# Towards SISO Bistatic Sensing for ISAC

Zhongqin Wang, *Member, IEEE*, J. Andrew Zhang, *Senior Member, IEEE*,
Kai Wu, *Member, IEEE*, Min Xu, *Member, IEEE*, Y. Jay Guo, *Fellow, IEEE*

*Abstract*—Integrated Sensing and Communication (ISAC) is a key enabler for next-generation wireless systems. However, real-world deployment is often limited to low-cost, single-antenna transceivers. In such bistatic Single-Input Single-Output (SISO) setup, clock asynchrony introduces random phase offsets in Channel State Information (CSI), which cannot be mitigated using conventional multi-antenna methods. This work proposes WiDFS 3.0, a lightweight bistatic SISO sensing framework that enables accurate delay and Doppler estimation from distorted CSI by effectively suppressing Doppler mirroring ambiguity. It operates with only a single antenna at both the transmitter and receiver, making it suitable for low-complexity deployments. We propose a self-referencing cross-correlation (SRCC) method for SISO random phase removal and employ delay-domain beamforming to resolve Doppler ambiguity. The resulting unambiguous delay-Doppler-time features enable robust sensing with compact neural networks. Extensive experiments show that WiDFS 3.0 achieves accurate parameter estimation, with performance comparable to or even surpassing that of prior multi-antenna methods, especially in delay estimation. Validated under single- and multi-target scenarios, the extracted ambiguity-resolved features show strong sensing accuracy and generalization. For example, when deployed on the embedded-friendly MobileViT-XXS with only 1.3M parameters, WiDFS 3.0 consistently outperforms conventional features such as CSI amplitude, mirrored Doppler, and multi-receiver aggregated Doppler.

*Index Terms*—ISAC, Bistatic Sensing, SISO, CSI, Clock Asynchrony, Lightweight Neural Networks, Activity Recognition

## I. INTRODUCTION

INTEGRATED Sensing and Communication (ISAC) is rapidly emerging as a key enabler for next-generation wireless systems by embedding sensing capabilities into existing communication infrastructure [1], [2], [3], [4]. ISAC enables a wide range of applications, such as environment monitoring [5], [6], human-computer interaction [7], [8], and healthcare [9]. As ISAC moves toward real-world deployment, bistatic Single-Input Single-Output (SISO) architectures, where the transmitter and receiver are separately deployed with only one antenna each, have attracted increasing attention due to their low hardware complexity. Such setups are common in low-power IoT nodes, smart sensors, and embedded edge platforms, while still supporting reliable communication functions. They are particularly well suited for applications such as in-home gesture control, behavior monitoring in elder-care settings, and occupancy detection for energy-efficient buildings. These make SISO-based bistatic ISAC a practical and scalable solution for ubiquitous sensing.

Zhongqin Wang, J. Andrew Zhang (Corresponding Author), Kai Wu, and Min Xu are with the School of Electrical and Data Engineering, University of Technology Sydney, Sydney 2007, Australia. E-mail:{zhongqin.wang, andrew.zhang, kai.wu, min.xu}@uts.edu.au

Y. Jay Guo is with the Global Big Data Technologies Centre, University of Technology Sydney, Sydney 2007, Australia. E-mail: jay.guo@uts.edu.au

Most existing CSI-based sensing work in ISAC focuses on Single-Input Multiple-Output (SIMO) or Multiple-Input Multiple-Output (MIMO) architectures, where the receiver is typically equipped with multiple antennas. This is primarily because, in bistatic setups where the transmitter and receiver are not synchronized, the CSI is distorted by timing offset (TO) and carrier frequency offset (CFO), resulting in time-varying random phase shifts. To address this issue, prior works have proposed many compensation techniques such as Cross-Antenna Cross-Correlation (CACC) [10], [11] and Cross-Antenna Signal Ratio (CASR) [9], [12], [13]. CACC performs by computing the conjugate multiplication of CSI from two Rx antennas, while CASR achieves by dividing CSI across the two antennas. These methods have enabled a range of sensing applications, including target tracking [14], [15], [16], activity recognition [17], [18], [19], and even water-level monitoring [20]. Recent advances in deep learning have given rise to a variety of data-driven sensing frameworks. Models leveraging attention mechanisms [21], domain adaptation [22], and spatiotemporal feature extraction [23] have significantly improved performance across diverse ISAC tasks.

Meanwhile, communication standardization is also progressing rapidly. The IEEE 802.11bf amendment introduces native sensing capabilities into Wi-Fi for CSI-based sensing [24]. This initiative represents a key milestone toward enabling commercial devices to support wireless sensing, facilitating the transition of ISAC technologies from research prototypes to real-world deployment. While existing research has established a strong foundation under multi-antenna settings, further efforts are still required to enable robust sensing in practical, low-complexity configurations such as bistatic SISO systems.

However, towards SISO bistatic sensing for ISAC, several challenges remain to be addressed:

*1) CSI random phase removal under SISO constraints.* In multi-antenna systems, the TO and CFO caused by clock asynchrony are identical across Rx antennas, enabling efficient methods such as CACC and CASR to fully eliminate phase offsets. However, these methods are not applicable in SISO configurations. The research on SISO-based sensing remains limited. Several methods [18], [25] typically assume strong line-of-sight (LoS) conditions, which may not hold in practical deployments, or involve high computational costs, making them unsuitable for embedded platforms. Therefore, it is essential to clean raw CSI without relying on multi-antenna diversity or ideal propagation assumptions.

*2) Doppler mirror ambiguity.* In ISAC systems, the limited bandwidth of communication signals, for example in LTE-based systems where it ranges from 1.4 MHz to 20 MHz, results in coarse range resolution. As a result, Doppler features play a crucial role in capturing target motion for accurate

sensing. However, many existing systems suffer from Doppler mirror ambiguity, where the energy is nearly symmetric across positive and negative Doppler frequencies, making it difficult to determine the true direction of motion. This commonly arises in CACC-based systems, which inherently induces symmetric signal components. While CASR avoids this symmetry, it is primarily designed for single-target scenarios and may introduce nonlinear distortions in feature extraction. The Doppler mirroring is often overlooked in prior work.

*3) Learning robust representations from wireless signals.* Unlike image and language tasks where abundant, high-quality datasets enable end-to-end learning directly from raw inputs, wireless sensing faces unique challenges. Due to multipath propagation, environmental variability, hardware diversity, and differing frequency bands, the received signals vary significantly across settings. Collecting diverse and large-scale labeled datasets is costly and challenging. While some works [26], [23] employ simulation data for training, their effectiveness in real-world deployment remains uncertain. As a result, directly feeding raw CSI into deep networks often leads to poor generalization. Moreover, common data augmentation techniques used in vision such as rotation or cropping lack theoretical validation in the context of wireless signals. Therefore, designing interpretable signal features and augmentation strategies is critical for enhancing model generalization.

In this work, we propose *WiDFS 3.0*, a real-time bistatic sensing system for SISO-based ISAC tasks. It features a separately deployed transmitter and receiver, each with only a single antenna, providing a low-cost ISAC solution. WiDFS 3.0 enables precise estimation of delay, Doppler, and ambiguity-resolved micro-Doppler under the impact of clock asynchrony. Based on the extracted high-quality features, lightweight neural networks deployable on embedded platforms can achieve strong sensing performance. WiDFS 3.0 also supports extension to multi-antenna setups for angular estimation and is compatible with different operating frequencies and bandwidths.

At first, to mitigate the impact of TO and CFO, we construct an energy-adjusted reference CSI from the original CSI itself, called *Self-Referencing Cross-Correlation (SRCC)*. Specifically, we apply an inverse fast Fourier transform (IFFT) to transform the CSI into the delay domain. A Gaussian window is then applied to strengthen the dominant path (a superposition of multiple paths) while suppressing weaker multipath components. The windowed result is transformed back to the frequency domain via Fourier transform (FFT) to reconstruct a new CSI. On this basis, cross-correlation between the raw and reference CSI is performed to remove TO and CFO while preserving the linear structure of delay and Doppler features, facilitating subsequent feature extraction.

Second, we adopt a beamforming-based approach to suppress Doppler mirroring. We first aggregate multiple consecutive CSI measurements to form a coherent processing interval (CPI) [1]. Within each CPI, we apply Capon beamforming [27] to estimate weights for each delay bin, suppress dynamic by-product noise, and extract beamformed Doppler signatures to construct a 2D delay-Doppler map. By concatenating the map across multiple CPIs, we construct an unambiguous delay-Doppler-time tensor for downstream sensing tasks.

Third, we use embedded-friendly lightweight neural networks like MobileViT-XXS [28] by inputting our unambiguous features to achieve accurate and robust sensing. Due to the coarse range resolution in communication systems, we compress the delay dimension to extract Doppler-time features for fine-grained sensing in single-target scenarios. In contrast, for multi-target tasks, preserving the delay dimension enables spatial discrimination among targets. We also introduce a Doppler-guided data augmentation strategy grounded in physical principles, providing meaningful variations during training.

Our main contributions are highlighted as follows:

*1)* We propose a low-complexity SISO SRCC based on energy-adjusted CSI for random phase removal in bistatic systems, without relying on ideal propagation assumptions.

*2)* We develop a lightweight delay-domain beamforming approach to suppress dynamic by-product noise and resolve Doppler ambiguity during feature extraction.

*3)* We introduce a Doppler-guided data augmentation strategy grounded in physical motion principles, improving model robustness and generalization.

*4)* We conduct extensive experiments validating the efficiency and effectiveness of WiDFS 3.0.

- On an edge platform (Raspberry Pi 4B), our unambiguous delay and Doppler feature extraction runs in just 8.5 milliseconds, demonstrating real-time capability.
- Despite using a single antenna at both ends, WiDFS 3.0 achieves a 2.05 m median range error (20 MHz, without smoothing), outperforming prior multi-antenna systems by 0.5–2 m, while offering Doppler mirror suppression performance comparable to multi-antenna approaches.
- The Doppler ambiguity-resolved features enable accurate and generalizable sensing. A compact model like MobileViT-XXS (1.3M parameters) achieves an F1 score of 0.928–0.938 on the single-target Widar 3.0 dataset [17], significantly outperforming the classical BVP baseline (0.849–0.859). On the multi-target WiMANS dataset [29], it achieves an F1 score of 0.659 for behavior recognition and a people counting accuracy of 0.629, surpassing the commonly-used CSI amplitude baseline.

## II. RELATED WORK

This section reviews representative methods for random phase removal, feature extraction, and data-driven learning.

### A. CSI Random Phase Removal

Many techniques [30] have been proposed to eliminate the CSI random phase offsets in bistatic sensing, typically categorized by their use of spatial (antenna), spectral (subcarrier), or temporal (Doppler) dimensions. One widely used approach is CACC [10], which requires two Rx antennas at least. By leveraging the fact that clock-induced phase distortions are identical across antennas, CACC can removes TO and CFO while preserving the linear relationships among delay, Doppler, and Angle of Arrival (AoA). However, the conjugate operation in CACC introduces symmetry in signal representations. Specifically, Doppler spectra often exhibit nearly identical energy at positive and negative frequencies,

making it difficult to infer true motion direction. To mitigate this, WiDFS [11] proposes Differential CACC (DCACC), which applies linear transformations across different antenna pairs to suppress mirror ambiguity, but it still requires three Rx antennas. Other multi-antenna works [31] exploit subspace methods to extract dominant AoA and build reference signals. Another popular method is CASR [9], [12], [13], which computes the CSI ratio between two Rx antennas to cancel out random phase offsets along with automatic gain control (AGC) variations. CASR is effective in Doppler estimation for single-target scenarios. However, its nonlinear formulation complicates accurate estimation of delay and AoA.

Recently, several compensation methods have started exploring SISO configurations by constructing reference signals in the spectral domain across subcarriers. SHARP [18] employs compressed sensing to recover a reference signal, but its performance depends on a strong LoS path, and its computational cost limits real-time applicability. The work [32] still relies on the presence of a LoS path to extract relative TO and CFO for CSI compensation. And cross-frequency cross-correlation (CFCC) [25] also requires a strong LoS condition to function effectively. In addition, linear regression-based methods [33] aim to estimate relative TO and CFO, though their effectiveness degrades in multipath environments due to nonlinear phase variations. In contrast, our SISO SRCC constructs an energy-adjusted CSI, sharing the same TO and CFO as the original CSI, without relying on LoS assumptions.

### B. Multi-Dimension Signal Feature Extraction

ISAC systems often extract motion-related features across multiple domains, such as Doppler (temporal), delay (spectral), and AoA (spatial). Many existing works extract only Doppler features to perform tasks such as human tracking and behavior sensing. For instance, PITrack [34] estimates target trajectories using Doppler shifts observed across multiple receivers. However, such methods are prone to accumulated trajectory drift over time. Other methods [35] focus solely on AoA estimation across multiple Rx antennas to localize the target, but tend to lose robustness in the presence of low spatial diversity.

More advanced approaches adopt joint estimation of Doppler, delay, and AoA, requiring only a receiver equipped with multiple Rx antennas. For example, Widar2.0 [10] and mD-Track [36] apply maximum likelihood estimation to jointly infer these parameters. While such optimization-based frameworks can achieve high accuracy, they typically incur high computational overhead and are sensitive to initialization. Recent works have explored alternatives such as mirrored-MUSIC [37] to reduce the complexity of multidimensional estimation, or tensor decomposition techniques [38]. WiDFS [11] leverages DCACC to suppress both random phase distortions and Doppler ambiguity, and then uses MUSIC-based high-resolution Doppler estimation as a basis to infer delay and AoA for precise target localization. WiDFS 2.0 [16] further enhances this pipeline by identifying multiple latent Doppler components from a tracked target, followed by delay and AoA estimation for each component. This method improves tracking accuracy. In this work, we focus on the SISO bistatic setup,

where only Doppler and delay dimensions are available. We introduce a lightweight delay-domain beamforming method to achieve unambiguous delay and Doppler estimation.

### C. Data-driven Sensing Techniques

Data-driven techniques are playing an increasingly important role in ISAC applications. Inspired by the success of end-to-end learning in vision and natural language processing, many studies [19], [39] directly feed raw or minimally processed CSI such as amplitude, real, or imaginary components into deep neural networks. These approaches focus on network design and aim to learn feature representations implicitly. While they often achieve high sensing accuracy in specific conditions, generalization across diverse environments remains a key challenge. To improve robustness, many work focuses on first removing random phase distortions and then extracting Doppler spectrograms for network input. For instance, SHARP [18] mitigates TO and CFO, applies short-time Fourier transform (STFT) to generate Doppler features, and uses a custom lightweight CNN for behavior recognition. Widar 3.0 [17] aligns multi-receiver Doppler profiles to form the BVP feature, which is then processed by a sequence modeling network.

