use $\gamma_{k,m}$ to approximate the instantaneous SNR. This is because the Doppler effect raised by human motion does not dominate the channel gain. Hence the worst communication throughput among all receivers is

$$R = \min_{k=1,\cdots,K} \frac{N_{c,k} T_s}{T_o} \sum_{m=1}^{M} \log_2 (1 + \gamma_{k,m}). \tag{16}$$

In order to characterize the A-T region, we first define the following weighted summation of recognition accuracy $A$ and worst communication throughput $R$.

$$f(w_A, w_R) = w_A A + w_R R, \tag{17}$$

where $w_A$ and $w_R$ are the weights. Given a pair of weight $(w_A, \ w_R)$, the optimal $A^*$ and $R^*$ maximizing the objective $f(w_A, w_R)$ could be obtained via the following optimization problem.

$$\mathcal{P} : (A^*, R^*) = \operatorname*{argmax}_{A,R,N_s,\mathbf{N_c}} \ f(w_A, w_R),$$

$$\text{s.t.} \quad \text{constraint in (16)}$$

$$\sum_{m=1}^{M} P_{k,m} \leq P, \ \forall k$$

$$N_s + \sum_{k=1}^{K} N_{c,k} = N. \tag{18}$$

Hence, the achievable A-T region can be characterized by $\{(A, R) = \operatorname{argmax} f(w_A, w_R)|w_A > 0, \ w_R > 0\}$. Moreover, the closed-form relationship between the recognition accuracy and communication throughput on the Pareto boundary is elaborated in the following theorem.

**THEOREM 1.** *Denote $\Theta^{-1}$ as the inverse function of $\Theta$, let (A\*, R\*) be a point on the Pareto boundary, then*

$$T_s \Theta^{-1}(A^*) + \left[ \sum_{k=1}^{K} \frac{T_o}{\sum_{\substack{m=1 \\ \gamma_{k,m} \geq \gamma_{k,0}}}^{N} \log_2 \left( \frac{\gamma_{k,m}}{\gamma_{k,0}} \right)} \right] R^* = N T_s. \tag{19}$$

*Proof.* Please refer to Appendix A. $\qquad\square$

## VI. SIMULATION AND EXPERIMENT

In this section, we verify the proposed SDP3 framework via both simulation and experiment. Specifically, a conference room with size $(3\,\mathrm{m}, \ 4.5\,\mathrm{m}, \ 3\,\mathrm{m})$ (i.e., length, width, height) is considered, where the lower left conner of the room is the origin of the coordinates and the radar is located as $(1.5, 1, 1)$. The sensing target is either an adult or a child at the location $(3, 4.2, 0)$ initially. The motions to be classified include the child/adult standing, child walking, child pacing, adult walking and adult pacing. In order to keep the consistency between the simulation and experiment, the experiment is conducted in the rooms with similar configuration. In both simulation and experiment, two carrier frequencies $f_c = 3.5\,\mathrm{GHz}$ and $60\,\mathrm{GHz}$ are considered. The total transmit power is $P = 1\,\mathrm{W}$ for both sensing and communication and the bandwidth is $B = 10\,\mathrm{MHz}$. A directional antenna with gain $P_t = 25\,\mathrm{dB}$ is used for sensing, which is directed to the sensing target. An omni-directional antenna is used for communication. The datasets of received sensing signals are generated via both experiment and simulation.

Specifically, the experiments are conducted in two scenarios and the human motion datasets are also generated shown in Fig. 5: (1) there is only one person in the room, who is the sensing target; (2) there are two persons in the room, wherein one of them is the sensing target and the other one raises sensing interference. The Scenario 1 is sensed on both 3.5GHz and 60GHz band, and the Scenario 2 is sensed on 60GHz band only. In fact, the Scenario 2 is common in practice. It refers to the situation with stronger sensing uncertainty. In the simulation dataset, we generate the samples of received sensing signal via both SDP3 data simulator and ray-tracing model, such that their performances can be compared. In both approaches, the adult and child have $B = 16$ primitives respectively. In the SDP3 data simulator, the IEEE 802.11ax/ay channels [35], [39] are adopted to generate the component $v_{k,j}(t)$ in (6). In both simulation and experiment datasets, 200 samples are generated for each motion.

### A. KL Divergence

In this part, the calibration of hyper-parameter $\rho$ for the motion of adult walking, as elaborated in Section III, is first demonstrated, which verifies the existence of sensing uncertainty.

