

RF-Prox: Radio-based Proximity Estimation of Non-directly Connected Devices

Journal:	<i>IEEE Internet of Things Journal</i>
Manuscript ID	Draft
Manuscript Type:	Special Issue on Integrated Sensing and Communications (ISAC) for 6G IoE
Date Submitted by the Author:	n/a
Complete List of Authors:	Gao, Yuchong Chi, Guoxuan; Tsinghua University, School of Software Yang, Zheng ; ZHENG YANG Cheng, Shijie; Tsinghua University, School of Software Wei, Zhiqing; Key Laboratory of Universal Wireless Communications, Ministry of Education, Beijing University of Posts and Telecommunications
Keywords:	Wireless Localization, Wireless Sensing, Domain Adaptation, Proximity Estimation, Transfer Learning, Spatio-temporal Encoder

RF-Prox: Radio-based Proximity Estimation of Non-directly Connected Devices

Yuchong Gao, *Student Member, IEEE*, Guoxuan Chi*, *Member, IEEE*,
Zheng Yang, *Fellow, IEEE*, Shijie Cheng, Zhiqing Wei, *Member, IEEE*,

Abstract—Recent years have witnessed an increasing number of mobile devices, posing a more diversified demand for device localization solutions. While existing wireless localization solutions can obtain the relative locations of connected devices, they fall short in estimating the spatial relationships between devices that are not directly connected. In response to this technical challenge, we introduce *RF-Prox*, the pioneering system designed for the proximity estimation of non-directly connected devices. *RF-Prox* leverages the analysis of received wireless signals to evaluate the spatial proximity between two devices. It incorporates an innovative multi-resolution spatio-temporal encoder that extracts multi-scale spatio-temporal features from complex-valued wireless signals. These features are subsequently analyzed and converted into a domain-adaptive proximity metric. To enhance the generalizability of *RF-Prox*, we adopt a transfer learning framework. Initially, *RF-Prox* is pre-trained using a substantial dataset from a source domain and subsequently fine-tuned with data from the target deployment domain, markedly diminishing the necessity for extensive data collection within the target domain. *RF-Prox* has been rigorously implemented and its efficacy evaluated in environments based on both Wi-Fi and cellular networks, encompassing indoor and outdoor scenarios. Our evaluation results indicate that, upon fine-tuning, *RF-Prox* exhibits an exemplary accuracy rate of 98.6% and 91.3% in identifying the most proximate device in indoor and outdoor settings, respectively. Remarkably, even without fine-tuning, the pre-trained model demonstrates an impressive zero-shot accuracy of 92.2% and 88.9%, respectively, showcasing its exceptional performance in both proximity estimation accuracy and domain generalizability.

Index Terms—Domain adaptation, proximity estimation, transfer learning, spatio-temporal encoder.

I. INTRODUCTION

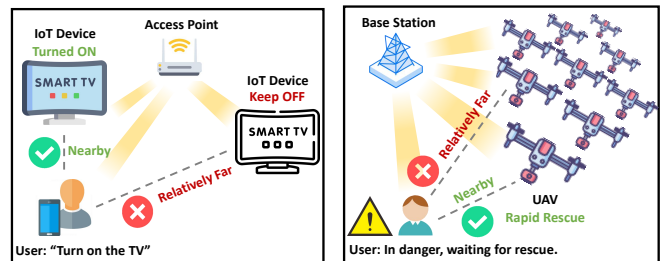
LOCATION awareness is a key enabler for a wide range of applications such as smart homes, augmented reality, and security monitoring [1]. With the increasing number of mobile devices, extensive research efforts have been devoted to wireless-based localization, which infers the devices' relative locations from ubiquitous radio signals.

Despite the advances in positioning technologies such as the Global Positioning System (GPS), which offers satisfactory outdoor positioning accuracy, its performance is significantly impaired under extreme weather conditions due to interference or obstruction of satellite signals with ground devices [2]. Moreover, rapid urbanization has led to an increased indoor time for individuals, escalating the demand for indoor

Yuchong Gao, Guoxuan Chi, Zheng Yang and Shijie Cheng are with the School of Software and BNRist, Tsinghua University, Beijing 100084, China (e-mail: gaoyc01@gmail.com, chiguoxuan@gmail.com, hmilyyz@gmail.com, chengsj23@mails.tsinghua.edu.cn).

Zhiqing Wei is with Key Laboratory of Universal Wireless Communications, Ministry of Education, Beijing University of Posts and Telecommunications (BUPT), Beijing 100876, China (e-mail: weizhiqing@bupt.edu.cn).

*Guoxuan Chi is the corresponding author.