Recently, large language models (LLMs) have also been explored for CSI understanding, leveraging their powerful capacity to model complex spatiotemporal relationships in multipath channels. Some works [23], [26] train LLMs using simulated CSI to overcome the scarcity of labeled real-world data. While promising, the practical effectiveness of these LLM-based methods remains to be fully validated. In this work, we extract physically grounded delay-Doppler-time features under SISO bistatic settings, where Doppler mirror ambiguity is explicitly suppressed. These interpretable features not only support lightweight neural network inference on embedded platforms, but also provide a strong foundation for generalizable and data-efficient sensing.

### III. SYSTEM OVERVIEW

In this work, we propose a lightweight SISO bistatic sensing scheme *WiDFS 3.0* for ISAC applications. The system is designed for real-world deployment with separately deployed transmitter and receiver, each equipped with only a single antenna, offering a low-cost and low-complexity solution. Despite the SISO constraint and clock asynchrony, WiDFS 3.0 achieves accurate delay, Doppler, and ambiguity-resolved micro-Doppler estimation. Moreover, the framework is extensible to Single-Input Multiple-Output (SIMO), Multiple-Input Multiple-Output (MIMO), and Multiple-Input Single-Output (MISO) configurations to further exploit Angle of Departure (AoD) and AoA information for enhanced sensing. The overall workflow of WiDFS 3.0 is illustrated in Fig. 1, which consists of the following key modules:

*1) Phase Compensation via CSI Reconstruction:* To mitigate the effects of TO and CFO caused by clock asynchrony, we construct an energy-adjusted reference CSI along the subcarrier dimension. Sharing the same TO and CFO as the original CSI, this reference enables effective phase cancellation through conjugate multiplication while preserving both temporal and spectral linear structures.
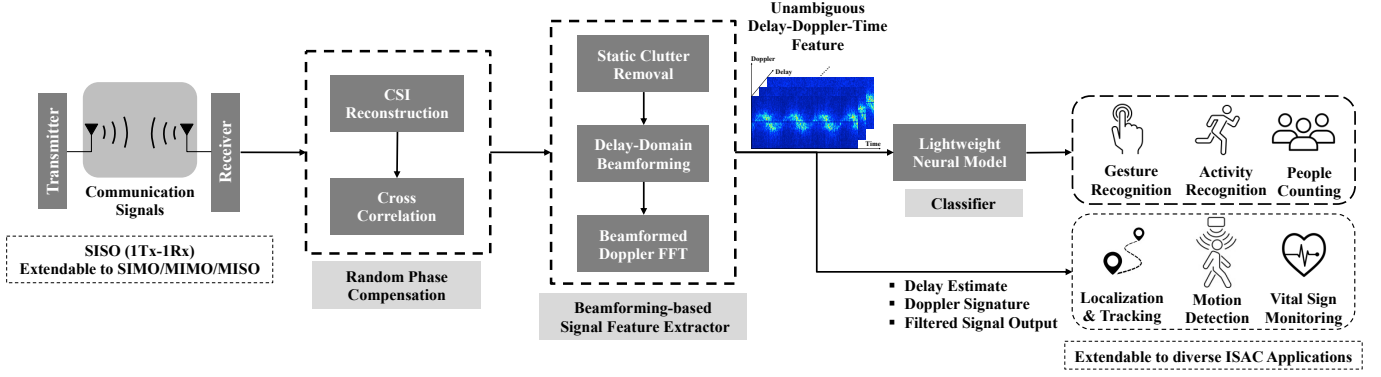
**Fig. 1:** System overview of WiDFS 3.0: A SISO bistatic sensing framework for ISAC.

*2) Beamforming-based Signal Feature Extractor:* We adopt a delay-domain beamforming approach to suppress mirror ambiguity and dynamic noise in Doppler estimation. By aggregating across multiple CPIs, we obtain an unambiguous delay-Doppler-time tensor for downstream sensing tasks.

*3) Lightweight Neural Network Classifier:* The unambiguous features are fed into off-the-shelf compact neural networks, such as MobileViT-XXS, which enable accurate activity recognition and are well-suited for embedded deployment. In addition, a Doppler-guided data augmentation strategy is introduced to enhance model generalization.

## IV. PHASE COMPENSATION VIA CSI RECONSTRUCTION

In this section, we present a CSI model in a SISO bistatic setup and reconstruct CSI values for random phase removal.

### A. SISO Bistatic CSI Model

In SISO bistatic systems, the transmitter and receiver are spatially separated and operate without a shared clock. Each side is equipped with a single antenna. Let $CSI_{i,j}$ denote the CSI measured at the $i$-th subcarrier and the $j$-th OFDM symbol (CSI sampling time). Within the short-duration CPI (around 100 milliseconds), the transmit power, automatic gain control settings, and the amplitude of each signal are typically assumed to be constant [1]. The CSI measurements are $\mathbf{csi} \in \mathbb{C}^{N \times M}$, where $N$ is the number of subcarriers and $M$ is the number of OFDM symbols. The CSI is modeled as:

$$CSI_{i,j} = \underbrace{e^{-\boldsymbol{J}\left(2\pi f_i \tau_j^{\text{TO}} + \phi_j^{\text{CFO}}\right)} e^{-\boldsymbol{J}\phi^h}}_{\text{Imperfect Signal Processing}} \left(H_i^S + H_{i,j}^X\right), \quad (1)$$

where

$$\begin{cases} H_i^S = \sum_{l_1} \rho_i^S[l_1] e^{-\boldsymbol{J}2\pi f_i \tau^S[l_1]} \\ H_{i,j}^X = \sum_{l_2} \rho_i^X[l_2] e^{-\boldsymbol{J}2\pi \left(f_i \tau^X[l_2] + f^D[l_2](j-1)\Delta t\right)} \end{cases} \quad (2)$$

These variables are explained as follows:

1) $\tau_j^{\text{TO}}$ and $\phi_j^{\text{CFO}}$ denote the random TO and CFO, caused by clock asynchrony between the transmitter and receiver.

2) $\phi^h$ denotes a hardware-induced phase offset, initialized by the phase-locked loop at power-up. Due to the randomness of the initial phase, it may vary across different power cycles.

3) $H_i^S$ represents the channel frequency response (CFR) associated with static paths. These include both the direct line-of-sight (LOS) or non-line-of-sight (NLOS) path between the transmitter and receiver, as well as other reflections from static environmental objects such as walls, floors, and furniture. Each static path is characterized by an attenuation $\rho_i^S[l_1]$ and a propagation delay $\tau^S[l_1]$. $f_i$ is the subcarrier frequency.

4) $H_{i,j}^X$ is the CFR contributed by reflections from moving objects. Each dynamic path has an attenuation $\rho_i^X[l_2]$, a delay $\tau^X[l_2]$, and a Doppler frequency shift (DFS) $f^D[l_2]$ caused by object motion. $\Delta t$ is the OFDM symbol interval.

### B. CSI Reconstruction via Delay-Domain Windowing

For the $j$-th OFDM symbol, we reconstruct the CSI by utilizing the frequency domain across subcarriers. The goal is to enhance the power of the dominant propagation path while suppressing weaker paths in the reconstructed CSI. Here, the dominant path typically represents a composite of closely-spaced strong reflections due to the limited delay resolution. To achieve that, we first transform the original CSI from the frequency domain to the time domain using IFFT. The time-domain Channel Impulse Response (CIR) is given by:

$$h_j(\tau) = \mathcal{F}^{-1}\{CSI_{i,j}\} = \sum_{i=0}^{N-1} CSI_{i,j} e^{\boldsymbol{j}2\pi f_i \tau}, \quad (3)$$

where $\tau$ denotes a discrete delay bin.

A Gaussian window function $\mathcal{G}(\tau)$ is then applied, centered at the peak energy bin $\tau_j^{\text{peak}}$, so the windowed CIR is

$$h_j'(\tau) = \mathcal{G}\left(\tau - \tau_j^{\text{peak}}\right) \cdot h_j(\tau), \quad (4)$$

where

$$\mathcal{G}(\tau - \hat{\tau}_j) = e^{-\left(\frac{\tau - \tau_j^{\text{peak}}}{2\sigma}\right)^2}. \quad (5)$$

The parameter $\sigma$ controls the width of the Gaussian window.

After that, we transform the windowed CIR back to the frequency domain using FFT to obtain the reconstructed CSI. This process is expressed as:

$$\mathcal{CSI}_{i,j} = \mathcal{F}\{h_j'(\tau)\} = \sum_{\tau=0}^{N-1} h_j'(\tau) e^{-\boldsymbol{j}2\pi f_i \tau}. \quad (6)$$

Compared to the original CSI, the reconstructed version offers an energy-adjusted representation in which the dominant path is enhanced while the power of secondary paths is effectively suppressed. Note that, overly small $\sigma$ may introduce spectral distortion in the frequency domain due to severe filtering, which may introduce additional noise in CSI reconstruction. Its impact will be analysed in the following.

### C. Random Phase Removal via Cross-Correlation

To removal the impact of TO and CFO, we perform cross correlation between the original and reconstructed CSI:

$$
\begin{aligned}
\Delta CSI_{i,j} &= CSI_{i,j}\overline{\mathcal{CSI}}_{i,j} \\
&= \left(H_i^S + H_{i,j}^X\right)\left(\overline{\mathcal{H}}_i^S + \overline{\mathcal{H}}_{i,j}^X\right) \\
&= \underbrace{H_i^S\overline{\mathcal{H}}_i^S}_{\text{static}} + \underbrace{H_i^S\overline{\mathcal{H}}_{i,j}^X + \overline{\mathcal{H}}_i^S H_{i,j}^X}_{\text{dynamic}} + \underbrace{H_{i,j}^X\overline{\mathcal{H}}_{i,j}^X}_{\text{by-product}}.
\end{aligned} \tag{7}
$$

This random phase removal process is referred to as *Self-Referencing Cross-Correlation (SRCC)*. It can simultaneously mitigate the impacts of clock asynchronism and hardware mismatch. When extended to multi-antenna systems, our SRCC can eliminate the need for complex antenna-specific correction in AoA estimation. It also preserves the linear relationship of signal features such as Doppler and delay.

In addition, Eq. (7) contains four cross-correlation terms. *(1)* $H_i^S\overline{\mathcal{H}}_i^S$ represents the static component, capturing the signal contribution between the fixed transmitter and receiver, as well as reflections from static objects in the environment. *(2)* $H_i^S\overline{\mathcal{H}}_{i,j}^X$ and $\overline{\mathcal{H}}_i^S H_{i,j}^X$ are cross-terms that encode dynamic variations of interest. However, they are complex conjugates of each other, exhibiting symmetric structures in delay and Doppler domains. Specifically, the DFS may result in nearly equal energy at positive and negative frequencies, making it difficult to determine the true direction of movement. In our SRCC, the windowing function can enhance the energy of $H_i^S$ (typically containing dominant paths) while suppressing $H_{i,j}^X$ (which includes weaker dynamic components). As a result, the windowed terms often satisfy $\|H_i^S\overline{\mathcal{H}}_{i,j}^X\| < \|\overline{\mathcal{H}}_i^S H_{i,j}^X\|$, which helps mitigate mirror ambiguity to some extent. However, narrower windows may amplify phase noise. Therefore, it is necessary to apply an appropriate Gaussian window and further exploit the delay domain to effectively suppress mirror ambiguity. *(3)* $H_{i,j}^X\overline{\mathcal{H}}_{i,j}^X$ is a by-product term that generally acts as noise and should be suppressed. Most prior works [11], [16], [40] simply neglect this term, assuming its energy is negligible compared with other components. However, this assumption does not always hold in practice, especially under weak LoS or NLoS conditions where dynamic multipath reflections may retain significant energy. In such cases, the by-product term can generate considerable interference, leading to degraded sensing performance.

### D. Impact of Windowing on CSI Reconstruction Accuracy

To quantify the impact of delay-domain windowing in CSI reconstruction, we analysis the phase variance of the reconstructed CSI across subcarriers using the Cramér-Rao Lower Bound (CRLB). Assuming additive Gaussian noise and unbiased estimation, the CRLB for estimating the phase $\phi_{i,j}$ of the reconstructed CSI at the $i$-th subcarrier and the $j$-th OFDM symbol is inversely proportional to the energy of the windowed CIR (see Appendix A for derivation):

$$
\mathrm{Var}\left(\phi_{i,j}\right) \geq \frac{\eta^2}{\left\|\mathcal{G}\left(\tau - \tau_j^{\text{peak}}\right)\cdot h_j\left(\tau\right)\right\|^2}, \tag{8}
$$

where $\eta^2$ is noise power and $\|\cdot\|^2$ is to compute signal power.

A wider window can reduce the phase noise in the reconstructed CSI by preserving more signal energy, while it introduces additional multipath components and increases ambiguity in the cross-correlation results. In contrast, a narrower window suppresses side lobes and reduces such ambiguity, but at the cost of higher phase noise due to reduced signal energy. We evaluate the practical impact of this trade-off in Sec. VIII.

## V. BEAMFORMING-BASED SIGNAL FEATURE EXTRACTOR

While deep learning networks are capable of learning rich representations, the generated features are often difficult to interpret. In this work, we adopt a signal-model-driven approach to extract interpretable features. Specifically, we employ a beamforming-based method to generate a delay-Doppler-time feature that captures motion dynamics under SISO bistatic setups, while maintaining low computational complexity.

### A. Dynamic Component Separation

Within each CPI window, we compute the average of the cross-correlation CSI $\Delta H_{i,j}$ along the time dimension at each subcarrier. The average is denoted as $\mathcal{U}_i$, given by:

$$
\mathcal{U}_i = \frac{1}{M}\sum_{j=1}^{M}\Delta CSI_{i,j} \approx H_i^S\overline{\mathcal{H}}_i^S, \tag{9}
$$

We then subtract the static term from $\Delta CSI_{i,j}$ to obtain the dynamic component:

$$
\mathcal{V}_{i,j} = \Delta CSI_{i,j} - \mathcal{U}_i \approx H_i^S\overline{\mathcal{H}}_{i,j}^X + \overline{\mathcal{H}}_i^S H_{i,j}^X + H_{i,j}^X\overline{\mathcal{H}}_{i,j}^X. \tag{10}
$$

This operation can suppress static clutter and highlight motion-induced changes. To enhance motion sensitivity, we normalize the dynamic component $\mathcal{V}_{i,j}$ by the static term $\mathcal{U}_i$ to get

$$
\mathcal{W}_{i,j} = \frac{\mathcal{V}_{i,j}}{\mathcal{U}_i} \approx \frac{H_{i,j}^X}{H_i^S} + \frac{\overline{\mathcal{H}}_{i,j}^X}{\overline{\mathcal{H}}_i^S} + \frac{H_{i,j}^X}{H_i^S}\frac{\overline{\mathcal{H}}_{i,j}^X}{\overline{\mathcal{H}}_i^S}. \tag{11}
$$

According to [11], [13], this ratio can amplify motion variations relative to the static background, effectively increasing the signal-to-noise ratio (SNR) for dynamic components.

### B. Delay-Domain Beamforming with By-Product Suppression

We first apply the MVDR (Minimum Variance Distortionless Response) algorithm [27] to estimate the delay profile, yielding a set of beamforming weights for each delay bin.