Specifically, the KL divergence between the experimental and SDP3-simulated data samples versus the hyper-parameter $\rho$ is shown in Fig. 6(a-c) for different carrier frequencies
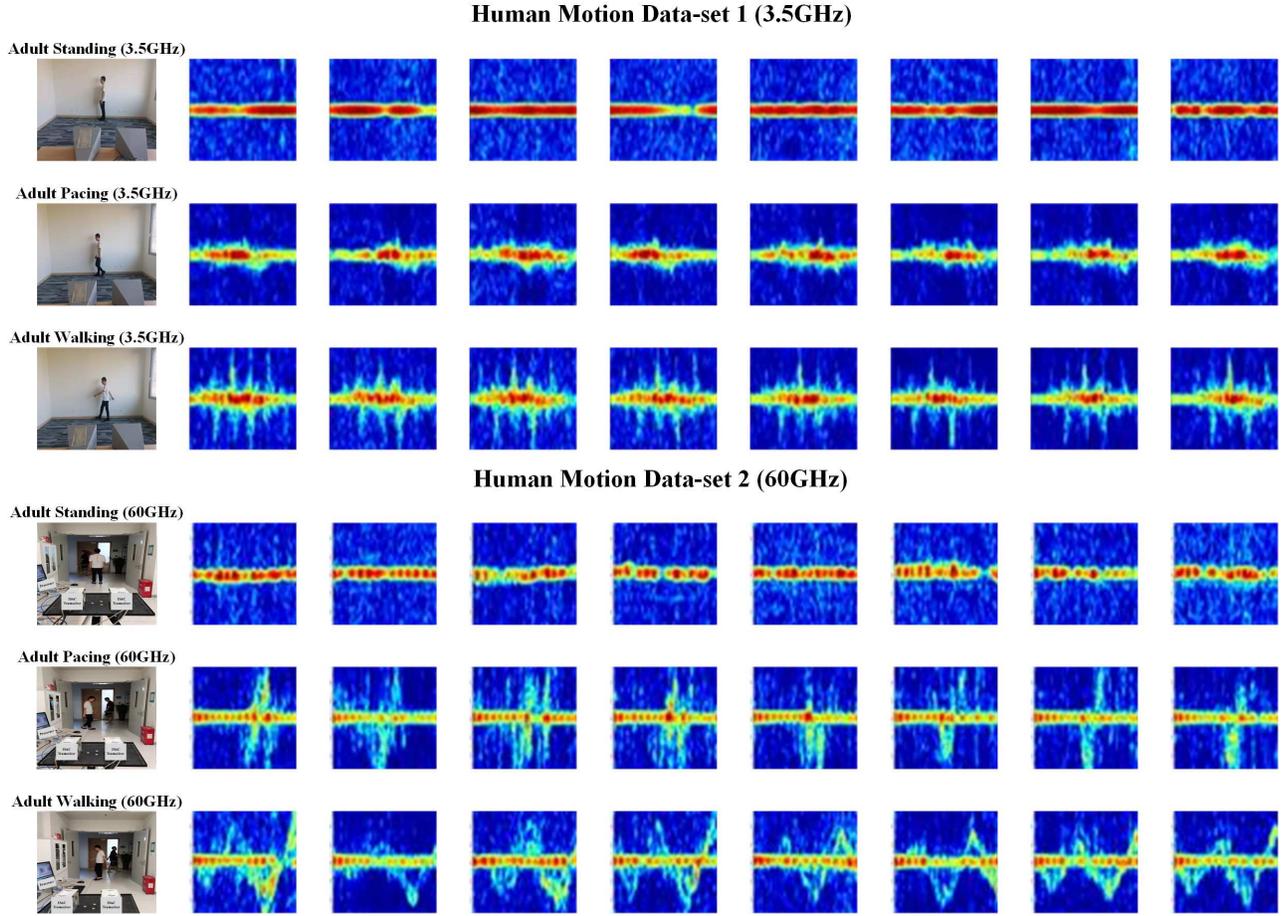
Fig. 5: The human motion datasets obtained by sensing an adult target in a conference room, where the carrier frequencies are 3.5 GHz (upper side) and 60 GHz (lower side) respectively. The scenarios with carrier frequencies 3.5GHz and 60GHz are referred to as the scenario 1 and 2 respectively. The height of the adult is 1.75m. The following three motions are tested: standing, pacing and walking. The FMCW is with 10 MHz bandwidth and 100 μs sweep time.