(a) Implicit control of IoT device (b) Proximity-based UAVs scheduling
Fig. 1. Illustration of two application scenarios of *RF-Prox*.

positioning services across commercial, governmental, and communication sectors. This surge in demand places higher requirements on the coverage of positioning services [3], highlighting the limitations of GPS in indoor environments and emphasizing the need for integrated solutions suitable for both outdoor and indoor wireless sensing in scenarios where GPS is unavailable.

Current wireless localization methods are primarily designed for devices with direct communication links, such as wireless access points (AP) and user equipment (UE). However, these methods fall short in determining spatial relationships between non-directly connected devices (e.g., UE and IoT devices), a capability that is pivotal for a range of emerging applications including implicit control of IoT devices and proximity-based unmanned aerial vehicle (UAVs) scheduling, as illustrated in Fig. 1.

One straightforward approach involves estimating the location of each device independently, then deducing their relative proximity from these estimates. However, geometric-based approaches that rely on channel parameters like angle-of-arrival (AoA) [4], [5], time-of-flight (ToF) [6], [7], and their fusion [8], [9], are prone to significant errors in non-line-of-sight (NLoS) conditions. On the other hand, the fingerprint-based localization technique, while effective, requires extensive labeled data collection and faces major generalization challenges across different domains [10], [11].

Unlike traditional device localization methods, our approach is inspired by the principle of "estimating by comparing", based on the observation that wireless devices nearby share similar signal propagation characteristics. By analyzing and comparing the spatio-temporal features encoded in the signals received from two devices, we can estimate their proximity. Nonetheless, actualizing this concept into a functional system presents formidable challenges. Firstly, the precise extraction of spatio-temporal features is complex, as conventional geometric parameters like Angle of Arrival (AoA) and Time of Flight (ToF) are plagued by significant inaccuracies in

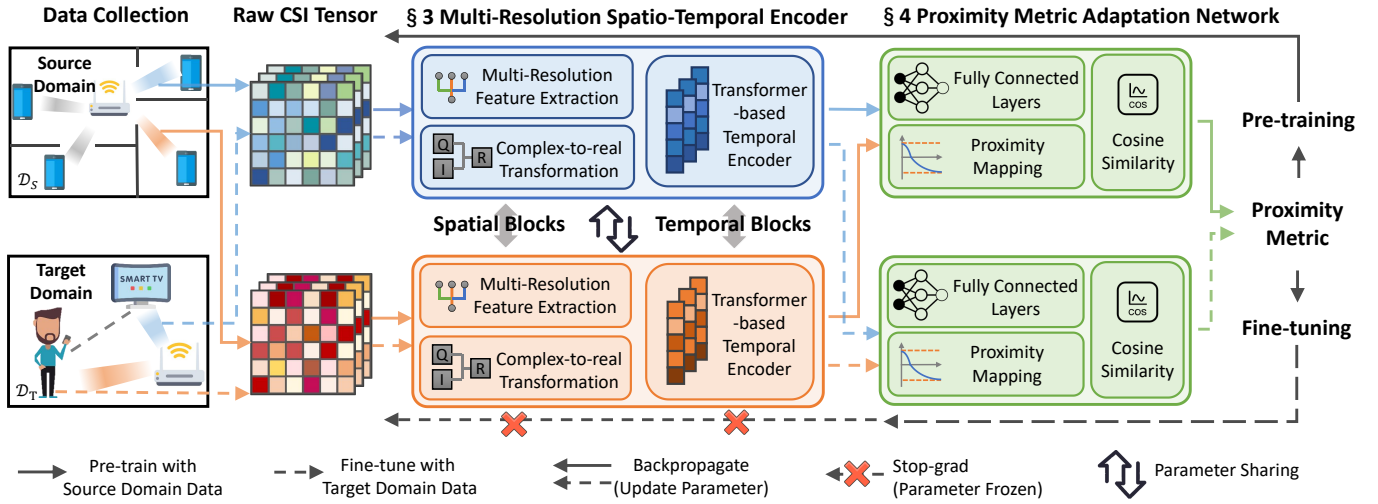


Fig. 2. An overview of the *RF-Prox*, where solid and dashed lines represent data collection from source domain and target domain, respectively, with blue and orange used to distinguish devices.

Non-Line-of-Sight (NLoS) conditions [9]. Secondly, devising a domain-adaptive proximity metric is essential, considering the variability in signal propagation due to differing scenario configurations and device placements, which directly affects proximity assessments.