*1) Observation Matrix.* As shown in Eq. (11), $\mathcal{W}_{i,j}$ contains two conjugate-mirrored dynamic components, i.e., $H_{i,j}^X/H_i^S$ and $\overline{\mathcal{H}}_{i,j}^X/\overline{\mathcal{H}}_i^S$, whereas the by-product term does not exhibit

such symmetry. To preserve this conjugate relationship in the beamforming process, we construct the observation matrix by concatenating both the original $\mathcal{W}_{i,j}$ and its conjugate counterpart $\overline{\mathcal{W}}_{i,j}$. This joint representation enables the MVDR beamformer to fully exploit the phase correlation between the mirrored dynamic components, ensuring coherent alignment of target-related phases while attenuating uncorrelated interference, including the by-product term. The observation matrix $\boldsymbol{\Lambda} \in \mathbb{C}^{N \times 2M}$ is formulated as:

$$\boldsymbol{\Lambda} = \begin{bmatrix} \mathcal{W}_{1,1} & \cdots & \mathcal{W}_{1,M} & \overline{\mathcal{W}}_{1,1} & \cdots & \overline{\mathcal{W}}_{1,M} \\ \mathcal{W}_{2,1} & \cdots & \mathcal{W}_{2,M} & \overline{\mathcal{W}}_{2,1} & \cdots & \overline{\mathcal{W}}_{2,M} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \mathcal{W}_{N,1} & \cdots & \mathcal{W}_{N,M} & \overline{\mathcal{W}}_{N,1} & \cdots & \overline{\mathcal{W}}_{N,M} \end{bmatrix}_{N \times 2M}.$$

(12)

*2) Steering Matrix.* A 2D steering matrix $\boldsymbol{A} \in \mathbb{C}^{N \times L}$ is then constructed, where each steering vector $\mathbf{a}(\Delta\tau_l)$ corresponds to a delay value $\Delta\tau_l$. Here, $\Delta\tau = \tau^X - \tau^S$ denotes the relative delay between a dynamic path and a static path. The steering matrix is expressed as:

$$\boldsymbol{A} = [\mathbf{a}(\Delta\tau_1), \mathbf{a}(\Delta\tau_2), \cdots, \mathbf{a}(\Delta\tau_L)]_{N \times L}, \quad (13)$$

where

$$\mathbf{a}(\Delta\tau) = \left[ e^{-j2\pi f_1 \Delta\tau}, e^{-j2\pi f_2 \Delta\tau}, \cdots, e^{-j2\pi f_N \Delta\tau} \right]^{\mathsf{T}}. \quad (14)$$

Since the observation matrix already incorporates both the original and conjugate dynamic components, the resulting delay-domain response becomes conjugate-symmetric about zero delay. This symmetry makes the delay profile on one side of the axis fully redundant with the other, allowing the beamforming search to be restricted to a single side (either $\tau \geq 0$ or $\tau \leq 0$) without loss of information. Here, we set the delay search range to $\tau \geq 0$.

*3) Covariance Matrix.* To obtain a more robust estimation of the covariance matrix, we apply forward-backward smoothing,

$$\widetilde{\boldsymbol{R}}_{\boldsymbol{\Lambda}} = \boldsymbol{R}_{\boldsymbol{\Lambda}} + \boldsymbol{J}\boldsymbol{R}_{\boldsymbol{\Lambda}}\boldsymbol{J} + \epsilon\boldsymbol{I}, \quad (15)$$

where $\boldsymbol{R}_{\boldsymbol{\Lambda}} = \boldsymbol{\Lambda}\boldsymbol{\Lambda}^{\mathsf{H}}$ is the covariance matrix, $\boldsymbol{J}$ is the exchange matrix of size $N \times N$, obtained by flipping the identity matrix from left to right, and $\epsilon\boldsymbol{I}$ is a small regularization term added to improve stability. The forward-backward smoothing can improve the effective rank of the matrix and strengthen the separation between the signal and noise subspaces [41].

*4) Beamforming Weights.* The beamforming weight vector $\mathbf{w}(\Delta\tau) \in \mathbb{C}^{N \times 1}$ for the delay bin $\Delta\tau$ can be computed as:

$$\mathbf{w}(\Delta\tau) = \frac{\widetilde{\boldsymbol{R}}_{\boldsymbol{\Lambda}}^{-1} \mathbf{a}(\Delta\tau)}{\mathbf{a}^{\mathsf{H}}(\Delta\tau) \widetilde{\boldsymbol{R}}_{\boldsymbol{\Lambda}}^{-1} \mathbf{a}(\Delta\tau)}, \quad (16)$$

which can suppress noise and interference from all other delays while preserving the amplitude and phase of the signal at the target delay $\Delta\tau$. The corresponding beamforming weights are then applied separately to the original and conjugate parts of the observation matrix to get

$$\begin{cases} \boldsymbol{X}^{\text{orig}}(\Delta\tau) = \mathbf{w}^{\mathsf{H}}(\Delta\tau) \boldsymbol{\Lambda}^{\text{orig}} \\ \boldsymbol{X}^{\text{conj}}(\Delta\tau) = \mathbf{w}^{\mathsf{H}}(\Delta\tau) \boldsymbol{\Lambda}^{\text{conj}} \end{cases}, \quad (17)$$

where $\boldsymbol{\Lambda}^{\text{orig}} \in \mathbb{C}^{N \times M}$ and $\boldsymbol{\Lambda}^{\text{conj}} = \overline{\boldsymbol{\Lambda}}^{\text{orig}} \in \mathbb{C}^{N \times M}$ are the first and last $M$ columns of the matrix $\boldsymbol{\Lambda} \in \mathbb{C}^{N \times 2M}$. The resulting outputs, $\boldsymbol{X}^{\text{orig}}(\Delta\tau) \in \mathbb{C}^M$ and $\boldsymbol{X}^{\text{conj}}(\Delta\tau) \in \mathbb{C}^M$, coherently enhance target-related signals due to their preserved conjugate symmetry, whereas the by-product term, lacking such symmetry, suffers phase misalignment between its two components and is thus suppressed.

*5) Beamformed Signal.* The final beamformed signal for Doppler processing is obtained by summing the two outputs:

$$\boldsymbol{X}(\Delta\tau) = \boldsymbol{X}^{\text{orig}}(\Delta\tau) + \boldsymbol{X}^{\text{conj}}(\Delta\tau). \quad (18)$$

Each element corresponding to the $j$-th OFDM symbol can be expressed as:

$$X_j(\Delta\tau) = \sum_l \rho[l] e^{-\boldsymbol{J}\left(2\pi f^D[l](j-1)\Delta t + \phi[l]\right)} + \mathcal{N}_j \quad (19)$$

where there may contain multiple targets at the delay bin, and $\mathcal{N}_j$ represents residual noise.

### C. Beamformed Doppler FFT

To reveal fine-grained motion patterns, such as limb swings, quick gestures, or subtle posture transitions, we perform a Doppler analysis on the beamformed signal at each delay bin $\Delta\tau$. We apply the FFT along the temporal dimension (i.e., across OFDM symbols) within a CPI as follows:

$$\begin{aligned} Y\left(f^D \mid \Delta\tau\right) &= \mathcal{F}\{X_j(\Delta\tau)\} \\ &= \sum_{j=0}^{M} X_j(\Delta\tau) e^{-\boldsymbol{j}2\pi f^D(j-1)\Delta t}. \end{aligned} \quad (20)$$

Repeating this process across delay bins yields a 2D delay-Doppler spectrum, where delay offers coarse spatial separation and Doppler captures motion velocity and periodicity. This joint spectrum enables fine-grained motion characterization.

### D. Spatiotemporal Feature Representation across CPIs

To capture continuous motion patterns, we extend the delay-Doppler spectrum over multiple CPIs. Specifically, a delay-Doppler frame is computed per CPI and stacked along the temporal axis to form a 3D tensor:

$$\boldsymbol{\mathcal{S}} \in \mathbb{R}^{L_{\text{Delay}} \times L_{\text{Doppler}} \times L_{\text{CPI}}}, \quad (21)$$

where each element $\mathcal{S}(\Delta\tau, f^D, \ell)$ in $\boldsymbol{\mathcal{S}}$ denotes the delay-Doppler-time feature at the delay bin $\Delta\tau$, DFS $f^D$, and CPI frame index $\ell$. Here, $L_{\text{delay}}$ is the number of delay bins, $L_{\text{doppler}}$ is the number of Doppler bins, and $L_{\text{CPI}}$ is the number of consecutive CPIs. Such a structured feature enables downstream neural models to recognize complex motion patterns.

### E. Algorithm Complexity Analysis

We analyze the computational complexity of the proposed delay-Doppler-time feature extraction pipeline.

*1) CSI Reconstruction:* Each OFDM symbol requires an IFFT of size $N$, yielding $\mathcal{O}(N \log N)$ per symbol and $\mathcal{O}(MN \log N)$ per CPI.

*2) Static Component Removal:* Temporal averaging and subtraction over the $N \times M$ matrix leads to $\mathcal{O}(NM)$ complexity.

*3) Delay-domain Beamforming:* For each delay bin, MVDR weight computation involves matrix inversion with cost $\mathcal{O}\left(N^3\right)$, and applying weights adds $\mathcal{O}\left(NM\right)$. For $L_{\text{Delay}}$ bins, the total cost is $\mathcal{O}\left(L_{\text{Delay}}N^3 + L_{\text{Delay}}NM\right)$.

*4) Doppler FFT:* Performing an FFT of size $M$ for each delay bin yields $\mathcal{O}\left(L_{\text{Delay}}M \log M\right)$.

The overall complexity is dominated by MVDR beamforming and can be approximated as $\mathcal{O}\left(L_{\text{Delay}}N^3\right)$. The number of delay bins is typically small due to limited bandwidth, keeping the computational load low. For example, based on Intel 5300 CSI data, our algorithm requires only 8.5 ms per CPI on a Raspberry Pi without applying any code-level optimization (see Section VIII for details).

## VI. LIGHTWEIGHT NETWORK AS CLASSIFIER

The proposed interpretable and unambiguous features are robust to variations in transceiver position, environment, hardware, and individual differences, which simplifies the learning task for the recognition network. Accordingly, a lightweight neural network suitable for deployment on embedded devices is employed to perform efficient classification.

### A. Input Representation

*1) Single-Target Scenario.* For a single moving target, we compress the feature over the delay dimension. This yields a micro-Doppler signature that preserves Doppler-time characteristics while discarding spatial variation,

$$\mathcal{S}\left(f^D, \ell\right) = \sum_{k=1}^{L_{\text{Delay}}} \mathcal{S}\left(\Delta\tau_k, f^D, \ell\right). \tag{22}$$

The resulting representation is a 2D Doppler-time map:

$$\boldsymbol{\mathcal{S}}' \in \mathbb{R}^{L_{\text{Doppler}} \times L_{\text{CPI}}}, \tag{23}$$

which simplifies the model input and emphasizes motion periodicity and velocity variation over time.

*2) Multiple-Target Scenario.* In the presence of multiple moving targets, directly ignoring the delay dimension may obscure spatially distinct motion patterns. Instead, we retain delay bins as separate spatial channels, forming the 3D tensor $\boldsymbol{\mathcal{S}}$ defined in Eq. (21). This representation allows the model to capture simultaneous motions at different distances and improves discrimination in crowded environments.

### B. Exploring Doppler Characteristics for Data Augmentation

In a bistatic system, the Doppler velocity $v^D$ is determined by the geometry relationship among the transmitter, receiver, and the moving target,

$$v^D = \frac{cf^D}{f_c} = \boldsymbol{v} \cdot \left(\frac{\boldsymbol{P} - \boldsymbol{P}_{\text{Tx}}}{\|\boldsymbol{P} - \boldsymbol{P}_{\text{Tx}}\|} + \frac{\boldsymbol{P} - \boldsymbol{P}_{\text{Rx}}}{\|\boldsymbol{P} - \boldsymbol{P}_{\text{Rx}}\|}\right), \tag{24}$$

where $\boldsymbol{v} = (v_x, v_y)$ is the 2D velocity vector of the moving target on the x-y plane, $\boldsymbol{P} = (x, y)$ is the corresponding 2D position, and $\boldsymbol{P}_{\text{Tx}} = (x_{\text{Tx}}, y_{\text{Tx}})$ and $\boldsymbol{P}_{\text{Rx}} = (x_{\text{Rx}}, y_{\text{Rx}})$ are the 2D coordinates of the transmitter and receiver, respectively. When a global coordinate system is defined with the transmitter placed at the origin, i.e., $\boldsymbol{P}_{\text{Tx}} = (0, 0)$, all positions

and velocities are expressed relative to this reference point. In the following, we leverage this relationship to propose several physically meaningful data augmentation strategies to enhance the robustness of our classifier. For multi-target scenarios, the same augmentation scheme is applied to the Doppler-time features across all delay bins.

*1) Location.* Changing the target position $\boldsymbol{P}$ alters the geometric projection vectors toward the transmitter and receiver, which causes a displacement along the Doppler dimension. For example, consider $\boldsymbol{v} = (0, 1)$ m/s, $\boldsymbol{P}_{\text{Tx}} = (0, 0)$ m, and $\boldsymbol{P}_{\text{Rx}} = (4, 0)$ m. When the target is at $\boldsymbol{P} = (2, 2)$ m, the Doppler velocity is $v^D \approx 1.41$ m/s, while at $\boldsymbol{P} = (3, 2)$ m, it increases slightly to $v^D \approx 1.45$ m/s. To simulate such location-induced variations, we apply Doppler-axis translations and affine transformations (e.g., stretching or scaling).

*2) Orientation.* Changing the direction of the velocity vector $\boldsymbol{v}$ while keeping its magnitude and the target position $\boldsymbol{P}$ fixed can alter the projection onto the transmitter-receiver axis. For example, with $\boldsymbol{P} = (2, 2)$ m, $\boldsymbol{P}_{\text{Tx}} = (0, 0)$ m, and $\boldsymbol{P}_{\text{Rx}} = (4, 0)$ m, consider two velocity vectors with the same magnitude $\|\boldsymbol{v}\| = 1$ m/s: upward motion $\boldsymbol{v} = (0, 1)$ m/s yields a projected Doppler velocity of $v^D \approx 1.41$ m/s, while diagonal motion $\boldsymbol{v} = \left(\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2}\right)$ m/s results in $v^D \approx 1.00$ m/s. These differences can be emulated via Doppler-axis translation, stretching, or mirroring to simulate reversed motion.

*3) Speed.* Scaling the magnitude of the velocity vector $\boldsymbol{v}$ directly influences the Doppler shift, since $v^D \propto \|\boldsymbol{v}\|$. Increasing the target's speed results in a broader Doppler spread, while decreasing it compresses the Doppler range. To simulate such effects, we scale the Doppler-time spectrogram along the Doppler dimension accordingly.

*4) Timing Variation.* Varying the starting time of a gesture or activity shifts the entire Doppler-time spectrogram along the temporal axis, without altering the underlying Doppler content. During augmentation, such variation can be emulated by applying temporal shifts along the time axis.