and scenarios. It can be seen that the KL divergence is sensitive to the value of $\rho$ of the proposed DAHC model, and the optimal values minimizing the KL divergence are $\rho^* = 0.997$, $0.997$ and $0.996$ for Scenario 1 on 3.5GHz band, Scenario 1 on 60GHz band and Scenario 2 on 60GHz band respectively. Note that the KL divergence measures the distance between two distributions, the minimum KL divergence means the best match between the experiment and the SDP3 data simulator. Note that the samples in Fig. 6(b) and (c) are obtained with the same carrier frequency but different scenarios. This demonstrates the interference from non-target person: it leads to smaller value of calibrated hyper-parameter $\rho$, thus stronger level of sensing uncertainty. In all the three figures Fig. 6(a-c), it can be observed that the KL divergence between the experiment and the ray-tracing model is constant and significantly larger than the calibrated DAHC model. This demonstrates that it is necessary to include the sensing uncertainty in the channel model to capture the potential interference from non-target motions.
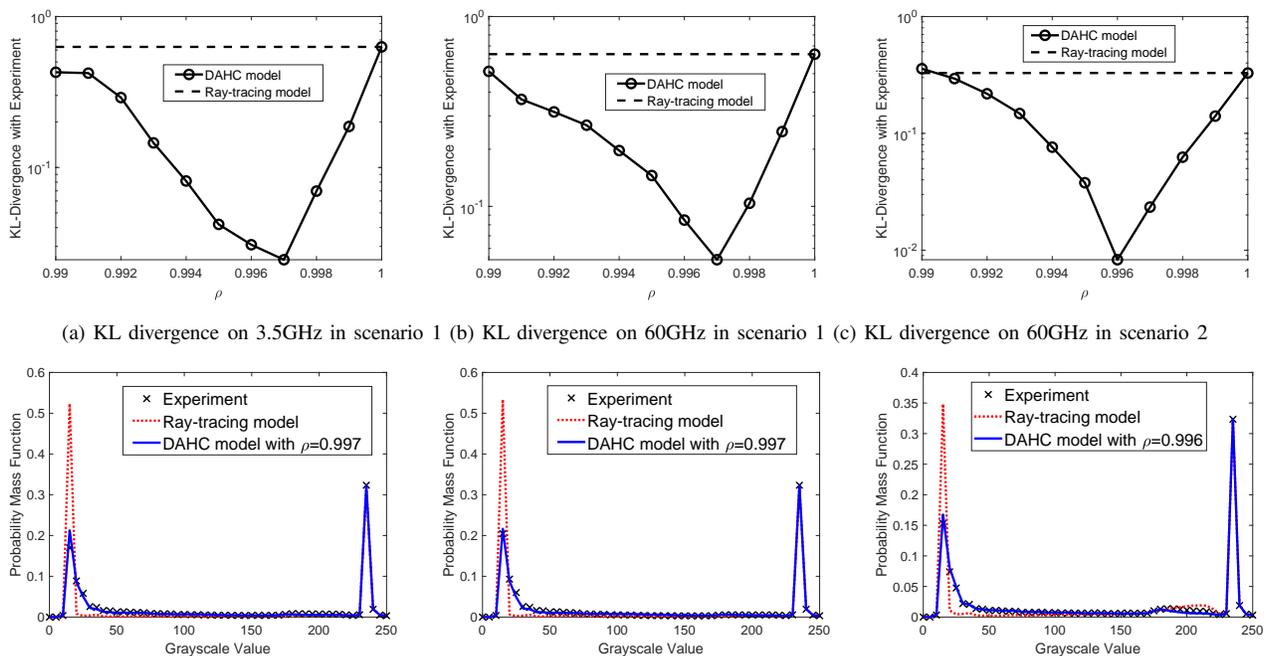
In order to further justify the superior performance of the proposed DAHC model in sensing data generation, the grayscale distribution of spectrogram is compared in Fig. 6(d-

f). Specifically, the spectrograms from experiment, ray-tracing model and DAHC model are generated for different carrier frequencies and scenarios. Then, the PMF of grayscale levels is calculated for each spectrogram. It can be observed that the grayscale PMFs from the experiment and proposed DHAC model match very well, and ray-tracing model fails to match the real experiment. This coincides with the comparison of KL divergence in Fig. 6(a-c).

*B. Recognition Accuracy*

In this part, we continue to show that the sensing dataset generated via SDP3 data simulator with the DAHC model could achieve the recognition accuracy close to the experiment dataset. The hyper-parameter $\rho$ calibrated in the above part is adopted. The simulation dataset via ray-tracing model is also investigated as a comparison.
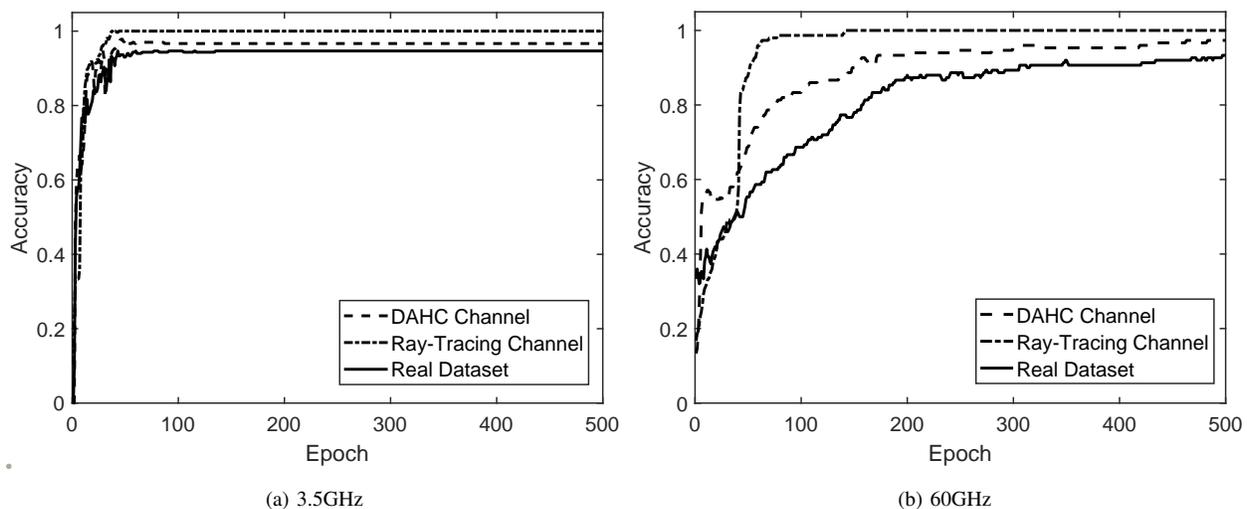
Particularly, 600 samples of sensing signals are picked up from experiment, proposed SDP3 data simulator, ray-tracing-based simulator respectively. Their spectrograms are used to train the DSN in Section IV-A respectively. The training of DSN is implemented via Momentum optimizer with a learning rate of 0.06 and a mini-batch size of 300. The recognition