To address these challenges, we introduce *RF-Prox*, the first proximity estimation system for wireless devices that are not directly connected, adaptable to a wide range of radio frequency signals. For accurate spatio-temporal feature extraction, we employ a data-driven strategy and create a complex-valued neural network module called Multi-Resolution Spatio-Temporal Encoder (MRSTE). This encoder excels at deriving multi-scale latent representations from the wireless signal and fusing them into a comprehensive feature vector that captures the spatio-temporal characteristics of the wireless channel. To establish a domain-adaptive proximity metric, we formulate a Proximity Metric Adaptation Network (PMAN), which compares the spatio-temporal features of two wireless channels to assess device proximity, incorporating domain adaptation techniques. In *RF-Prox*, we utilize a transfer learning approach [12], enabling model pre-training on source domain data followed by minimal fine-tuning with target domain data, significantly diminishing the need for extensive target domain data collection while maintaining broad generalization capabilities.

RF-Prox's efficacy is validated through extensive evaluations across over 9,000 domains, encompassing Wi-Fi-based indoor environments and cellular-based outdoor scenarios, with more than 1,000,000 data samples gathered. The evaluation results highlight that a fine-tuned *RF-Prox* achieves remarkable accuracy of 98.6% and 91.3% in identifying the most proximate device in indoor and outdoor scenarios, respectively. Impressively, even without fine-tuning, the pre-trained model demonstrates substantial zero-shot accuracy, reaching 92.2% and 88.9%, underscoring its exceptional performance in proximity estimation accuracy and domain adaptability. Additionally, *RF-Prox*'s superior spatial-awareness capability is evidenced through its application of a sorting metric, the Normalized Discounted Cumulative Gain (NDCG), in both scenarios.

We summarize our contributions as follows.

- We propose *RF-Prox*, the first proximity estimation system for non-directly connected wireless devices. *RF-Prox* shows the domain-adaptive capability and can be easily deployed in any target domain environment, making a promising step towards integrated sensing and communication.
- Our proposed Multi-Resolution Spatio-Temporal Encoder is a pioneering attempt at applying complex-valued neural networks to wireless sensing. The multi-resolution design and transformer-based temporal processing model have unique advantages and can also be integrated into other types of wireless sensing applications.
- The transfer learning mechanism adopted by our system has been proven effective, providing a new approach to enhance the generalizability of data-driven wireless systems.
- We implement and evaluate *RF-Prox* on both Wi-Fi-based indoor environments and cellular-based outdoor scenarios, which showcases the practicality and effectiveness of deploying *RF-Prox* in target domain scenarios. We made our code, data and pre-trained models publicly available¹ to facilitate the research community.

Compared to our prior conference version [13], we have made significant enhancements include expanding the system to support full-scene proximity estimation across different wireless signal types, thereby improving universality, scaling, and facilitating Integrated Sensing and Communication (ISAC) in the 6G era. In our module design, we have seamlessly integrated transformer-based temporal analysis with CNN-based spatial feature extraction, offering improved adaptability to device mobility and bolstered support for analyzing temporal data in practical applications. Additionally, we've replaced the previous ordered MLP-based feature mapping module with an unordered cosine similarity module, and extend the transfer learning of simulation-to-reality to a more pervasive transfer task under source and target domains. To validate these enhancements, we have comprehensively revised all experiments, incorporated cellular-based outdoor scenarios, redesigned and retrained our models, ensuring they are fine-tuned for the new

¹Our project is available [here](#).

case. We've also conducted a full evaluation and optimization of these additions. To foster transparency and collaboration, we have made all related code publicly available.

The rest of this paper is organized as follows. We begin with an overview of *RF-Prox* in Section II, followed by the detailed design of the MRSTE in Section III and PMAN in Section IV. Our implementation and evaluation of *RF-Prox* are shown in Section V, followed by the related work in Section VI, and the conclusion in Section VII.

II. SYSTEM OVERVIEW

RF-Prox is a device proximity estimation system based on the wireless signal, which integrates two pivotal components: the *Multi-Resolution Spatio-Temporal Encoder (MRSTE)* and the *Proximity Metric Adaptation Network (PMAN)*. As depicted in Fig. 2, the workflow of *RF-Prox* initiates by capturing the Channel State Information (CSI) from wireless links between two distinct mobile devices. This CSI data is then processed through the MRSTE module to extract relevant features. The MRSTE module employs a complex-valued neural network, incorporating residual convolution blocks to derive multi-resolution latent representations from the CSI's real and imaginary components. Following this, a complex-to-real transformation is performed, converting these complex representations into real-valued spatial features for subsequent analysis. The MRSTE concludes with a transformer-based temporal module, which compresses and extracts latent temporal information. After processing through MRSTE, the domain-agnostic spatio-temporal features are amalgamated and fed into the PMAN module's fully connected layers for a comprehensive analysis and comparison. The proximity metric between two devices is then determined using cosine similarity, after transformation by a elaborately designed proximity mapping function.