*5) Noise Injection.* In practice, Doppler measurements are affected by abnormal noise and multipath interference. To enhance the model robustness, we simulate these effects by injecting Gaussian noise into the Doppler-time spectrogram.

### C. Lightweight Network Selection

Many real-time ISAC sensing tasks require not only accurate motion recognition but also low computational overhead to enable deployment on resource-constrained edge platforms such as wireless routers and embedded IoT devices. While numerous customized neural architectures [42], [43] have been proposed to enhance recognition performance by increasing model complexity and parameter count, such designs often come at the cost of computational efficiency. In this work, we adopt several representative lightweight backbone networks, including MobileViT [28], MobileNetV2 [44], SqueezeNet [45], ShuffleNet [46], and shallow CNN, all of which are well-suited for on-device deployment. For comparison, we also include ResNet18 [47] and a simple MLP as baselines. Specifically, the MLP consists of two fully connected layers with ReLU activation and dropout, sharing the extracted features for

classification. The CNN comprises three convolutional blocks with pooling, followed by flattening and classification heads. It is worth noting that if these lightweight models can achieve strong sensing performance, it is reasonable to expect further improvements with stronger networks.

### D. Loss Function

We adopt different loss functions for single-target and multi-target activity recognition tasks.

*1) Single-Target Scenario.* For single-target activity recognition, we adopt the categorical cross-entropy (CE) loss:

$$\mathcal{L}_{\text{single}} = \text{CE}\left(\hat{\boldsymbol{y}}, \boldsymbol{y}\right) = -\sum_{j=1}^{C} y_j \log\left(\hat{y}_j\right), \qquad (25)$$

where $C$ is the number of activity classes, $\boldsymbol{y} \in \{0,1\}^C$ is the one-hot ground-truth label, and $\hat{\boldsymbol{y}} \in [0,1]^C$ is the predicted probability distribution over activity classes. The model uses a 2D Doppler-Time tensor $\mathcal{S}'$ as input.

*2) Multi-Target Scenario.* We formulate multi-target activity recognition as a joint task comprising two components: (1) a $K$-class multi-label classification task for activity recognition, and (2) a $(C+1)$-class classification task for people counting (including the case when no person is present), where $K$ denotes the number of predefined activity categories, and $C$ denotes the maximum number of individuals. Each input sample may contain multiple individuals performing different activities simultaneously. The ground truth consists of a binary vector $\boldsymbol{y} \in \{0,1\}^K$ indicating the presence of each activity, and a scalar count label $\boldsymbol{c} \in \{0,1,\ldots,C\}$ representing the number of active individuals. The model takes a 3D Delay-Doppler-Time tensor $\mathcal{S}$ as input. The total loss is:

$$\mathcal{L}_{\text{multi}} = \frac{1}{K}\sum_{j=1}^{K}\text{BCE}\left(\hat{y}_j, y_j\right) + \lambda_c \text{CE}\left(\hat{\boldsymbol{c}}, \boldsymbol{c}\right), \qquad (26)$$

where $\text{BCE}\left(\cdot\right)$ denotes binary cross-entropy for each activity label, $\text{CE}\left(\cdot\right)$ is the CE loss for people count, and $\lambda_c$ is a hyperparameter. At inference, a fixed threshold of 0.5 is applied to each $\hat{y}_j$ to determine the presence of activity $j$, while the predicted number of individuals is obtained by selecting the class with the highest confidence in $\hat{\boldsymbol{c}}$.

## VII. IMPLEMENTATION

### A. CSI Datasets

*1) Object Tracking Datasets:* We validate the effectiveness of our SISO feature extraction using WiFi and LTE object tracking datasets, collected in indoor environments under LOS conditions with a transmitter-receiver separation of approximately 2 m. The target moves at a speed of 1 m/s.

- *WiFi:* The WiFi dataset from WiDFS [11] and WiDFS2.0 [16] is collected using Intel 5300 NICs over a 20 MHz bandwidth at 5 GHz. CSI from 30 subcarriers is sampled at 1 kHz using the Linux 802.11n CSI Tool [48]. A single target walks repeatedly along elliptical, V-shaped, and rectangular trajectories, with ground truth provided by a millimeter-wave radar sensor.

- *LTE:* The 3.1 GHz LTE dataset [14] is collected using an NI Massive MIMO testbed serving as the base station (BS) and a USRP device acting as the user equipment (UE). Pilot streaming is implemented via LabVIEW Communications 2.0. The extracted CSI contains 100 active subcarriers with an effective frequency resolution of 180 kHz (i.e., resource block-level spacing) over a 20 MHz bandwidth. A person repeatedly walks back and forth along a linear trajectory, facing the tranceiver side.

*2) Activity Recognition Datasets:* To evaluate the effectiveness of our SISO features for activity classification, we adopt single-target and multi-target activity recognition datasets, both collected using Intel 5300 NICs.

- *Single-Target:* The Widar3.0 dataset [17] is collected at 5 GHz using multiple receivers, each equipped with three antennas. *Dataset 1* contains six human-computer interaction gestures, including Clap (25,039 samples), Draw-O (26,543 samples), Draw-Zigzag (34,046 samples), Push & Pull (25,799 samples), Slide (29,248 samples), and Sweep (25,049 samples), performed by 16 participants across 5 environments and 5 orientations, totalling over 1,65,724 samples. *Dataset 2* contains numeric drawing gestures (0-9), collected across 5 locations and 5 orientations by 2 users. Each class contains 3,900 samples, leading to a total of 39,000 samples.

- *Multi-Target:* The WiMANS dataset [29] is a multi-user WiFi CSI dataset captured using a 3Tx-3Rx antenna configuration at both 2.4 GHz and 5 GHz. CSI is sampled at 1 kHz, with each 3-second segment containing 3000 CSI frames in a $3 \times 3 \times 30$ format each. It covers 9 daily activities, including *Nothing, Walking, Rotation, Jumping, Waving, Lying Down, Picking Up, Sitting Down, and Standing Up*, performed by 6 users across 3 environments and 5 locations, totaling 11,286 samples.

### B. SISO Feature Extraction Parameters

We use the following parameters for feature extraction:
- *CPI:* Each CPI has 128 CSI samples (i.e., OFDM symbols), with a sliding step of 32 between adjacent CPIs.
- *CSI Reconstruction:* The number of IFFT bins is 128. A Gaussian window with $\sigma = 64$ is applied.
- *Delay Beamforming:* MVDR beamforming is applied over a $[0, 32)$ m delay range with 1 m resolution, producing 32 bins. For single- and multi-target activity recognition datasets, the resolution is reduced to 4 m per bin, yielding 8 bins to limit data and computation.
- *Weighted Doppler FFT:* A Doppler FFT of size 128 is applied to each delay bin, covering a Doppler frequency range of $[-150, 150]$ Hz. For 5 GHz signals, this corresponds to a Doppler velocity range of $[-8, 8]$ m/s.

### C. Model Training Parameter

All models are trained for 512 epochs with a batch size of 128 and an initial learning rate of 0.001 using the Adam optimizer. A step scheduler halves the learning rate every 196 epochs. Input features are Z-score normalized, and the
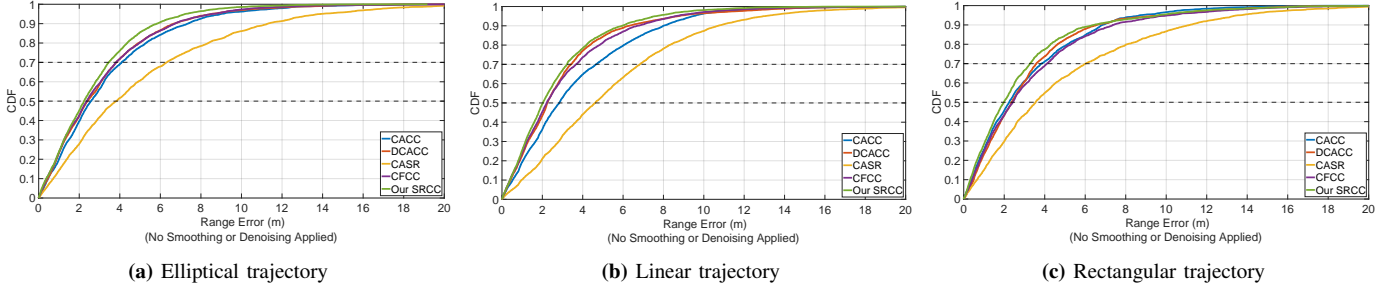
**(a)** Elliptical trajectory      **(b)** Linear trajectory      **(c)** Rectangular trajectory

**Fig. 2:** CDF of range estimation error across three trajectories.

hyperparameter $\lambda_c$ in Eq. (26) is set to 0.2. For the single-target dataset, a 70%-30% train-test split is used, while an 80%-20% split is adopted for the multi-target dataset. All features are rescaled to a fixed temporal length of 64 via interpolation, downsampling longer and upsampling shorter sequences. In the delay-Doppler-time representation, the delay dimension is treated as the channel dimension for batch training. Models using delay-Doppler-time or Doppler-time features employ the augmentation strategies in Section VI.B. To address slight class imbalance in *Widar 3.0 Dataset 1*, class weighting is applied to the cross-entropy loss, with weights inversely proportional to class frequencies.

### D. Baselines

*1) Delay-Doppler Feature Extraction:* We compare our beamforming-based WiDFS 3.0 with 1Tx-1Rx SRCC *(hereafter referred to as SRCC for simplicity in the following experiments)* with other approaches, all of which mitigate TO and CFO using lightweight signal processing techniques:

- *CACC (1Tx-2Rx)* [10]: A two-antenna method using conjugate multiplication between Rx antennas.
- *DCACC (1Tx-3Rx)* [11]: A three-antenna method that first applies CACC to eliminate random phase offsets, followed by a differential CACC (DCACC) transformation across three Rx antennas to remove Doppler ambiguity.
- *CASR (1Tx-2Rx)* [9], [12]: A two-antenna method based on CSI ratio between Rx antennas.
- *CFCC (1Tx-1Rx)* [25]: A single-antenna method based on frequency-domain correlation across subcarriers.

All the above baseline methods estimate the target's delay and Doppler velocity using 2D FFT, without smoothing or post-processing, in order to obtain raw measurements.

*2) Activity Recognition:* We evaluate activity recognition performance using both delay-Doppler-time (multi-target scenario) and Doppler-time (single-target scenario) features. To further benchmark performance, we also include the BVP feature [17], derived from a multi-receiver setup. In addition, we assess the effectiveness of compact neural network architectures. For the single-target scenario, we report accuracy, precision, recall, and F1-score. For the multi-target scenario, we use class-wise precision, recall, and F1-score to evaluate per-class performance under more complex conditions.

## VIII. EXPERIMENTAL RESULTS

This section presents a comprehensive evaluation of WiDFS 3.0 in terms of feature quality and classification performance.

**TABLE I:** Range (i.e., delay) estimation error (m) at 50% and 70% percentiles. Our SRCC achieves the lowest range estimation error across all trajectories. *No smoothing or denoising is applied.*

| Percentile | Track | CACC | DCACC | CASR | CFCC | Our SRCC |
|---|---|---|---|---|---|---|
| 50% | Ellipse | 2.57 m | 2.40 m | *3.83 m* | 2.37 m | **2.16 m** |
| | Linear | 2.84 m | 2.25 m | *4.61 m* | 2.21 m | **2.02 m** |
| | Rectangle | 2.24 m | 2.36 m | *3.55 m* | 2.42 m | **1.98 m** |
| 70% | Ellipse | 4.10 m | 3.82 m | *6.30 m* | 3.82 m | **3.42 m** |
| | Linear | 4.67 m | 3.42 m | *6.83 m* | 3.66 m | **3.22 m** |
| | Rectangle | 3.84 m | 3.59 m | *6.02 m* | 4.12 m | **3.27 m** |

### A. Delay-Doppler Feature Extraction Quality

*1) Delay:* Fig. 2 and Table. I report the raw range estimation error (computed by converting delay to distance using the speed of light for intuitive interpretation) at the 50% and 70% percentiles across three trajectories. Fig. 3 visualizes the raw range variations for a segment of the elliptical trajectory. Our WiDFS 3.0 with the 1Tx-1Rx SRCC method achieves the lowest error, with the median errors of 2.16 m, 2.02 m, and 1.98 m for the ellipse, linear, and rectangle trajectories, respectively. Even at the 70% percentile, our SRCC maintains the lowest estimation errors with 3.42 m, 3.22 m, and 3.27 m, respectively. In contrast, the 1Tx-2Rx CASR exhibits the highest estimation errors due to nonlinear distortions across subcarriers. The 1Tx-2Rx CACC suffers from limited mirror suppression. The 1Tx-3Rx DCACC and 1Tx-1Rx CFCC have similar levels of accuracy by using linear transformations without introducing nonlinear artifacts. Overall, our beamforming-based pipeline can achieve better noise suppression.