(a) KL divergence on 3.5GHz in scenario 1 (b) KL divergence on 60GHz in scenario 1 (c) KL divergence on 60GHz in scenario 2

(d) PMF comparison on 3.5GHz in scenario 1 (e) PMF comparison on 60GHz in scenario 1 (f) PMF comparison on 60GHz in scenario 2

Fig. 6: The calibration results of DAHC model for walking in scenario 1 and 2.



(a) 3.5GHz

(b) 60GHz

Fig. 7: Comparison of recognition accuracy among the datasets generated by real experiments, DAHC channel model and ray-tracing model with different carrier frequency.

accuracies of the trained DSNs are tested by other samples of the datasets. As a result, the recognition accuracy versus training epoch is illustrated in Fig. 7. It is observed that compared with the ray-tracing-based simulator, the curve of the SDP3 data simulator is closer to that of experiment. For example, it is shown that the accuracy gap are 2% versus 5.33% at 3.5GHz and 4% versus 6.7% at 60GHz respectively. This demonstrates that the DAHC model can simulate the motion recognition performance in real experiment better than the ray-tracing model. Hence, the dataset generated by the SDP3 data simulator could be used for the analysis of sensing-

communication performance tradeoff. Thus, the significant effort of extensive real scenario experiment can be saved.

### C. Verification of SDP3 Performance Predictor

The performance of the SDP3 performance predictor, as proposed in Section IV-B, is demonstrated in this part, where the motion recognition accuracy versus the number of sensing slots $N_s$ is approximated analytically. We adopt the samples of received sensing signals generated by SDP3 data simulator. Both 3.5GHz and 60GHz bands are considered and five different motions, including child/adult standing, child
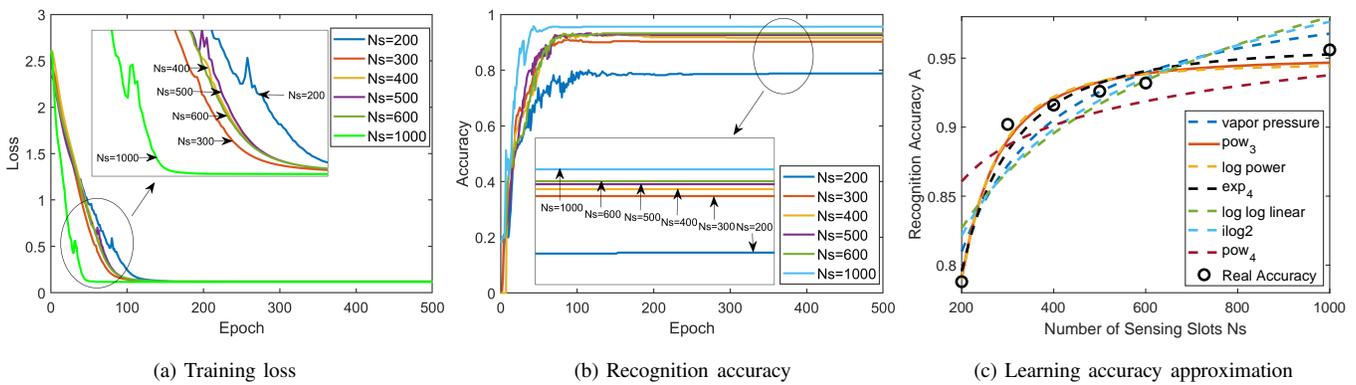
Fig. 8: The motion recognition performance of DSN and approximation result.

walking, child pacing, adult walking and adult pacing, are classified. The samples in the dataset are generated with $N_s = 1000$. The received signals of the first 200, 300, 400, 500, 600, 1000 sensing time slots are used for DSN training and accuracy testing, respectively, yielding the motion recognition accuracies for different values of $N_s$. Specifically, the training procedure is implemented via Adam optimizer with a learning rate of 0.01 and a mini-batch size of 500. The convergence of the training loss for all the values of $N_s$ is shown in Fig. 8(a). Then the trained models are tested by 500 new samples respectively, and the corresponding recognition accuracies can be obtained as shown in Fig. 8(b). It can be observed that larger $N_s$ leads to better recognition accuracy.