RF-Prox adopts a transfer learning framework, initiating with a model pre-trained in a source domain, which can then be directly applied and fine-tuned within a target domain as necessary. During the pre-training phase, a substantial dataset from the source domain \mathcal{D}_S is utilized to enhance the model's ability to generalize and extract domain-independent features. After pre-training, the model can be fine-tuned with only a minial dataset from the target domain \mathcal{D}_T . By leveraging the transfer learning mechanism, *RF-Prox* can be efficiently adapted to new environments, facilitating its practical deployment in varied scenarios.

III. MULTI-RESOLUTION SPATIO-TEMPORAL ENCODER

In this section, we introduce the *Multi-Resolution Spatio-Temporal Encoder (MRSTE)*, designed to extract the domain-independent spatio-temporal information embedded in the CSI. As depicted in Figure 3, MRSTE takes the complex-valued CSI tensor as input and transforms it into multi-resolution latent spaces via paralleled residual convolution blocks. Latent representations with different resolutions are then fused by channel concatenation. After passing through a fully connected layer, the fused complex-valued representation is converted to a real-valued spatial feature. The converted spatial features of a time-series are then sent to the Transformer block for temporal feature extraction, which could be further used for robust device proximity estimation.

Compared with geometric-based algorithms [5], [8], the MRSTE adopts a data-driven approach, analyzing signal statistical information in high-dimensional space. This approach significantly enhances the efficacy of the system in diverse settings, including both indoor and outdoor environments, particularly in scenarios afflicted by Non-Line-of-Sight (NLoS) conditions.

A. CSI Preliminary

Considering the phenomenon of multipath propagation, the wireless channel can be modeled as a function of frequency f and time t , expressed as

$$H(f, t) = \sum_{l=1}^L \alpha_l(t, f) e^{-j2\pi f \tau_l(t)}, \quad (1)$$

where L denotes the number of multipath components, $\alpha_l(t, f)$ embodies the complex attenuation factor, and $\tau_l(t)$ the propagation delay corresponding to the l -th path, respectively. CSI represents a discretized sampling of the channel response [14], with frequency domain samples positioned on specific OFDM subcarriers, time domain samples corresponding to each received packet, and spatial domain samples for each radio chain (i.e., Tx-Rx pair), rendering CSI a complex-valued tensor $\mathbf{H} \in \mathbb{C}^{T \times S \times A}$, with T , S , and A indicating the number of time samples, subcarriers, and radio chains, respectively.

B. Complex-valued Network for CSI Processing

Previous studies have often utilized processed CSI data, such as the short-time Fourier transform and ToF-AoA spectrogram [15], [16], as inputs to classification network models for learning, or have divided the CSI into its real and imaginary components for independent processing within deep neural networks [17]. Conversely, our approach leverages the raw, unmodified CSI to mine richer spatio-temporal information. Therefore, we embrace the concept of the complex-valued neural network, incorporating novel elements such as complex-valued linear and convolutional layers into the MRSTE.

To start with, a linear transformation for a CSI matrix $\mathbf{H} = \mathbf{H}_r + j\mathbf{H}_i$ with complex-valued weight $\mathbf{W} = \mathbf{W}_r + j\mathbf{W}_i$ can be decomposed into several real-valued transformations:

$$\text{Linear}(\mathbf{H}; \mathbf{W}) = \begin{bmatrix} \Re(\mathbf{WH}) \\ \Im(\mathbf{WH}) \end{bmatrix} = \begin{bmatrix} \mathbf{W}_r & -\mathbf{W}_i \\ \mathbf{W}_r & \mathbf{W}_i \end{bmatrix} \begin{bmatrix} \mathbf{H}_r \\ \mathbf{H}_i \end{bmatrix}. \quad (2)$$

Similarly, given a complex kernel $\mathbf{C} = \mathbf{C}_r + j\mathbf{C}_i$, the convolution operation $\mathbf{C} * \mathbf{H}$ on the complex domain can also be equivalently written into the following form:

$$\text{Conv}(\mathbf{H}; \mathbf{C}) = \begin{bmatrix} \Re(\mathbf{C} * \mathbf{H}) \\ \Im(\mathbf{C} * \mathbf{H}) \end{bmatrix} = \begin{bmatrix} \mathbf{C}_r & -\mathbf{C}_i \\ \mathbf{C}_r & \mathbf{C}_i \end{bmatrix} * \begin{bmatrix} \mathbf{H}_r \\ \mathbf{H}_i \end{bmatrix}. \quad (3)$$

Research has affirmed [18] the feasibility of implementing dropout, batch normalization, and activation mechanisms directly within the complex domain by independently manipulating the real and imaginary components of the input. This approach ensures that each complex module within the MRSTE is a synthesis of operations conducted in the real domain, thereby preserving the differentiability across the entirety of the MRSTE module.