*2) Doppler:* Fig. 4 illustrates the extracted Doppler profiles for a single moving target. Except for the 1Tx-2Rx CACC, other methods show similar Doppler values. The CACC suffers from strong Doppler mirroring, where the power at $+f^D$ and $-f^D$ is nearly symmetric, making it challenging to distinguish the correct Doppler bin. Additionally, the antenna pair selection for CACC makes it highly sensitive: when the chosen antenna pair has large signal energy asymmetry, doppler mirror suppression can be effective; however, in cases with small energy differences, the mirror problem become severe. Even with delay-domain filtering, the suppression capability remains limited due to the inherently narrow bandwidth. In contrast, 1Tx-3Rx DCACC and 1Tx-2Rx CASR utilize spatial domain across Rx antennas to achieve more robust mirror suppression. Our 1Tx-1Rx SRCC and CFCC can improve Doppler clarity by incorporating delay-domain filtering.
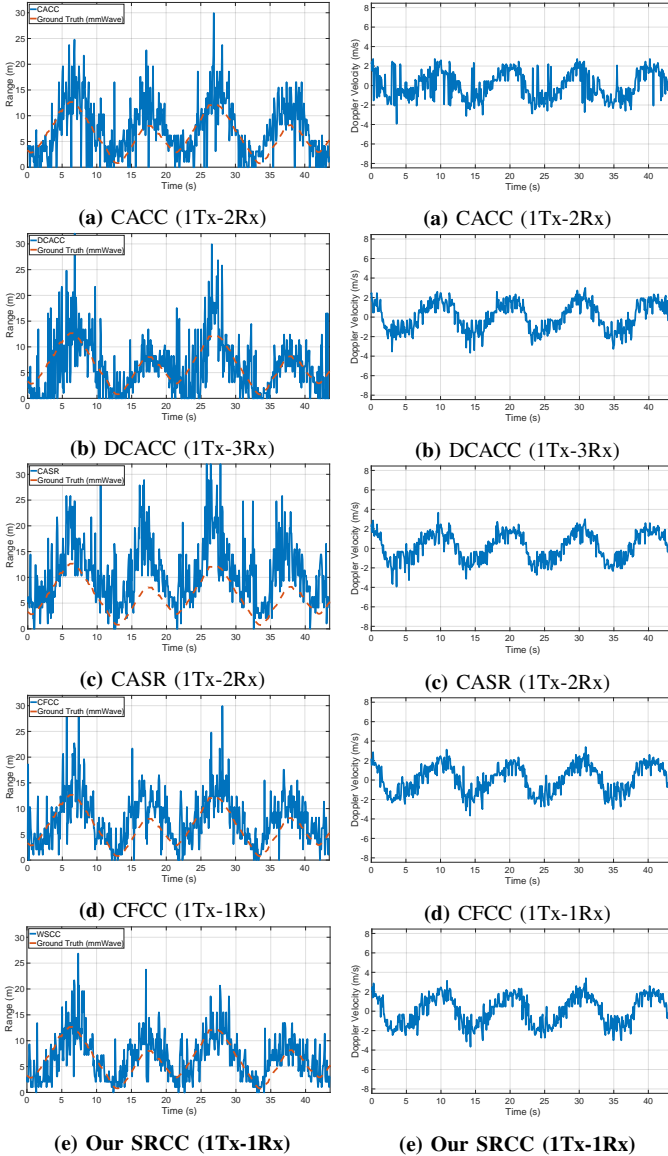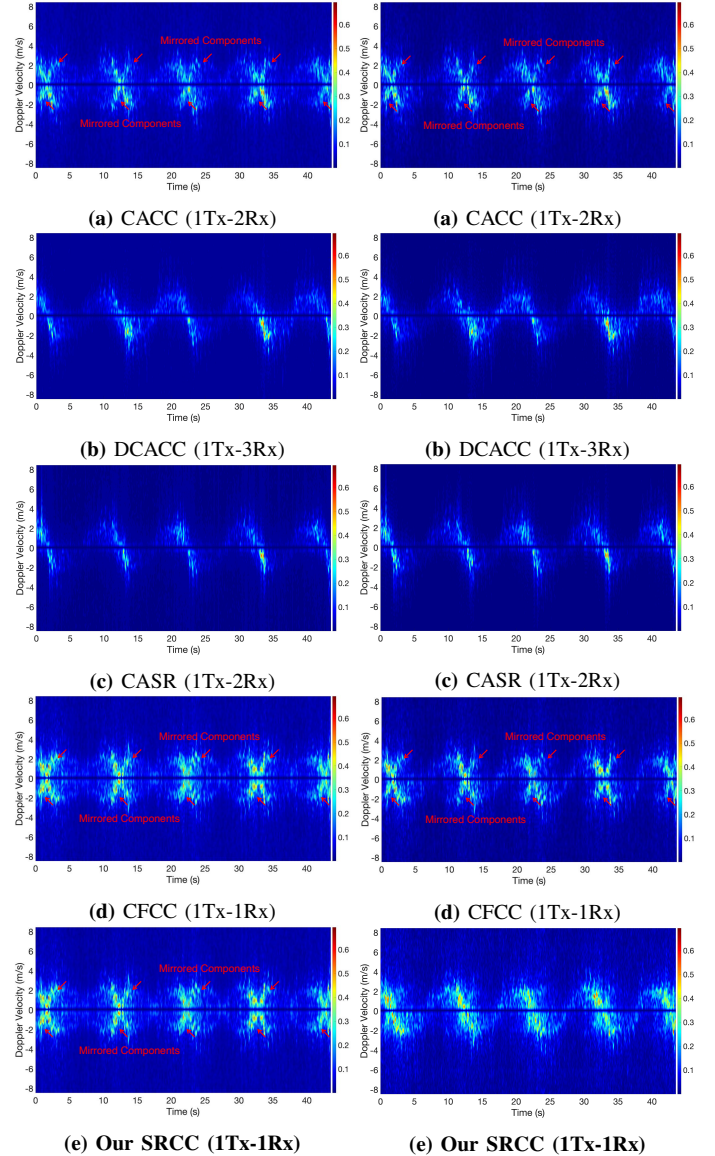
**(a)** CACC (1Tx-2Rx)



**(a)** CACC (1Tx-2Rx)



**(a)** CACC (1Tx-2Rx)



**(a)** CACC (1Tx-2Rx)



**(b)** DCACC (1Tx-3Rx)



**(b)** DCACC (1Tx-3Rx)



**(b)** DCACC (1Tx-3Rx)



**(b)** DCACC (1Tx-3Rx)



**(c)** CASR (1Tx-2Rx)



**(c)** CASR (1Tx-2Rx)



**(c)** CASR (1Tx-2Rx)



**(c)** CASR (1Tx-2Rx)



**(d)** CFCC (1Tx-1Rx)



**(d)** CFCC (1Tx-1Rx)



**(d)** CFCC (1Tx-1Rx)



**(d)** CFCC (1Tx-1Rx)



**(e) Our SRCC (1Tx-1Rx)**



**(e) Our SRCC (1Tx-1Rx)**



**(e) Our SRCC (1Tx-1Rx)**



**(e) Our SRCC (1Tx-1Rx)**

**Fig. 3:** Range  **Fig. 4:** Doppler (w. delay)  **Fig. 5:** Micro-Doppler (w.o. delay)  **Fig. 6:** Micro-Doppler (w. delay)

*3) Micro-Doppler:* Fig. 5 and Fig. 6 present the micro-Doppler signatures extracted without (w.o.) and with (w.) delay-domain filtering, respectively. Without delay filtering, obvious mirrored Doppler components appear in CACC, as well as in 1Tx-1RX CFCC and SRCC. These symmetric artifacts around the zero-Doppler axis severely compromise the interpretation of target dynamics. After applying delay-domain filtering, the three methods exhibit improved suppression of the mirrored components. However, CACC and CFCC still retain residual mirrored energy, whereas our SRCC yields significantly cleaner micro-Doppler signatures. This demonstrates the effectiveness of our delay-aligned beamforming in isolating motion-induced Doppler features and eliminating aliased components. For multi-antenna DCACC and CASR, both leverage spatial domain information to suppress Doppler mirroring, producing clean spectrograms without symmetric distortions. However, CASR suffers from nonlinear distortions in the subcarrier domain, resulting in degraded range estimation. In addition, Fig. 7 compares 1Tx-1Rx CFCC and SRCC

on LTE data at 3.1 GHz. Compared to 5 GHz WiFi, the LTE signals have a longer wavelength, resulting in smaller Doppler shifts and lower sensing sensitivity. Consistent with the previous WiFi-based results, both methods initially exhibit severe Doppler mirroring when delay filtering is not applied. With delay filtering, CFCC achieves partial suppression, while our SRCC achieves much stronger attenuation of mirrored components, demonstrating robustness across different signals.

In summary, DCACC effectively suppresses Doppler mirroring and enables accurate range estimation in multi-antenna systems, but AoA estimation requires complex antenna calibration [11]. In contrast, our SISO SRCC delivers comparable Doppler quality while inherently avoiding inter-antenna phase inconsistencies caused by clock asynchrony and hardware diversity, making it well-suited for compact, low-cost, and calibration-free deployments.

*4) Impact of Gaussian Window Width:* Fig. 8 and Table. II show the impact of the Gaussian window width $\sigma$ on delay and Doppler estimation. The smaller $\sigma$ facilitates the
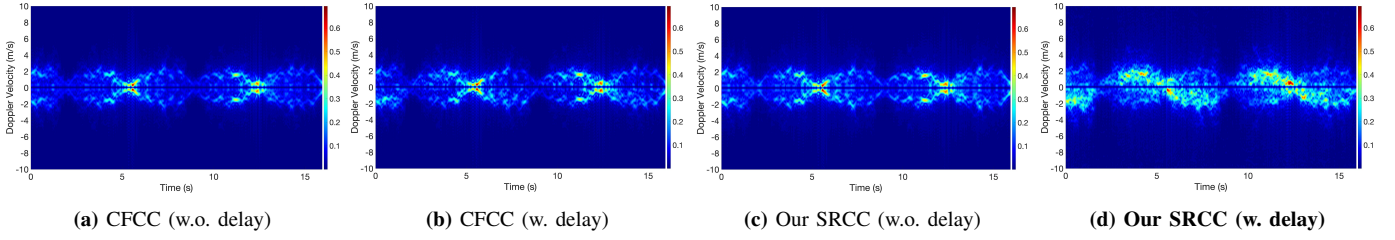
**(a)** CFCC (w.o. delay)  **(b)** CFCC (w. delay)  **(c)** Our SRCC (w.o. delay)  **(d) Our SRCC (w. delay)**

**Fig. 7:** Comparison of 1Tx-1Rx micro-Doppler signatures under 3.1 GHz LTE signals.



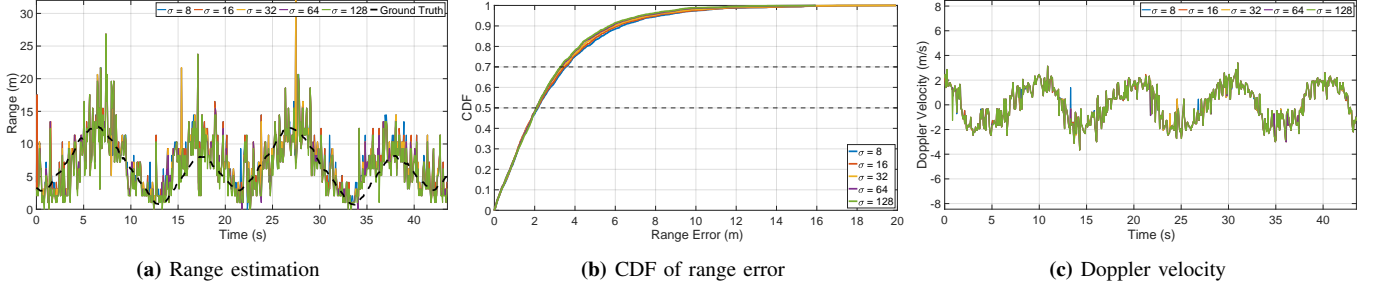**(a)** Range estimation  **(b)** CDF of range error  **(c)** Doppler velocity

**Fig. 8:** Impact of Gaussian window width ($\sigma$)

**TABLE II:** Range estimation errors (m) of our SRCC at the 50% and 70% percentiles. Overall, increasing the Gaussian window width tends to reduce the delay estimation error.

| Metric | $\sigma = 8$ | $\sigma = 16$ | $\sigma = 32$ | $\sigma = 64$ | $\sigma = 128$ |
|---|---|---|---|---|---|
| 50% | 2.16 m | 2.14 m | 2.08 m | **2.08 m** | 2.10 m |
| 70% | 3.55 m | 3.48 m | 3.37 m | 3.33 m | **3.26 m** |

**TABLE III:** Feature generation latency on different platforms without any code-level optimization (milliseconds).

| Platform | Avg. Latency (ms) | Std. Dev. (ms) |
|---|---|---|
| Raspberry Pi 4B (8GB) | 8.5 | 4.3 |
| MacBook Pro 2019 (Intel i7, 2.6 GHz) | 1.2 | 0.46 |

isolation of dominant paths. However, insufficient smoothing introduces more noise in the range estimates. Fig. 8 (a) and (b) shows that increasing $\sigma$ yields a slight improvement in range accuracy. Fig. 8 (c) shows that Doppler estimation remains relatively stable across a wide range of $\sigma$ values. The larger $\sigma$ combines more adjacent delay components, leading to reduced purity of the constructed signal. Despite this, our delay-domain beamforming effectively suppresses mirror artifacts.

*5) Feature Generation Overhead:* Table. III summarizes the average latency on two platforms. On a Raspberry Pi 4B (8GB), the average latency is 8.5 ms with a standard deviation (Std.) of 4.3 ms, without using any optimization such as parallelization or JIT compilation [49]. The MacBook Pro (Intel i7, 2.6 GHz) achieves a much lower latency of 1.2 ms. As discussed in Section V. E, the main computational cost arises from MVDR weight estimation across delay bins. These results illustrate that our method supports real-time execution, even on resource-constrained edge devices.

### B. Single-Target Sensing Performance

We evaluat WiDFS 3.0 under single-target settings, focusing on the impact of network architecture, input features, data augmentation, and generalization, all based on our interpretable and unambiguous micro-Doppler inputs.

**TABLE IV:** Performance comparison of Dataset 1. *Acc.*: Accuracy, *Prec.*: Macro Precision, *Rec.*: Macro Recall, *F1*: Macro F1-score.

| Feature | Model | Acc. | Prec. | Rec. | F1 |
|---|---|---|---|---|---|
| Our SRCC (1Tx–1Rx) | CNN | 0.780 | 0.779 | 0.780 | 0.779 |
| | MLP | 0.796 | 0.794 | 0.794 | 0.792 |
| | SqueezeNet | 0.879 | 0.876 | 0.879 | 0.877 |
| | ShuffleNetV2 | 0.920 | 0.918 | 0.919 | 0.919 |
| | MobileNetV2 | 0.933 | 0.931 | 0.932 | 0.931 |
| | ResNet18 | 0.934 | 0.933 | 0.933 | 0.933 |
| | **MobileViT-XXS** | **0.939** | **0.938** | **0.938** | **0.938** |
| DCACC (1Tx–3Rx) | **MobileViT-XXS** | **0.991** | **0.991** | **0.991** | **0.991** |
| BVP (Multi-Receiver) | MobileViT-XXS | 0.850 | 0.849 | 0.849 | 0.849 |

**TABLE V:** Detailed classification report of each activity class (Dataset 1).

| Class | Precision | Recall | F1-score | Test Samples (#) |
|---|---|---|---|---|
| Clap | 0.938 | 0.946 | 0.942 | 7510 |
| Draw-O | 0.932 | 0.941 | 0.936 | 7965 |
| **Draw-Zigzag** | **0.969** | **0.976** | **0.973** | **10215** |
| Push&Pull | 0.937 | 0.937 | 0.937 | 7740 |
| Slide | 0.922 | 0.898 | 0.910 | 8774 |
| Sweep | 0.931 | 0.932 | 0.931 | 7514 |