Hence, the curving fitting elaborated in Section IV can be proceeded. The optimized tuning parameters of all the curve models are list in Table III, and their comparison with the simulated accuracies is illustrated in Fig. 8(c). It can be observed that the models of pow$_3$, log power and exp$_4$ match the simulated accuracy very well. This is because they possess higher curvature compared with other models. With the best-fitting model pow$_3$, the tuning parameters $\alpha$, $\beta$ and $\gamma$ are obtained by minimizing the MSE as in (14), yielding $(\alpha, \beta, \gamma) = (61906, 2.4297, 0.9499)$. Thus, the relation between the recognition accuracy and number of sensing slots is approximated as

$$A = 0.9499 - 61906 \times N_s^{-2.4297}, \quad (20)$$

### D. Sensing-Communication Tradeoff

Finally, the SDP3 performance optimizer is utilized to illustrate the tradeoff between sensing and communication performance. The result with one sensing target and K=3 communication receivers is shown in Fig. 9, where the carrier frequency is 3.5GHz. Two distributions of receivers are considered: the communication receivers are at the locations (1.5,3,1), (2.5,1,1) and (0.5,1,1) for the case 1, and at (1.5,2.8,1), (2.3,1,1) and (0.7,1,1) for the case 2.

It can be observed that the accuracy-throughput (A-T) region of case 1 is a subset of that of case 2. This is because of the better communication channel in case 2. Both of two regions consist of three zones: 1) sensing saturation zone (bottom); 2) communication saturation zone (left); and 3)

TABLE III: Parametric Learning Curve Models Fitting Results

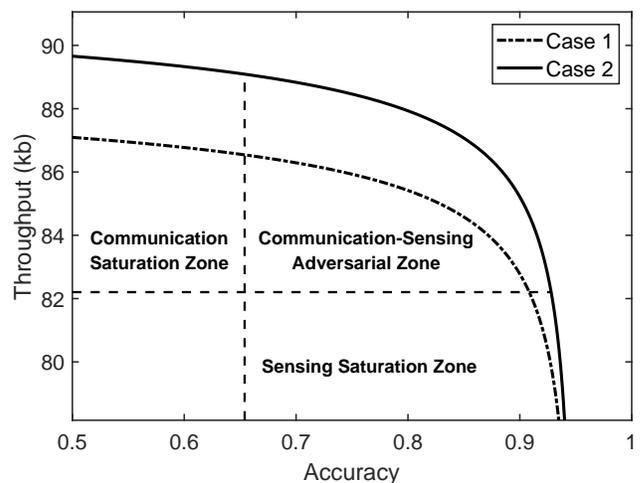| Name | Parameters | MSE |
|---|---|---|
| vapor pressure | $\alpha = 0.0117, \beta = -44.5180$ | 0.0017 |
| pow$_3$ | $\alpha = 6.1906e4, \beta = 2.4297$ $\gamma = 0.9499$ | 3.383e-4 |
| log power | $\alpha = 0.9460, \beta = 4.7438$ $\gamma = -2.9235$ | 3.7018e-4 |
| exp$_4$ | $\alpha = 2.9129, \beta = 6.9568$ $\gamma = 0.2082, \epsilon = 0.9576$ | 5.4693e-4 |
| log log linear | $\alpha = 0.2347, \beta = 1.0423$ | 0.0038 |
| ilog2 | $\alpha = 3.5228, \beta = 1.4863$ | 0.003 |
| pow$_4$ | $\alpha = 98.4911, \beta = -9.5892e3$ $\gamma = 1.2176, \epsilon = 0.1117$ | 0.0065 |



Fig. 9: The A-T region for case 1 and 2 when $K = 3$, $P = 1$ Watt, $T_o = 50\mu s$, $T_s = 250\mu s$ and $N = 4000$.

sensing-communication adversarial zone (upper right). In the sensing saturation zone, reducing the communication performance can hardly improve the sensing performance, but a

slight decrease of the sensing performance can significantly improve the communication throughput. This is because excessive wireless resources have been allocated to sensing task. The situation is the opposite in the communication saturation zone. Therefore, for a practical ISAC system, it is not desirable to enter the sensing or communication saturation zone. On the other hand, both sensing and communication performance varies sensitively with respect to each other in the sensing-communication adversarial zone. This is because the wireless resource allocation between the two functionalities is balanced. The proposed accuracy-throughput region analysis facilitates the search of the best Pareto point with heterogeneous quality-of-service requirements.