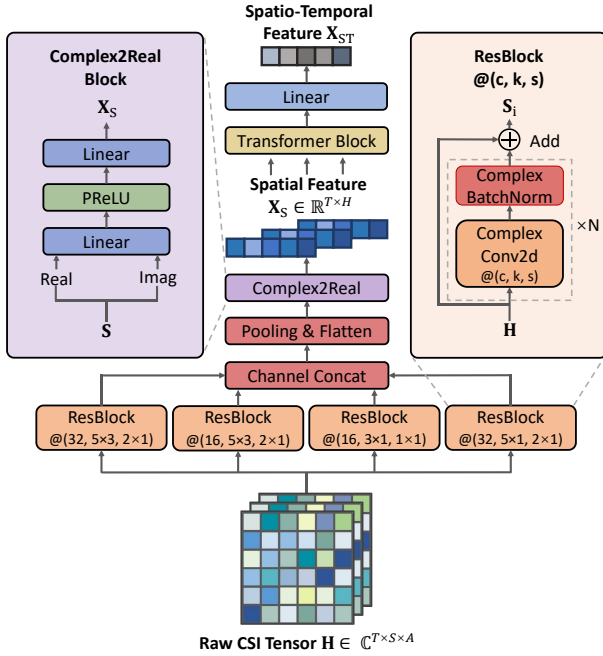


Fig. 3. Illustration of Multi-Resolution Spatio-Temporal Encoder.

C. Multi-Resolution Feature Extraction

The core design of MRSTE is the principle of fusing CSI features across multiple resolutions, a strategy that has demonstrated its efficacy within the domain of computer vision [19]. The underlying logic for this methodology is that variations in antenna spacing, when measuring the Angle-of-Arrival (AoA) from CSI, can introduce a balance between resolution and the operational range [20].

Fig. 3 showcases the MRSTE's structure, which includes four residual convolution blocks, each characterized by unique output channels, kernel dimensions, and stride lengths, thereby forging four concurrent processing pathways. Let us represent a residual block as $\text{ResBlock}(\cdot)$ and define \mathbf{C}_i as the parameter set corresponding to the i -th residual block. For an input CSI tensor \mathbf{H} , the output feature from the i -th block can be articulated as

$$\mathbf{S}_i = \text{ResBlock}(\mathbf{H}; \mathbf{C}_i), \quad i = 0, 1, 2, 3, \quad (4)$$

where the residual block is basically a convolution with shortcut connection [21], which makes the model easier to train by solving the gradient disappearance problem during the training for better expressive ability.

$$\text{ResBlock}(\mathbf{H}; \mathbf{C}_i) = \text{BatchNorm}(\text{Conv}(\mathbf{H}; \mathbf{C}_i)) + \mathbf{H}. \quad (5)$$

Features extracted from parallel residual blocks are then concatenated along the channel dimension and fuse into a latent representation $\mathbf{S} = \text{Concat}(\mathbf{S}_0, \mathbf{S}_1, \mathbf{S}_2, \mathbf{S}_3)$. The concatenated \mathbf{S} contains multi-level features of the CSI input, which greatly improves the receptive field of MRSTE and thus enhances the generalization performance of *RF-Prox*.

D. Complex-to-Real Transformation

After processing CSI with paralleled residual blocks, multi-level features can be extracted. In order to transform the complex-valued latent representation \mathbf{S} to the real-valued spatial feature $\mathbf{X}_S \in \mathbb{R}^{T \times H}$, where H is hidden dimension, we

design a complex-to-real transformation module C2R, which applies two linear operations on the real and imaginary part:

$$\begin{aligned} \mathbf{X}_S &= \text{C2R}(\mathbf{S}; \mathbf{W}^R, \mathbf{W}^I) \\ &= \text{PReLU}(\text{Linear}(\Re(\mathbf{S}), \mathbf{W}^R) + \text{Linear}(\Im(\mathbf{S}), \mathbf{W}^I)), \end{aligned} \quad (6)$$

where \mathbf{W}^R and \mathbf{W}^I are the real-valued linear weights. Note that for better expressive ability of the model, the PReLU activation function is leveraged to add non-linear factors to the features:

$$\mathbf{Y} = \text{PReLU}(\mathbf{X}) = \begin{cases} \mathbf{X}, & \text{if } X \geq 0 \\ \theta \mathbf{X}, & \text{if } X < 0, \end{cases} \quad (7)$$

where θ is a learnable parameter.