*1) Overall Performance across Lightweight Models:* We next present detailed evaluations on two datasets of Widar 3.0, including Dataset 1 for human-computer interaction activities, and Dataset 2 for digit gesture recognition. To ensure a fair comparison, all models are trained and evaluated on each dataset using a 70%-30% train-test split, with uniform sampling across different gestures. Each receiver captures CSI from three Rx antennas, and in each training epoch, we randomly select the micro-Doppler feature from one of the antennas. In addition, we visualize the extracted micro-Doppler features for six gestures, shown in Fig.9. Different hand movements induce distinct Doppler signatures, enabling the model
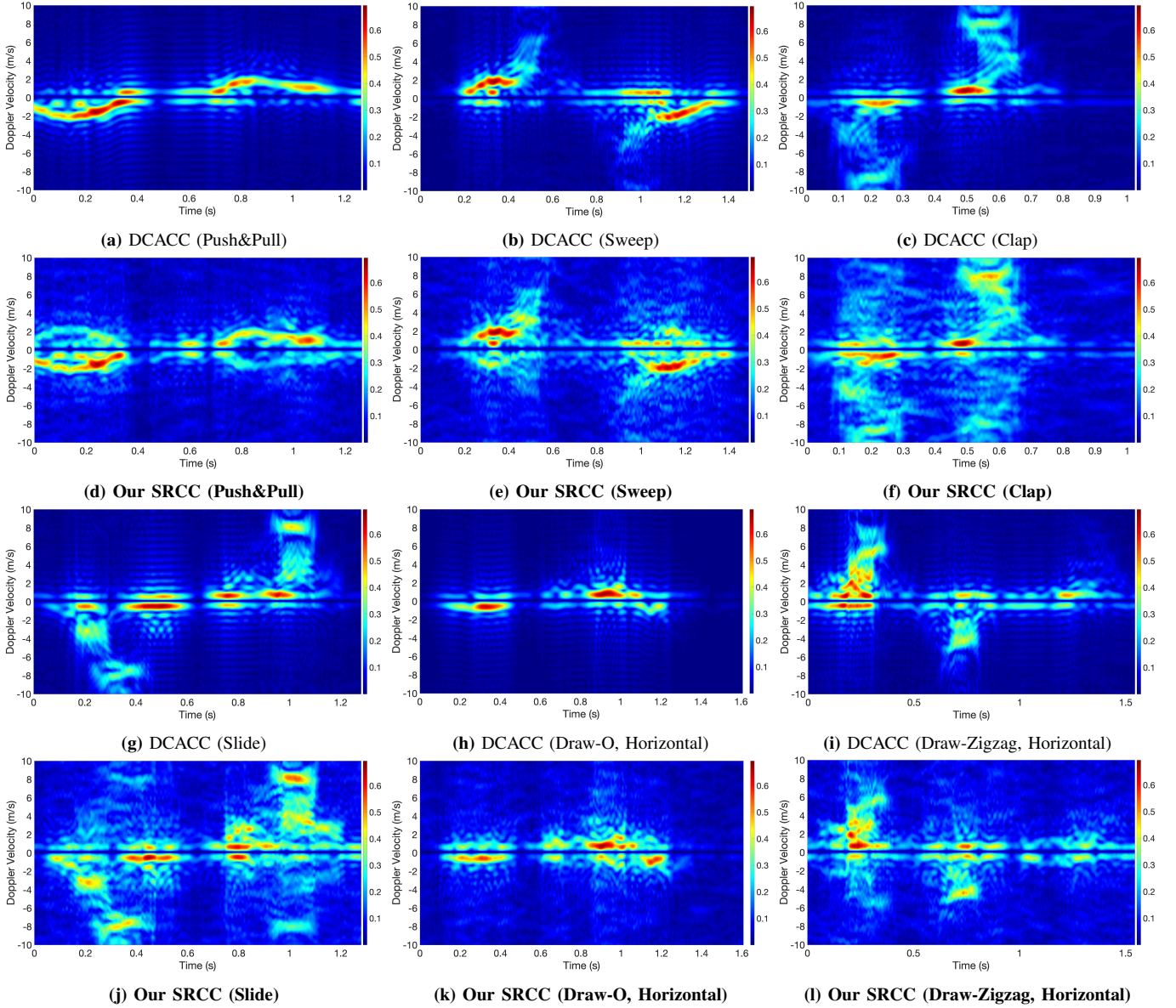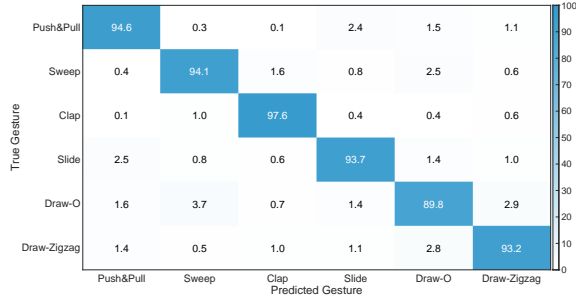
**Fig. 9:** Micro-Doppler signatures for six gestures (Dataset 1) using 1Tx-3Rx DCACC (top) and our 1Tx-1Rx SRCC method (bottom), showing effective suppression of mirror Doppler components in SISO configurations.

to capture different motion dynamics for classification. The 1Tx-3Rx DCACC method, which effectively removes Doppler mirroring artifacts, is used as a reference for comparison. Our 1Tx-1Rx SRCC exhibit significant suppression of mirror Doppler components.
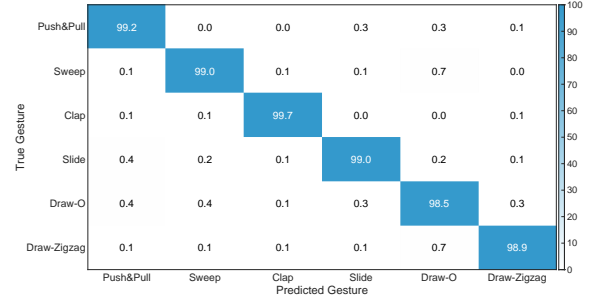
- *Dataset 1*. Table. IV shows that our 1Tx-1Rx SRCC achieves strong and consistent performance across different lightweight backbones, including MobileNetV2, ShuffleNetV2, SqueezeNet, ResNet18, and MobileViT-XXS. Among them, MobileViT-XXS achieves the best macro F1-score (0.938) and accuracy (0.939), slightly outperforming ResNet18, MobileNetV2, and ShuffleNetV2. However, the MLP and CNN models show limited performance. The MLP processes features independently and lacks temporal inductive bias, limiting its ability to model dynamic gestures. The shallow CNN is constrained by local receptive fields, making it

difficult to capture long-range dependencies. In contrast, the attention-based MobileViT model enables global context modeling and spatial preservation, leading to better recognition performance. Table. V further details the per-class classification performance on six human-computer interaction gestures. All classes achieve high precision and recall, with the Draw-Zigzag gesture yielding the highest F1-score of 0.973. This may be partially attributed to its relatively higher number of training samples.

- *Dataset 2*. As summarized in Table. VI and Table VII, our WiDFS 3.0 still maintains high performance, achieving a macro F1-score of 0.928. We can see that digits like Draw-1, Draw-7 and Draw-8 achieve F1-scores above 0.96, while more ambiguous digits like Draw-2 and Draw-6 yield slightly lower scores. Overall, these results validate the robustness of WiDFS 3.0. The interpretable micro-Doppler features extracted from SISO configura-
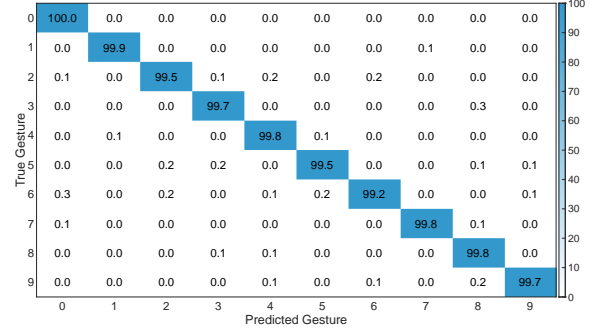
**(a)** Our SRCC (1Tx-1Rx)      **(b)** DCACC (1Tx-3Rx)

Fig. 10: Confusion matrix of human-computer interaction gestures (Dataset 1).



**(a)** Our SRCC (1Tx-1Rx)      **(b)** DCACC (1Tx-3Rx)

Fig. 11: Confusion matrix of numeric drawing gestures (Dataset 2).

TABLE VI: Performance comparison of DataSet 2.

| Input Feature | Model | Accuracy | Precision | Recall | F1-score |
|---|---|---|---|---|---|
| Our SRCC (1Tx–1Rx) | MobileViT-XXS | **0.928** | **0.928** | **0.928** | **0.928** |
| DCACC (1Tx–3Rx) | | **0.997** | **0.997** | **0.997** | **0.997** |
| BVP (Multi-Receiver) | | 0.859 | 0.861 | 0.859 | 0.859 |

TABLE VII: Detailed classification report of each digit gesture class (DataSet 2).

| Class | Precision | Recall | F1-score | Test Samples (#) |
|---|---|---|---|---|
| Draw-0 | 0.932 | 0.944 | 0.938 | 1170 |
| **Draw-1** | **0.996** | **0.985** | **0.991** | **1170** |
| Draw-2 | 0.907 | 0.867 | 0.886 | 1170 |
| Draw-3 | 0.926 | 0.927 | 0.926 | 1171 |
| Draw-4 | 0.920 | 0.938 | 0.929 | 1170 |
| Draw-5 | 0.892 | 0.908 | 0.900 | 1169 |
| Draw-6 | 0.894 | 0.891 | 0.893 | 1170 |
| Draw-7 | 0.967 | 0.957 | 0.962 | 1170 |
| Draw-8 | 0.957 | 0.974 | 0.965 | 1170 |
| Draw-9 | 0.891 | 0.891 | 0.891 | 1170 |

tions maintain strong motion discriminability, enabling lightweight models to achieve high performance .

- *Performance Comparison.* Table. IV and Table. VI report the recognition performance using DCACC-based features, achieving macro F1-scores of 0.991 and 0.997 for Dataset 1 and 2, respectively. The corresponding confusion matrices are shown in Fig. 10 (a) and Fig. 11 (a). The super high recognition accuracy underscore the effectiveness of unambiguous micro-Doppler features in capturing motion dynamics with high discriminability. Our method achieves comparable sensing performance while effectively suppressing Doppler mirroring. However, due to the limitations of single-antenna configurations, the mirror components cannot be entirely eliminated. However, the BVP feature yields lower accuracy on both datasets, primarily due to its reliance on multi-receiver Doppler aggregation without explicit ambiguity suppression, using raw CACC for random phase removal. This makes it more susceptible to mirrored components and inter-receiver calibration errors.

- *Data Augmentation* We use MobileViT-XXS to evaluate the impact of data augmentation due to its highest classification performance described above. Fig. 12 shows the

augmentation strategy yields more than 1% improvement in macro F1-score compared to the case without augmentation. The proposed physics-guided augmentation is beneficial for improving model performance.

*2) Generalization Performance:* We evaluate the generalization of our feature by testing on samples from one target subset while training on all others.

- *Impact of location:* We train the model on data from four target positions and testing it on the remaining one. For example, when Location #1 is used for testing, the samples from the other locations (#2-#8) are used for training, covering various face orientations, receivers, and room. Here, Locations #6, #7, and #8 contain significantly fewer samples (around 3,000 samples only), while Locations #1-#5 have more balanced sample sizes. For fair comparison, we conduct the experiments on Locations
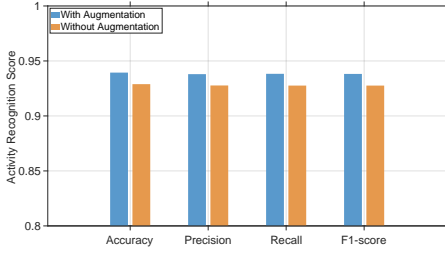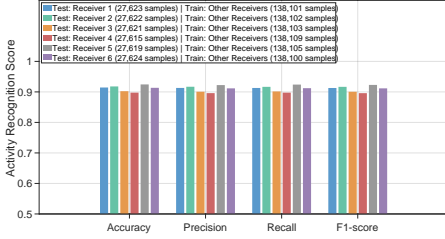
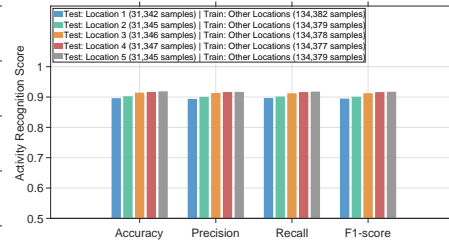Fig. 12: Impact of data augmentation



Fig. 13: Impact of location



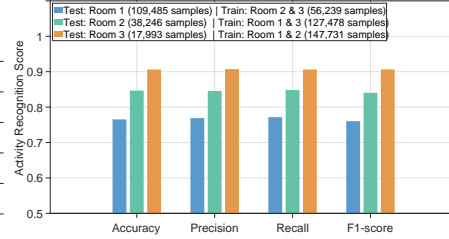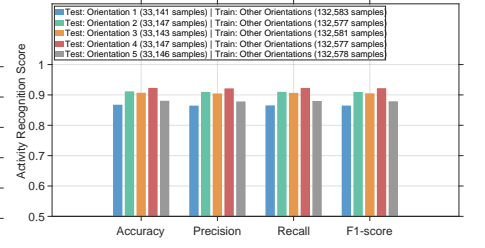Fig. 14: Impact of orientation



Fig. 15: Impact of receiver



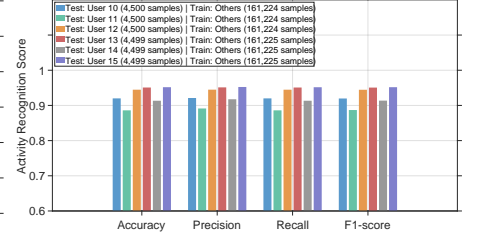Fig. 16: Impact of environment



Fig. 17: Impact of user

**TABLE VIII:** Performance comparison of different input features (Dataset 1) under a strict generalization setting: the model is trained on samples from Location #1 and Receiver #1 (5,224 samples), and tested on samples from all other locations and receivers (160,500 samples).

| Model | Input Feature | Minimum Setup | Accuracy | Macro Precision | Macro Recall | Macro F1-score | Accuracy Gap |
|---|---|---|---|---|---|---|---|
| | CSI Amplitude [50] | 1Tx–1Rx | 0.366 | 0.362 | 0.357 | 0.355 | 0.510 |
| | BVP [17] | Multi-Receiver | 0.422 | 0.440 | 0.417 | 0.418 | 0.454 |
| | CACC (w.o. delay) [10] | 1Tx–2Rx | 0.601 | 0.601 | 0.600 | 0.594 | 0.275 |
| | CACC (w. delay) | 1Tx–2Rx | 0.614 | 0.618 | 0.612 | 0.613 | 0.262 |
| MobileViT-XXS | CASR (w. delay) [9], [12] | 1Tx–2Rx | 0.853 | 0.852 | 0.851 | 0.852 | 0.023 |
| | **DCACC (w. delay) [11], [16]** | **1Tx–3Rx** | **0.876** | **0.875** | **0.873** | **0.873** | **/** |
| | CFCC (w. delay) [25] | 1Tx–1Rx | 0.565 | 0.565 | 0.561 | 0.558 | 0.311 |
| | **Our SRCC (w. delay)** | **1Tx–1Rx** | **0.767** | **0.764** | **0.765** | **0.763** | **0.109** |

#1-#5. As shown in Fig. 13, the sample distribution across the five locations is relatively uniform, achieving macro F1-scores ranging from 0.917 to 0.918. The results show the robustness of our SISO micro-Doppler feature in handling spatial shifts of the target's position.

- *Impact of orientation:* We evaluate the model's generalization under different body orientations. As shown in Fig. 14, the model consistently achieves similar sensing accuracy across all five orientations, with F1-scores above 0.88. The experiment result shows the SISO micro-Doppler feature can also maintain strong orientation invariance, enabling reliable activity recognition even under significant changes in a user's facing direction.