## VII. CONCLUSION

This paper proposes an SDP3 framework for the study of sensing-communication tradeoff with the particular application of human motion recognition. Specifically, the SDP3 data simulator with a data-driven hybrid channel model is proposed to generate the received sensing signals in a virtual environment. The SDP3 performance predictor is then introduced to approximate the motion recognition accuracy via analytical expression with the simulated dataset of sensing signals. Finally, the recognition accuracy and communication throughput tradeoff is characterized by the SDP3 performance optimizer. It is demonstrated that the dataset generated by the SDP3 data simulator matches the experiment dataset in KL divergence, grayscale PMF and motion recognition accuracy. Hence, the sensing-communication tradeoff can be investigated without extensive experiments. It is also shown that the sensing and communication performance is balanced in the sensing-communication adversarial zone of the A-T region, where both performance varies sensitively with respect to each other.

## APPENDIX A
## PROOF OF THEOREM 1

Appling the water-filling policy to the power constraints, the optimal worst communication throughput among all receivers is

$$R = \min_{k=1,\cdots,K} \frac{N_{c,k}T_s}{T_o} \sum_{\substack{m=1 \\ \gamma_{k,m} \geq \gamma_{k,0}}}^{N} \log_2\left(\frac{\gamma_{k,m}}{\gamma_{k,0}}\right), \quad (21)$$

where

$$\sum_{m=1}^{M} \left(\frac{1}{\gamma_{k,0}} - \frac{1}{\gamma_{k,m}}\right) = 1, \quad (22)$$

Then the $\mathcal{P}_1$ could be

$$\mathcal{P}_1 : \max_{A,R,N_s,\mathbf{N_c}} \quad (A, R),$$
$$\text{s.t.} \quad (21)$$
$$A = \Theta(C), \quad N_s + \sum_{k=1}^{K} N_{c,k} = N. \quad (23)$$

Given a fixed $N_s = N_s^*$, the optimal $\mathbf{t}$ is derived using KKT optimality conditions. Specifically, the Lagrangian of problem $\mathcal{P}$ is given by

$$L = -R^* + \sum_{k=1}^{K} \eta_k \left( R^* - \frac{N_{c,k}^* T_s}{T_o} \sum_{\substack{m=1 \\ \gamma_{k,m} \geq \gamma_{k,0}}}^{N} \log_2\left(\frac{\gamma_{k,m}}{\gamma_{k,0}}\right) \right)$$
$$+ \varrho \left( N_s^* + \sum_{k=1}^{K} N_{c,k}^* - N \right), \quad (24)$$

where $\{\eta_k \geq 0, \varrho\}$ are Lagrangian multipliers. By letting $\partial L/\partial R^* = 0$ and $\partial L/\partial N_{c,k}^* = 0$, we have

$$\sum_{k=1}^{K} \eta_k = 1, \quad -\frac{\eta_k T_s}{T_o} \sum_{\substack{m=1 \\ \gamma_{k,m} \geq \gamma_{k,0}}}^{N} \log_2\left(\frac{\gamma_{k,m}}{\gamma_{k,0}}\right) + \varrho = 0. \quad (25)$$

Now we will prove that $\eta_k \neq 0$ for any $k$ by contradiction. In particular, assume that $\eta_j = 0$ for some $j$. Putting $\eta_j = 0$ into the second equation of (25) yields $\varrho = 0$. Putting $\varrho = 0$ into the second equation of (25) with $k \neq j$ yields $\eta_k = 0$ for any $k \neq j$. Lastly, based on $\eta_j = 0$ and $\eta_k = 0$ for $k \neq j$, we have $\sum_{k=1}^{K} \eta_k = 0$. This contradicts to $\sum_{k=1}^{K} \eta_k = 1$ of the first equation of (25). Therefore, $\eta_k \neq 0$ for any $k$. Using the above result and the complementary slackness condition

$$\eta_k \left( R^* - \frac{N_{c,k}^* T_s}{T_o} \sum_{\substack{m=1 \\ \gamma_{k,m} \geq \gamma_{k,0}}}^{N} \log_2\left(\frac{\gamma_{k,m}}{\gamma_{k,0}}\right) \right) = 0, \quad (26)$$

the following equality is obtained

$$N_{c,1}^* \sum_{\substack{m=1 \\ \gamma_{1,m} \geq \gamma_{1,0}}}^{N} \log_2\left(\frac{\gamma_{1,m}}{\gamma_{1,0}}\right) = \cdots$$
$$= N_{c,K}^* \sum_{\substack{m=1 \\ \gamma_{K,m} \geq \gamma_{K,0}}}^{N} \log_2\left(\frac{\gamma_{K,m}}{\gamma_{K,0}}\right). \quad (27)$$

Combining the above result with the constraint $N_s^* + \sum_{k=1}^{K} N_{c,k}^* = N$ yields

$$N_{c,k}^* = \frac{N - N_s^*}{\left( \sum_{k=1}^{K} \frac{1}{\sum_{\substack{m=1 \\ \gamma_{k,m} \geq \gamma_{k,0}}}^{N} \log_2\left(\frac{\gamma_{k,m}}{\gamma_{k,0}}\right)} \right) \sum_{\substack{m=1 \\ \gamma_{k,m} \geq \gamma_{k,0}}}^{N} \log_2\left(\frac{\gamma_{k,m}}{\gamma_{k,0}}\right)}. \quad (28)$$

According to (21)

$$A^* = \Theta(N_s^*), \quad (29)$$

$$R^* = \frac{(N - N_s^*)T_s}{\sum_{k=1}^{K} T_o \left( \sum_{\substack{m=1 \\ \gamma_{k,m} \geq \gamma_{k,0}}}^{N} \log_2\left(\frac{\gamma_{k,m}}{\gamma_{k,0}}\right) \right)^{-1}}. \quad (30)$$

Rearranging equations (29)–(30), the proof is completed.