E. Transformer-based Temporal Encoder

Through the complex-to-real transformation module C2R, a time-series of real-valued spatial features are extracted. In order to extract the temporal features carried by the device as it moves, we use the transformer encoder module with the attention mechanism [22]. We first embed the series of \mathbf{X}_S into a high-dimensional representation $\mathbf{X}_E = \text{FC}(\mathbf{X}_S)$ with fully connected layers $\text{FC}(\cdot)$. To introduce temporal information, we use absolute positional encoding for \mathbf{X}_E to get embedded positional information \mathbf{P} , with each element as follows:

$$\begin{aligned} \mathbf{P}(\text{pos}, 2i) &= \sin\left(\text{pos}/10000^{2i/d}\right), \\ \mathbf{P}(\text{pos}, 2i+1) &= \cos\left(\text{pos}/10000^{2i/d}\right), \end{aligned} \quad (8)$$

where pos represents the sequence element's ordinal position, and i refers to the dimension index within the embedding space. Then, we formulate the transformer's encoded input as $\mathbf{X}_{\text{PE}} = \mathbf{X}_E + \mathbf{P}$.

In order to better learn the relationships between the elements inside the sequence, we use the attention mechanism to linearly transform the input \mathbf{X}_{PE} into queries, keys and values:

$$\begin{aligned} \mathbf{Q} &= \mathbf{X}_{\text{PE}} \mathbf{W}_Q, \\ \mathbf{K} &= \mathbf{X}_{\text{PE}} \mathbf{W}_K, \\ \mathbf{V} &= \mathbf{X}_{\text{PE}} \mathbf{W}_V, \end{aligned} \quad (9)$$

where \mathbf{W}_Q , \mathbf{W}_K and \mathbf{W}_V are the weight matrices used for linear transformation. Attention features \mathbf{X}_A are obtained as

$$\begin{aligned} \mathbf{X}_A &= \text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \\ &= \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right) \mathbf{V}, \end{aligned} \quad (10)$$

where d_k is the dimension of keys.

After residual connection and layer normalization, the spatio-temporal features \mathbf{X}_{ST} in the CSI can be finally extracted as

$$\mathbf{X}_{\text{ST}} = \text{LayerNorm}(\mathbf{X}_A + \mathbf{X}_{\text{PE}}), \quad (11)$$

which encapsulates the spatio-temporal relationship between the device and the access point. Thus, the spatio-temporal features corresponding to two different terminal devices can be further utilized to determine their proximity relationship, which will be detailed in Section IV.

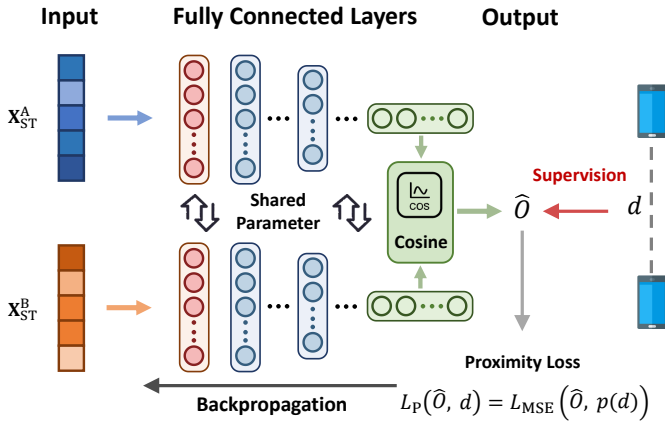


Fig. 4. An overview of the Proximity Metric Adaptation Network.

IV. PROXIMITY METRIC ADAPTATION NETWORK

A. Metric Network Design

In this section, we introduce the *Proximity Metric Adaptation Network (PMAN)*. As illustrated in Fig. 4, PMAN transforms the spatio-temporal features extracted from two wireless devices \mathbf{X}_{ST} to their proximity metric with domain adaptation capability. PMAN is implemented based on the fully connected layers, cosine similarity and an elaborately designed proximity loss function $L_P(\cdot)$. When being deployed in a new environment with different landscapes, building structures and device deployment, the PMAN learns the appropriate transformation and adapts to the new environment with only a small amount of fine-tuning data.

The PMAN takes the spatio-temporal features from two wireless devices as inputs, which are denoted as \mathbf{X}_{ST}^A and \mathbf{X}_{ST}^B respectively. Two features first pass through the fully connected layers $FC(\cdot)$ to get the latent representation for cosine similarity calculation $\hat{\mathbf{O}}$:

$$\hat{\mathbf{O}} = \frac{FC(\mathbf{X}_{ST}^A) \cdot FC(\mathbf{X}_{ST}^B)}{\max(\|FC(\mathbf{X}_{ST}^A)\|_2 \cdot \|FC(\mathbf{X}_{ST}^B)\|_2, \epsilon)}, \quad (12)$$

where ϵ is a small value to avoid division by zero.