- *Impact of receiver:* This experiment is to assess sensitivity to receiver placement. As shown in Fig. 15, our scheme maintains consistently high performance across all receivers, with F1-scores exceeding 0.91 in each case, which can demonstrate the model's strong generalization to varying transceiver positions.

- *Impact of environment:* As shown in Fig. 16, the model achieves the lowest F1-score of 0.74 when tested on Room 1, while the highest F1-score of 0.93 is observed when testing on Room 3. Here, the number of training samples in Room 1 is approximately half that of the other environments, leading to a data imbalance that contributes to the reduced performance. As the training data increases, the test accuracy tends to improve accordingly. In our feature extraction pipeline, static clutter removal is

employed to reduce the sensitivity of the extracted micro-Doppler features to environmental differences.

- *Impact of User:* To evaluate the generalization ability across different users, we select a subset of six users (User #10 to #15) since each user almost contributes an equal number of training and testing samples, thereby eliminating the impact of data imbalance. As shown in Fig. 17, our feature consistently achieves high F1-scores (0.88–0.95) across all users, indicating strong generalization to user-specific variations.

*3) Comparison of input features:* To examine how different input features affect generalization, we construct a strict train-test split: the model is trained only on samples from Location #1 and Receiver #1 (5,224 samples), and tested on samples from all other locations and receivers (160,500 samples). A summary of these comparisons is provided in Table VIII.

- *Commonly-used Features.* We first compare against three commonly used representations: CSI Amplitude, BVP, and Raw Doppler (directly extracted from CACC without delay filtering). These baselines yield F1-scores of only 0.355, 0.418, and 0.594, respectively. The CSI Amplitude feature is highly sensitive to absolute signal strength, which can vary significantly across hardware setups and environments, resulting in poor generalization. The BVP feature is designed to aggregate Doppler information across multiple receivers to improve motion representation. However, it may be less robust to unseen configurations due to its reliance on consistent receiver geometry.

**TABLE IX:** Performance comparison of different models on activity recognition and people counting.

| Method | Model | Activity Recognition | | | People Counting |
|---|---|---|---|---|---|
| | | Precision | Recall | F1-score | Accuracy |
| Our SRCC (1Tx–1Rx) | MLP | 0.720 | 0.405 | 0.507 | 0.493 |
| | CNN | 0.665 | 0.400 | 0.491 | 0.584 |
| | MobileNetV2 | 0.687 | 0.606 | 0.642 | 0.569 |
| | ShuffleNetV2 | 0.714 | 0.577 | 0.636 | 0.571 |
| | SqueezeNet | 0.700 | 0.562 | 0.622 | 0.551 |
| | ResNet18 | 0.706 | 0.590 | 0.641 | 0.574 |
| | **MobileViT-XXS** | **0.705** | **0.621** | **0.659** | **0.629** |
| DCACC (1Tx–3Rx) | **MobileViT-XXS** | **0.712** | **0.620** | **0.662** | **0.618** |
| CASR (1Tx–2Rx) | MobileViT-XXS | 0.665 | 0.101 | 0.174 | 0.383 |

**TABLE X:** Per-class performance on activity recognition.

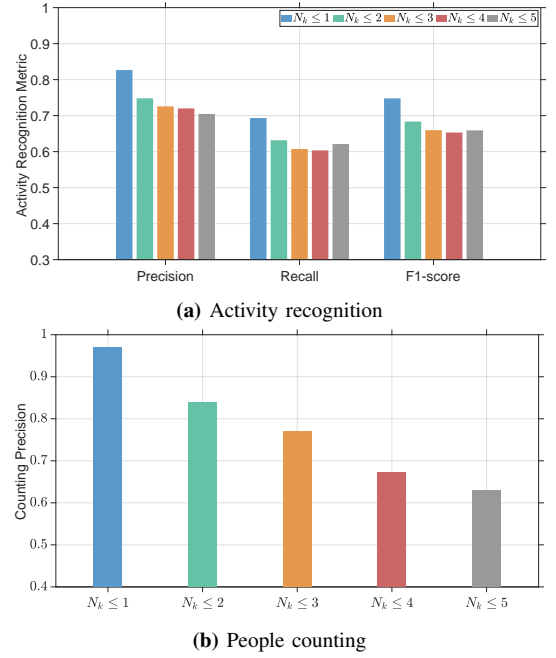| Activity | Precision | Recall | F1-score | Test Samples (#) |
|---|---|---|---|---|
| Jumping | 0.718 | 0.653 | 0.684 | 490 |
| Lying Down | 0.724 | 0.689 | 0.706 | 460 |
| Nothing | 0.648 | 0.489 | 0.558 | 519 |
| Picking Up | 0.701 | 0.554 | 0.619 | 487 |
| Rotation | 0.651 | 0.674 | 0.662 | 481 |
| Sitting Down | 0.655 | 0.564 | 0.607 | 482 |
| Standing Up | 0.632 | 0.579 | 0.604 | 468 |
| **Walking** | **0.927** | **0.809** | **0.864** | **487** |
| Waving | 0.687 | 0.578 | 0.628 | 483 |

Raw Doppler (CACC w.o. delay) suffers from strong mirror-Doppler ambiguity, which undermines its ability to distinguish motion direction.

- *Enhanced Features.* We evaluate Doppler-enhanced features, including multi-antenna based CACC, CASR, and DCACC, as well as single-antenna based CFCC and SRCC, all incorporating delay-domain filtering. Among them, 1Tx-3Rx DCACC achieves the highest performance with a macro F1-score of 0.873, benefiting from effective Doppler mirror suppression. CASR also performs well (F1-score: 0.852) using two Rx antennas. In contrast, CACC and CFCC show moderate gains (F1-scores: 0.613 and 0.558, respectively), showing that delay filtering helps reduce mirror artifacts but is not sufficient. Notably, our 1Tx-1Rx SRCC method achieves a strong F1-score of 0.763, outperforming several multi-antenna baselines. Overall, these results demonstrate that Doppler mirror suppression is key to improving generalization. Our SRCC achieves the optimal trade-off between mirror suppression and hardware simplicity, offering strong generalization with only a single antenna.

### C. Multi-Target Sensing Performance

We evaluate the performance of WiDFS 3.0 in complex multi-target scenarios using the WiMANS dataset. The following experiments use a 3D delay-Doppler-time tensor as input, preserving temporal motion dynamics and spatial separability.

*1) Overall Performance across Lightweight Models:* For our 1Tx-1Rx SRCC method, we randomly select a feature



**(a)** Activity recognition



**(b)** People counting

**Fig. 18:** Impact of the number of people on sensing performance.

from one transmitter-receiver antenna pair in each training epoch. Table XI presents the results for various lightweight models. We can see that MobileViT-XXS achieves the highest F1-score of 0.659 for activity recognition and the best people counting accuracy of 0.629. MobileNetV2 and ShuffleNetV2 also perform competitively, with F1-scores exceeding 0.63. For comparison, the 1Tx-3Rx DCACC configuration further improves recognition performance, benefiting from enhanced Doppler mirror suppression. It achieves an F1-score of 0.662 in activity recognition and a people counting accuracy of 0.618. In contrast, the 1Tx-2Rx CASR yields a significantly lower F1-score of only 0.174 due to poor recall, highlighting its limited suitability in complex multi-target environments. In addition, per-class activity recognition results are detailed in Table X. The *Walking* class, which exhibits the largest body movements, achieves the highest F1-score across all classes. Overall, these results demonstrate the effectiveness and robustness of our 1Tx-1Rx SRCC feature. However, multi-target recognition remains inherently challenging due to the limited bandwidth, which constrains delay resolution, and overlapping Doppler signatures from multiple individuals further reduce per-target separability, making accurate activity and people count estimation more difficult.

*2) Input Feature Comparison:* To assess the effectiveness of different input features, we adopt an corss-environment train-test split: the model is trained on data collected in the *classroom* and *empty room*, and tested on the unseen *meeting room* environment. MobileViT-XXS is used for all comparisons. As shown in Table XI, the extracted delay-Doppler-time feature based on our 1Tx-1Rx SRCC method achieves an F1-score of 0.588 and a people counting accuracy of 0.385, outperforming all 2D baselines such as Amplitude-Time (F1: 0.473), Delay-Time (F1: 0.427), and Doppler-Time (F1: 0.543). These 2D features suffer from dimensionality reduction: Doppler-Time and Delay-Time projections discard

**TABLE XI:** Performance comparison of different features under a cross-environment train-test split.

| Model | Method | Input Feature | Activity Recognition | | | People Counting |
| --- | --- | --- | --- | --- | --- | --- |
| | | | Precision | Recall | F1-score | Accuracy |
| MobileViT-XXS | Our SRCC (1Tx–1Rx) | Amplitude-Time | 0.571 | 0.413 | 0.473 | 0.385 |
| | | Delay-Time | 0.484 | 0.387 | 0.427 | 0.336 |
| | | Doppler-Time | 0.618 | 0.490 | 0.543 | 0.373 |
| | | **Delay-Doppler-Time** | **0.654** | **0.537** | **0.588** | **0.385** |
| | DCACC (1Tx–3Rx) | **Delay-Doppler-Time** | **0.713** | **0.508** | **0.591** | **0.436** |
| | CASR (1Tx–2Rx) | Delay-Doppler-Time | 0.646 | 0.094 | 0.162 | 0.309 |

either spatial or spectral information, while Amplitude-Time lacks physically meaningful motion encoding. In contrast, the 3D feature preserves full spatiotemporal information, leading to improved generalization across environments. We also compare with other multi-antenna methods. The 1Tx-3Rx DCACC approach achieves the highest performance, owing to its enhanced suppression of Doppler mirroring through spatial diversity. In contrast, the 1Tx-2Rx CASR method performs poorly (F1: 0.162), as it is primarily designed for single-target scenarios and struggles to generalize in complex multi-target cases. Overall, these results demonstrate the effectiveness of the proposed SRCC method, which achieves a good trade-off between performance and hardware simplicity. Maintaining both delay and Doppler dimensions is crucial for generalization, especially under multi-target conditions.

*3) Impact of Number of People:* We use the MobileViT-XXS model to evaluate the impact of the number of people. As shown in Fig. 18, both activity recognition and people counting performance tend to degrade as the number of people increases. Specifically, the F1-score for activity recognition drops from 0.748 (for $N_k \leq 1$) to 0.659 (for $N_k \leq 5$), while counting accuracy decreases from 0.969 to 0.629. This degradation is expected due to increased spectral overlap, occlusion, and motion aliasing. With more users, the temporal and spatial separability of delay-Doppler-time patterns diminishes, making it more difficult for the network to distinguish fine-grained motion features.

## IX. CONCLUSION

This work presents WiDFS 3.0, a practical and lightweight bistatic sensing framework for SISO-based ISAC systems. It can operate with a single antenna at both the transmitter and receiver, while remaining extensible to multi-antenna setups. We propose a SRCC technique for efficient CSI random phase removal and introduce a delay-domain beamforming pipeline that produces a robust 3D delay-Doppler-time representation. This feature preserves key motion cues while suppressing interference and ambiguity, and can be effectively utilized by compact neural networks. Extensive experiments show that WiDFS 3.0 consistently outperforms conventional methods and exhibits strong feature generalization in both single- and multi-target scenarios. Overall, it offers a scalable, cost-efficient, and robust sensing solution that brings ISAC closer to practical deployment.

## APPENDIX
### CRLB DERIVATION FOR CSI PHASE ESTIMATION WITH WINDOWING

Assume the reconstructed CSI $\mathcal{CSI}_{i,j}$ is distorted by additive complex Gaussian noise $\mathcal{N}(0, \eta^2)$. We treat the phase $\phi_{i,j} = \angle\mathcal{CSI}_{i,j}$ as the parameter to be analyzed. Under the complex Gaussian noise model, the Fisher Information for $\phi_{i,j}$ is given by:

$$\mathcal{I}(\phi_{i,j}) = \frac{2|\mathcal{CSI}_{i,j}|^2}{\eta^2}, \tag{27}$$

so the CRLB becomes:

$$\mathrm{Var}(\hat{\phi}_{i,j}) \geq \frac{1}{\mathcal{I}(\phi_{i,j})} = \frac{\eta^2}{2\|\mathcal{CSI}_{i,j}\|^2}. \tag{28}$$

From the energy consistancy of the FFT, we approximate $\|\mathcal{CSI}_{i,j}\|^2$ by

$$\|\mathcal{CSI}_{i,j}\|^2 \propto \left\| \mathcal{G}(\tau - \tau_j^{\mathrm{peak}}) \cdot h_j(\tau) \right\|^2. \tag{29}$$

Therefore, the phase estimation variance is bounded by:

$$\mathrm{Var}(\hat{\phi}_{i,j}) \geq \frac{\eta^2}{2 \left\| \mathcal{G}(\tau - \tau_j^{\mathrm{peak}}) \cdot h_j(\tau) \right\|^2}. \tag{30}$$

## REFERENCES

[1] J. A. Zhang, F. Liu, C. Masouros, R. W. Heath, Z. Feng, L. Zheng, and A. Petropulu, "An overview of signal processing techniques for joint communication and radar sensing," *IEEE Journal of Selected Topics in Signal Processing*, vol. 15, no. 6, pp. 1295–1315, 2021.

[2] J. A. Zhang, K. Wu, X. Huang, Y. J. Guo, D. Zhang, and R. W. Heath, "Integration of radar sensing into communications with asynchronous transceivers," *IEEE Communications Magazine*, vol. 60, no. 11, pp. 106–112, 2022.

[3] F. Liu, C. Masouros, and Y. C. Eldar, *Integrated Sensing and Communications.* Springer, 2023.

[4] S. Lu, F. Liu, Y. Li, K. Zhang, H. Huang, J. Zou, X. Li, Y. Dong, F. Dong, J. Zhu *et al.*, "Integrated sensing and communications: Recent advances and ten open challenges," *IEEE Internet of Things Journal*, 2024.

[5] K. Wu, Z. Wang, S.-L. Chen, J. A. Zhang, and Y. J. Guo, "Isac: From human to environmental sensing," *arXiv preprint arXiv:2507.13766*, 2025.

[6] Z. Wang, J. A. Zhang, K. Wu, and Y. J. Guo, "Water level sensing via communication signals in a bi-static system," *arXiv preprint arXiv:2505.19539*, 2025.

[7] Y. He, J. Liu, M. Li, G. Yu, and J. Han, "Forward-compatible integrated sensing and communication for wifi," *IEEE Journal on Selected Areas in Communications*, vol. 42, no. 9, pp. 2440–2456, 2024.

[8] F. Miao, Y. Huang, Z. Lu, T. Ohtsuki, G. Gui, and H. Sari, "Wi-fi sensing techniques for human activity recognition: Brief survey, potential challenges, and research directions," *ACM Computing Surveys*, vol. 57, no. 5, pp. 1–30, 2025.