## REFERENCES

[1] G. Li, S. Wang, J. Li, R. Wang, X. Peng, and T. X. Han, "Wireless sensing with deep spectrogram network and primitive based autoregressive hybrid channel model," in *Proc. SPAWC*, Lucca, Italy, Sep. 2021, pp. 481–485.

[2] F. Liu *et al.*, "Integrated sensing and communications: Toward dual-functional wireless networks for 6G and beyond," *IEEE J. Sel. Areas Commun.*, vol. 40, no. 6, pp. 1728–1767, Jun. 2022.

[3] B. K. Chalise, M. G. Amin, and B. Himed, "Performance tradeoff in a unified passive radar and communications system," *IEEE Signal Process. Lett.*, vol. 24, no. 9, pp. 1275–1279, Sep. 2017.

[4] B. K. Chalise and B. Himed, "Performance tradeoff in a unified multi-static passive radar and communication system," in *Proc. (RadarConf)*, Oklahoma City, OK, USA, Apr. 2018, pp. 653–658.

[5] X. Liu, T. Huang, N. Shlezinger, Y. Liu, J. Zhou, and Y. C. Eldar, "Joint transmit beamforming for multiuser MIMO communications and MIMO radar," *IEEE Trans. Signal Process.*, vol. 68, pp. 3929–3944, 2020.

[6] F. Liu, Y. -F. Liu, A. Li, C. Masouros, and Y. C. Eldar, "Cramér-Rao bound optimization for joint radar-communication beamforming," *IEEE Trans. Signal Process.*, vol. 70, pp. 240–253, Dec. 2021.

[7] D. W. Bliss, "Cooperative radar and communications signaling: The estimation and information theory odd couple," in *Proc. (RadarConf)*, Cincinnati, OH, USA, May 2014, pp. 50–55.

[8] A. R. Chiriyath, B. Paul, G. M. Jacyna, and D. W. Bliss, "Inner bounds on performance of radar and communications co-existence," *IEEE Trans. Signal Process.*, vol. 64, no. 2, pp. 464–474, Jan. 2016.

[9] A. Aguileta, R. Brena, O. Mayora, E. Molino-Minero-Re, and L. Trejo, "Multi-sensor fusion for activity recognition: A survey," *Sensors*, vol. 19, no. 17, Sep. 2019.

[10] M. G. Amin, Y. D. Zhang, F. Ahmad, and K. C. D. Ho, "Radar signal processing for elderly fall detection: The future for in-home monitoring," *IEEE Signal Process. Mag.*, vol. 33, no. 2, pp. 71–80, Mar. 2016.

[11] Y. Kim and H. Ling, "Human activity classification based on micro doppler signatures using a support vector machine," *IEEE Trans. Geosci. Remote Sens.*, vol. 47, no. 5, pp. 1328–1337, May. 2009.

[12] Y. Kim and T. Moon, "Human detection and activity classification based on Micro-Doppler signatures using deep convolutional neural networks," *IEEE Geoscience and Remote Sensing Letters.*, vol. 13, no. 1, pp. 8-12, Jan. 2016.

[13] F. Wang, W. Gong, and J. Liu, "On spatial diversity in WiFi-based human activity recognition: A deep learning-based approach," *IEEE Internet Things J.*, vol. 6, no. 2, pp. 2035-2047, Apr. 2019.

[14] Dazhuo Wang, Jianfei Yang, Wei Cui, Lihua Xie, and Sumei Sun, "Multimodal CSI-based human activity recognition using GANs," *IEEE Internet Things J.*, vol. 8, no. 24, pp. 17345-17355, Dec. 2021.

[15] K. Niu, F. Zhang, X. Wang, Q. Lv, H. Luo, and D. Zhang, "Understanding WiFi signal frequency features for position-independent gesture sensing," *IEEE Trans. Mobile Comput.*, early access, Mar. 2021. DOI: 10.1109/TMC.2021.3063135.

[16] Y. Zhang *et al.*, "Widar3.0: Zero-effort cross-domain gesture recognition with Wi-Fi," *IEEE Trans. Pattern Anal. Mach. Intell.*, early access, Aug. 2021. DOI: 10.1109/TPAMI.2021.3105387.

[17] C. Lin, J. Hu, Y. Sun, F. Ma, L. Wang, and G. Wu, "WiAU: An accurate device-free authentication system with resnet," in *Proc. 15th Annu. IEEE Int. Conf. Sens. Commun. Netw.*, Hong Kong, China, 2018, pp. 1–9.