In environments where indoor or complex outdoor conditions prevail, the use of Euclidean distance for quantifying device proximity may not accurately reflect the perceived nearness due to potential obstructions, such as walls and buildings. Consequently, our system adopts the Shortest Non-Blocking Distance (SNBD) as the ground truth label. The SNBD represents the minimum path length between wireless devices that does not intersect any physical barriers. This definition assists in aligning the proximity metric more closely with user perception. To optimize model learning during the backpropagation phase, we introduce a novel proximity loss function, $L_P(\cdot)$. This function is computed using a batch of predicted proximity outputs, $\hat{\mathbf{O}}$, and the corresponding SNBD values, \mathbf{d} .

$$L_P(\hat{\mathbf{O}}, \mathbf{d}) = L_{MSE}(\hat{\mathbf{O}}, p(\mathbf{d})) = \frac{1}{N} \|\hat{\mathbf{O}} - p(\mathbf{d})\|_2^2, \quad (13)$$

where $L_{MSE}(\cdot)$ indicates the mean square error, N is the batch size and $p(\mathbf{d}) = \tanh(-\log(\alpha \mathbf{d}))$, where α is the elastic parameter to control the distribution and steepness of the function $p(\cdot)$.

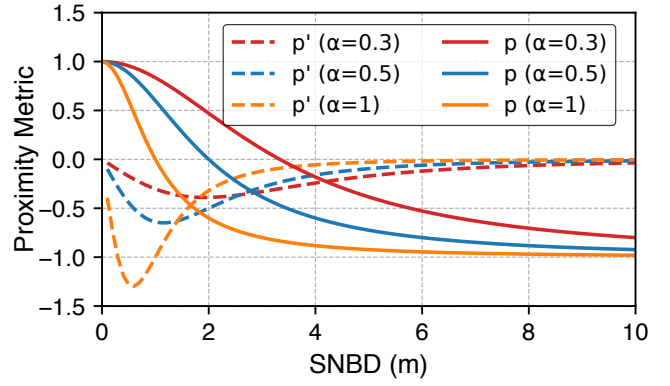


Fig. 5. Illustration of the proximity mapping function varied by α .

The proximity estimation mechanism within our system is designed to be more sensitive to devices in closer proximity than to those further away. Our evaluations indicate that, within indoor settings, relevant proximity distances typically span from 0 to 4 meters. Outdoor scenarios, however, present a more variable range of interest. The model's focus is modulated by an elasticity parameter, α , which adjusts according to the environmental context. For instance, within indoor settings, Fig. 5 demonstrates how different α values influence the mapping function, $p(\cdot)$, and its derivative, $p'(\cdot)$. An α value of 0.5 is chosen to ensure a steep gradient in the function between 0 and 4 meters, with a plateau beyond this range, thereby optimizing the model for indoor applications.

To sum up, the mapping function $p(\cdot)$ guides the model to pay more attention to the distance range of interest, helping PMAN converge quickly and perform better.

B. Transfer Learning

Our goal is to develop a full-scenario proximity detection system for various wireless signals. However, in some complex scenarios, data acquisition becomes very difficult and consumes a lot of manpower and resources. In order to solve this problem, we propose a pre-training and fine-tuning strategy to fill the performance gap from the source domain to the target domain, which greatly enhances the model's domain adaptation capability.

Pre-training in source domain. We build both indoor environments and outdoor scenarios on MATLAB for integrated sensing and communication based on the ray tracing model and collected labeled CSI data under thousands of different deployment cases. pre-training with a large amount of collected data from source domain helps the MRSTE module to learn domain generalized spatio-temporal features without overfitting specific structures. During the pre-training process, all the parameters in *RF-Prox* are jointly optimized. Suppose our pre-training model is M_P . Denote the source domain dataset as \mathcal{D}_S , and the set of parameter in MRSTE and PMAN as Θ_M and Θ_P respectively, the optimization (a.k.a backpropagation) process can be written as

$$\{\Theta_M, \Theta_P\} = \arg \min_{\{\Theta_M, \Theta_P\}} \sum_{(\mathbf{H}, \mathbf{d}) \sim \mathcal{D}_S} \frac{1}{|\mathcal{D}_S|} L_P(M_P(\mathbf{H}; \Theta_M, \Theta_P), \mathbf{d}). \quad (14)$$

Fine-tuning for real-world application. The pre-trained model based on the source domain data has a certain gen-