[9] Y. Feng, Y. Xie, D. Ganesan, and J. Xiong, "Lte-based pervasive sensing across indoor and outdoor," in *Proceedings of the 19th ACM Conference on Embedded Networked Sensor Systems*, 2021, pp. 138–151.

[10] K. Qian, C. Wu, Y. Zhang, G. Zhang, Z. Yang, and Y. Liu, "Widar2. 0: Passive human tracking with a single wi-fi link," in *ACM MobiSys*, 2018, pp. 350–361.

[11] Z. Wang, J. A. Zhang, M. Xu, and J. Guo, "Single-target real-time passive wifi tracking," *IEEE Transactions on Mobile Computing*, vol. 2, no. 6, pp. 3724–3742, 2023.

[12] Z. Ni, J. A. Zhang, K. Wu, and R. P. Liu, "Uplink sensing using csi ratio in perceptive mobile networks," *IEEE Transactions on Signal Processing*, vol. 71, pp. 2699–2712, 2023.

[13] Y. Hu, K. Wu, J. A. Zhang, W. Deng, and Y. J. Guo, "Performance bounds and optimization for csi-ratio-based bi-static doppler sensing in isac systems," *IEEE Transactions on Wireless Communications*, vol. 23, no. 11, pp. 17461–17477, 2024.

[14] K. Chen, J. A. Zhang, Z. Wang, and Y. J. Guo, "Development of an uplink sensing demonstrator for perceptive mobile networks," in *2023 22nd International Symposium on Communications and Information Technologies (ISCIT)*. IEEE, 2023, pp. 191–196.

[15] J. Pegoraro, J. O. Lacruz, T. Azzino, M. Mezzavilla, M. Rossi, J. Widmer, and S. Rangan, "Jump: Joint communication and sensing with unsynchronized transceivers made practical," *IEEE Transactions on Wireless Communications*, vol. 23, no. 8, pp. 9759–9775, 2024.

[16] Z. Wang, J. A. Zhang, H. Zhang, M. Xu, and J. Guo, "Passive human tracking with wifi point clouds," *IEEE Internet of Things Journal*, vol. 12, no. 5, pp. 5528–5543, 2025.

[17] Y. Zhang, Y. Zheng, K. Qian, G. Zhang, Y. Liu, C. Wu, and Z. Yang, "Widar3. 0: Zero-effort cross-domain gesture recognition with wi-fi," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 11, pp. 8671–8688, 2021.

[18] F. Meneghello, D. Garlisi, N. D. Fabbro, I. Tinnirello, and M. Rossi, "Sharp:environment and person independent activity recognition with commodity ieee 802.11 access points," *IEEE Transactions on Mobile Computing*, vol. 22, no. 10, pp. 6160–6175, 2023.

[19] X. Zheng, K. Yang, J. Xiong, L. Liu, and H. Ma, "Pushing the limits of wifi sensing with low transmission rates," *IEEE Transactions on Mobile Computing*, vol. 23, no. 11, pp. 10265–10279, 2024.

[20] Z. Ni, J. A. Zhang, and R. P. Liu, "Deep learning based water level sensing with interference suppression for isac systems," *IEEE Wireless Communications Letters*, no. 99, pp. 1–1, 2025.

[21] F. Luo, S. Khan, B. Jiang, and K. Wu, "Vision transformers for human activity recognition using wifi channel state information," *IEEE Internet of Things Journal*, vol. 11, no. 17, pp. 28111–28122, 2024.

[22] J. Zhang, J. Xue, Y. Li, and S. L. Cotton, "Leveraging online learning for domain-adaptation in wi-fi-based device-free localization," *IEEE Transactions on Mobile Computing*, vol. 24, no. 8, pp. 7773–7787, 2025.

[23] A. Liu, W. Jiang, S. Huang, and Z. Feng, "Multi-modal integrated sensing and communication in internet of things with large language models," *IEEE Internet of Things Magazine*, pp. 1–9, 2025.

[24] R. Du, H. Hua, H. Xie, X. Song, Z. Lyu, M. Hu, Y. Xin, S. McCann, M. Montemurro, T. X. Han *et al.*, "An overview on ieee 802.11 bf: Wlan sensing," *IEEE Communications Surveys & Tutorials*, vol. 27, no. 1, pp. 184–217, 2024.

[25] Y. Hu, K. Wu, J. Andrew Zhang, W. Deng, and Y. Jay Guo, "Cross-frequency sensing in bistatic isac systems," *IEEE Transactions on Wireless Communications*, pp. 1–1, 2025.

[26] T. Yang, P. Zhang, M. Zheng, Y. Shi, L. Jing, J. Huang, and N. Li, "Wirelessgpt: A generative pre-trained multi-task learning framework for wireless communication," *IEEE Network*, pp. 1–1, 2025.

[27] R. J. Mailloux, *Phased array antenna handbook*. Artech house, 2017.

[28] S. Mehta and M. Rastegari, "Mobilevit: Light-weight, general-purpose, and mobile-friendly vision transformer," in *International Conference on Learning Representations*, 2022. [Online]. Available: https://openreview.net/forum?id=vh-0sUt8HlG

[29] S. Huang, K. Li, D. You, Y. Chen, A. Lin, S. Liu, X. Li, and J. A. McCann, "Wimans: A benchmark dataset for wifi-based multi-user activity sensing," in *European Conference on Computer Vision*. Springer, 2024, pp. 72–91.

[30] K. Wu, J. Pegoraro, F. Meneghello, J. A. Zhang, J. O. Lacruz, J. Widmer, F. Restuccia, M. Rossi, X. Huang, D. Zhang, G. Caire, and Y. J. Guo, "Sensing in bistatic isac systems with clock asynchronism: A signal processing perspective," *IEEE Signal Processing Magazine*, vol. 41, no. 5, pp. 31–43, 2024.

[31] S. Dong, J. Zhao, Z. Lu, J. A. Zhang, T. Yang, and J. Deng, "Signal subspace tracking for aoa estimation in isac systems," in *2024 IEEE Wireless Communications and Networking Conference (WCNC)*. IEEE, 2024, pp. 1–6.

[32] K. Wu, J. A. Zhang, X. Huang, and Y. J. Guo, "A low-complexity csi-based wifi sensing scheme for los-dominant scenarios," in *ICC 2023 - IEEE International Conference on Communications*, 2023, pp. 2747–2752.

[33] N. Tadayon, M. T. Rahman, S. Han, S. Valaee, and W. Yu, "Decimeter ranging with channel state information," *IEEE Transactions on Wireless Communications*, vol. 18, no. 7, pp. 3453–3468, 2019.

[34] K. Niu, X. Wang, F. Zhang, R. Zheng, Z. Yao, and D. Zhang, "Rethinking doppler effect for accurate velocity estimation with commodity wifi devices," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 7, pp. 2164–2178, 2022.

[35] X. Zhang, Y. Zhang, G. Liu, and T. Jiang, "Autoloc: Toward ubiquitous aoa-based indoor localization using commodity wifi," *IEEE Transactions on Vehicular Technology*, vol. 72, no. 6, pp. 8049–8060, 2023.

[36] Y. Xie, J. Xiong, M. Li, and K. Jamieson, "md-track: Leveraging multi-dimensionality for passive indoor wi-fi tracking," in *The 25th Annual International Conference on Mobile Computing and Networking*, 2019, pp. 1–16.

[37] Z. Ni, J. A. Zhang, X. Huang, K. Yang, and J. Yuan, "Uplink sensing in perceptive mobile networks with asynchronous transceivers," *IEEE Transactions on Signal Processing*, vol. 69, pp. 1287–1300, 2021.

[38] J. Du, M. Han, Y. Chen, L. Jin, H. Wu, and F. Gao, "Tensor decompositions for integrated sensing and communications," *IEEE Communications Magazine*, vol. 62, no. 9, pp. 128–134, 2024.

[39] Z. Yang, Y. Zhang, K. Qian, and C. Wu, "Slnet: A spectrogram learning neural network for deep wireless sensing," in *20th USENIX Symposium on Networked Systems Design and Implementation (NSDI 23)*, 2023, pp. 1221–1236.

[40] 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," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 1, no. 3, pp. 1–22, 2017.

[41] J. Li, P. Stoica, and Z. Wang, "On robust capon beamforming and diagonal loading," *IEEE transactions on signal processing*, vol. 51, no. 7, pp. 1702–1715, 2003.

[42] C. Chen, G. Zhou, and Y. Lin, "Cross-domain wifi sensing with channel state information: A survey," *ACM Computing Surveys*, vol. 55, no. 11, pp. 1–37, 2023.

[43] V. K. Singh, A. Walecha, A. Gera, R. Jay, A. Bhattacharya, and M. Maity, "Slim-sense: A resource efficient wifi sensing framework towards integrated sensing and communication," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 9, no. 1, pp. 1–33, 2025.

[44] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4510–4520.

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

[46] N. Ma, X. Zhang, H.-T. Zheng, and J. Sun, "Shufflenet v2: Practical guidelines for efficient cnn architecture design," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 116–131.

[47] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[48] D. Halperin, W. Hu, A. Sheth, and D. Wetherall, "Tool release: Gathering 802.11 n traces with channel state information," *ACM SIGCOMM computer communication review*, vol. 41, no. 1, pp. 53–53, 2011.

[49] S. K. Lam, A. Pitrou, and S. Seibert, "Numba: A llvm-based python jit compiler," in *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC*, 2015, pp. 1–6.

[50] J. Yang, X. Chen, H. Zou, C. X. Lu, D. Wang, S. Sun, and L. Xie, "Sensefi: A library and benchmark on deep-learning-empowered wifi human sensing," *Patterns*, vol. 4, no. 3, 2023.
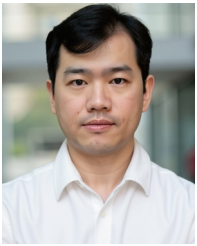
**Zhongqin Wang** is presently working as a Postdoctoral Research Fellow in the School of Electrical and Data Engineering at the University of Technology Sydney. He worked as a Lecturer at the School of Information Engineering, Capital Normal University, Beijing China, from 2022 to 2023. He attained his Ph.D. degree from the University of Technology Sydney, Australia, in 2021, and a M.S. degree from Nanjing University of Posts and Telecommunications, Nanjing, China, in 2014. His research interests include Radio Sensing and ISAC.

**J. Andrew Zhang** (M'04-SM'11) received the B.Sc. degree from Xi'an JiaoTong University, China, in 1996, the M.Sc. degree from Nanjing University of Posts and Telecommunications, China, in 1999, and the Ph.D. degree from the Australian National University, Australia, in 2004.

Currently, Dr. Zhang is a Professor in the School of Electrical and Data Engineering, University of Technology Sydney, Australia. He was a researcher with Data61, CSIRO, Australia from 2010 to 2016, the Networked Systems, NICTA, Australia from 2004 to 2010, and ZTE Corp., Nanjing, China from 1999 to 2001. Dr. Zhang's research interests are in the area of signal processing for wireless communications and sensing, with a focus on integrated sensing and communications. He has published more than 300 papers in leading journals and conference proceedings, and has won 6 best paper awards for his work, including in IEEE ICC2013. He is a recipient of CSIRO Chair's Medal and the Australian Engineering Innovation Award in 2012 for exceptional research achievements in multi-gigabit wireless communications.

**Kai Wu** (Member, IEEE) received the B.E. degree from Xidian University, Xi'an, China, in 2012, and the Ph.D. degree from Xidian University in 2019 and from the University of Technology Sydney (UTS), Sydney, Australia, in 2020. From Nov 2017 to Nov 2018, he was a visiting scholar at DATA61, Commonwealth Scientific and Industrial Research Organisation (CSIRO).

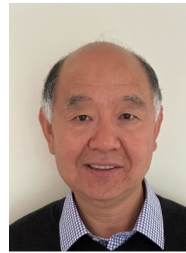He is currently a Lecturer with the School of Electrical and Data Engineering (SEDE) and the Global Big Data Technologies Centre (GBDTC) at UTS. He is the system architect of the TPG-UTS Networking Sensing Lab. His research interests include space/time/frequency signal processing and its applications in radar and communications and their joint designs. He published an authored book on joint communications and sensing (JCAS), aka integrated sensing and communications (ISAC), in December 2022.

His UTS Ph.D. degree was awarded "Chancellor's List 2020." His Xidian PhD thesis was awarded the "Best Ph.D. Thesis Award 2019" by the Chinese Institute of Electronics. He was awarded the Exemplary Reviewer for IEEE TRANSACTIONS ON COMMUNICATIONS, 2021. He is a Tutorial Speaker of WCNC'20, ICC'20, ISCIT'23, and RadarConf'23, presenting JCAS fundamentals and advancement. He was the TPC and special session (Co-)Chair/Member of numerous international conferences, e.g., ICC'20-23 and ISCIT'23. He is serving as the EiC Assistant for the IEEE ISAC-ETI Newsletter. He is an Associate Editor for IEEE Trans. on Mobile Computing, and has been a Guest Editor for the Special Issues in IEEE Journals.

**Min Xu** (M'10) is currently a Professor at University of Technology Sydney. She received the B.E. degree from the University of Science and Technology of China, Hefei, China, in 2000, the M.S. degree from National University of Singapore, Singapore, in 2004, and the Ph.D. degree from University of Newcastle, Callaghan NSW, Australia, in 2010. Her research interests include multimedia data analytics, computer vision and machine learning. She has published over 100 research papers in high quality international journals and conferences. She has been invited to be a member of the program committee for many international top conferences, including ACM Multimedia Conference and reviewers for various highly-rated international journals, such as IEEE Transactions on Multimedia, IEEE Transactions on Circuits and Systems for Video Technology and much more. She is an Associate Editor of Journal of Neurocomputing.

**Y . Jay Guo** (Fellow' 2014) received a Bachelor's Degree and a Master's Degree from Xidian University in 1982 and 1984, respectively, and a Ph.D Degree from Xian Jiaotong University in 1987, all in China. His current research interests include 6G antennas, mm-wave and THz communications and sensing systems as well as big data technologies such as deep machine learning and digital twin. He has published six books and over 800 research papers, and he holds 27 international patents.

Jay is a Fellow of the Australian Academy of Engineering and Technology, Royal Society of New South Wales and IEEE. He has won a number of the most prestigious Australian national awards. Together with his students and postdocs, he has won numerous best paper awards at international conferences such as IEEE AP-S, EuCAP and ISAP. He was a recipient of the prestigious 2023 IEEE APS Sergei A. Schelkunoff Transactions Paper Prize Award.

Jay is a Distinguished Professor and the founding Director of Global Big Data Technologies Centre (GBDTC) at the University of Technology Sydney (UTS), Australia. He is the founding Technical Director of the New South Wales (NSW) Connectivity Innovation Network (CIN). He is also the Founding Director of the TPG-UTS Network Sensing Lab. Before joining UTS in 2014, Prof Guo served as a Research Director in CSIRO for over nine years. Prior to CSIRO, he held various senior technology leadership positions in Fujitsu, Siemens and NEC in the U.K.