[18] L. Zhang, C. Wang, and D. Zhang, "Wi-pigr: Path independent gait recognition with commodity wi-fi," *IEEE Trans. Mobile Comput.*, vol. 21, no. 9, pp. 3414-3427, Sep. 2022.

[19] A. V. Oppenheim and R. W. Schafer, *Discrete-Time Signal Processing*. USA: Prentice Hall Press, 1989.

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

[21] S. Wang, Y.-C. Wu, M. Xia, R. Wang, and H. V. Poor, "Machine intelligence at the edge with learning centric power allocation," *IEEE Trans. Wireless Commun.*, vol. 19, no. 11, pp. 7293–7308, Jul. 2020.

[22] L. Zhou *et al.*, "Learning centric wireless resource allocation for edge computing: Algorithm and experiment," *IEEE Trans. Veh. Technol.*, vol. 70, no, 1, pp. 1035–1040, Jan. 2021.

[23] T. Domhan, J. T. Springenberg, F. Hutter, "Speeding up automatic hyperparameter optimization of deep neural networks by extrapolation of learning curves," in *Proc. IJCAI*, Buenos Aires, Argentina, Jul. 2015, pp. 3460–3468.

[24] M. Johnson, P. Anderson, M. Dras, and M. Steedman, "Predicting classification error on large datasets from smaller pilot data," in *Proc. ACL*, Melbourne, Australia, Jul. 2018, pp. 450–455.

[25] C.-X. Wang, J. Bian, J. Sun, W. Zhang, and M. Zhang, "A survey of 5G channel measurements and models," *IEEE Commun. Surveys Tuts.*, vol. 20, no. 4, pp. 3142–3168, 4th Quart., 2018.

[26] V. Nurmela *et al.*, METIS Channel Models, document ICT-317669/D1.4, METIS, New York, NY, USA, Jul. 2015.

[27] A. Ö. Kaya, L. J. Greenstein, and W. Trappe, "Characterizing indoor wireless channels via ray tracing combined with stochastic modeling," *IEEE Trans. Wireless Commun.*, vol. 8, no. 8, pp. 4165–4175, Aug. 2009.

[28] A. A. M. Saleh and R. Valenzuela, "A statistical model for indoor multipath propagation," *IEEE J. Sel. Areas Commun.*, vol. 5, no. 2, pp. 128–137, Feb. 1987.

[29] S. Y. Seidel and T. S. Rappaport, "Site-specific propagation prediction for wireless in-building personal communication system design," *IEEE Trans. Veh. Technol.*, vol. 43, no. 4, pp. 879–891, Nov. 1994.

[30] P. van Dorp and F. C. A. Groen, "Human walking estimation with radar," *IEE Proc. Radar Sonar Navig.*, vol. 150, no. 5, pp. 356–365, Oct. 2003.

[31] S. J. Orfanidis, *Electromagnetic Waves and Antennas*. Rutgers University, 2002.

[32] J. W. Crispin and A. L. Maffett, "Radar cross-section estimation for simple shapes," *Proc. IEEE*, vol. 53, no. 8, pp. 833–848, Aug. 1965.

[33] 3GPP, "Study on channel model for frequencies from 0.5 to 100 GHz," 3GPP TR 38.901 V16.1.0, Dec. 2019.

[34] S. Jaeckel, L. Raschkowski, K. Börner, and L. Thiele, "QuaDRiGa: A 3-D multi-cell channel model with time evolution for enabling virtual field trials," *IEEE Trans. Antennas Propag.*, vol. 62, no. 6, pp. 3242-256, Jun. 2014.

[35] A. Maltsev *et al.*, Channel Models for IEEE 802.11ay, document 802.11-15/1150r9, IEEE, New York, NY, USA, 2016.

[36] S. Rosati, G. E. Corazza, and A. Vanelli-Coralli, "OFDM Channel Estimation Based on Impulse Response Decimation: Analysis and Novel Algorithms," *IEEE Trans. Commun.*, vol. 60, no. 7, pp. 1996-2008, July 2012.

[37] M. Zhang *et al.*, "Channel models for WLAN sensing systems," IEEE 802.11 Documents, Sep. 2021. [Online]. Available: https://mentor.ieee.org/802.11/documents?is_dcn=Meihong.

[38] Y. Wang, Z. Shi, L. Huang, Z. Yu and C. Cao, "An extension of spatial channel model with spatial consistency," in *2016 IEEE 84th Vehicular Technology Conference (VTC-Fall)*, Montreal, QC, Canada, Sep. 2016, pp. 1-5.

[39] J. Liu et al., *IEEE 802.11ax Channel Model Document*. IEEE Standard 802.11-14/0882r4, Sep. 2014, pp. 1–10.

[40] R. Boulic, M. N. Thalmann, and D. Thalmann, "A global human walking model with real-time kinematic personification," *Vis. Comput.*, vol. 6, no. 6, pp. 344–358, Nov. 1990.

[41] Motion Research Laboratory, Carnegie Mellon University, http://mocap.cs.cmu.edu.