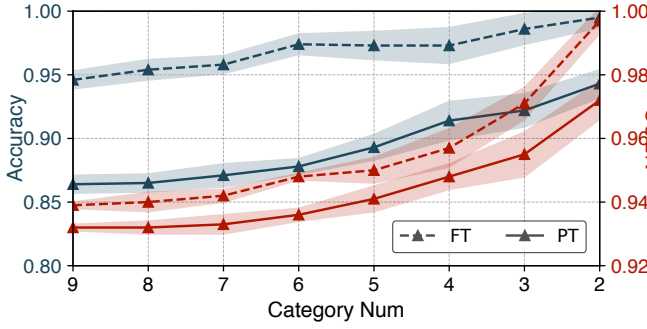


Fig. 6. Top-1 Accuracy and NDCG for Wi-Fi scenario

eralization capability for domain transfer. To achieve better transfer performance from source domain to target domain, a small amount of target domain data is collected to fine-tune the pre-trained model. During the fine-tuning process, the parameters of the MRSTE are frozen and only the parameters of the PMAN are optimized. the optimization process can be written as

$$\Theta_P = \arg \min_{\Theta_P} \sum_{(\mathbf{H}, \mathbf{d}) \sim \mathcal{D}_T} \frac{1}{|\mathcal{D}_T|} L_P(M_F(\mathbf{H}; \Theta_M, \Theta_P), \mathbf{d}), \quad (15)$$

where M_F refers to the fine-tuning model, and \mathcal{D}_T indicates the target domain dataset.

Upon completion of the pre-training and fine-tuning processes, we get a pre-trained model M_P with high generalizability, and a fine-tuned model M_P which adapts to a specific target domain environment.

V. EVALUATION

A. Experimental Methodology

1) *Experimental Scenarios*: To rigorously assess the performance of *RF-Prox*, we conducted comprehensive evaluations through two distinct case studies focused on proximity detection: one within indoor settings involving UEs and IoT devices utilizing Wi-Fi signals, and the other in outdoor UAV scenarios leveraging cellular signals. For these studies, we constructed source and target domains employing the MATLAB Communication Toolbox and the Deep MIMO toolkit [23]. This construction involved setting more than 30 distinct environments, each featuring 300 varied access point/base station (AP/BS) deployment scenarios. Within each scenario, over 20 UEs/UAVs were maneuvered across various locations and orientations to gather a comprehensive dataset of temporal Channel State Information (CSI). For the evaluation within the target domain, we generated novel scenarios with differing AP/BS configurations, where 3-10 UEs/UAVs at varying locations were permitted to move, thereby enabling the collection of corresponding CSI data.

2) *System Implementation*: Within Wi-Fi-enabled indoor settings, *RF-Prox* is configured with one AP and multiple client devices operating at 5.6 GHz, where both the transmitter and receiver are outfitted with three antennas each, spaced at $\lambda/2$, thereby forming a 3×3 antenna array. Conversely, for cellular-based outdoor scenarios, *RF-Prox* encompasses one BS and multiple UAVs operating at 200 GHz, with the BS equipped with a 4×4 antenna array (spaced at $\lambda/2$) and UAVs equipped with a single antenna, culminating in a 1×16 array

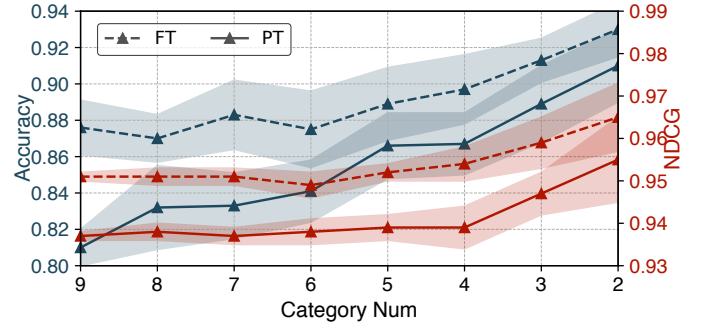


Fig. 7. Top-1 Accuracy and NDCG for cellular scenario

configuration. The target domain's evaluation was conducted under scenarios with a Signal-to-Noise Ratio (SNR) of 30 dB.

RF-Prox employs a hybrid programming approach, utilizing both MATLAB and Python to facilitate rapid and efficient processing. Specifically, MATLAB is utilized for the collection and preprocessing of CSI data, while a python-based deep learning model is leveraged for real-time proximity estimation.

3) *Comparative Methods*: To thoroughly benchmark *RF-Prox*'s efficacy, we juxtaposed it against two state-of-the-art Wi-Fi-based localization methods: SpotFi [5] and mD-Track [8]. This comparison was executed by substituting our MRSTE module with each alternative, thereby highlighting the superiority of our module design.

In scenarios featuring $N+1$ UEs/UAVs, with one randomly selected as the reference, the **Top-1 Accuracy** metric reflects the success rate of identifying the most proximate device amongst the remaining N . Additionally, the **NDCG** [24], a prevalent ranking metric within sorting problems ranging from 0 to 1, evaluates the accuracy of the proximity estimation results in a ranked order.

B. Overall Performance

In this section, we investigate the comprehensive efficacy of *RF-Prox* across both Wi-Fi-enabled indoor environments and cellular-based outdoor scenarios, denoted as W and C, respectively. The performance of the system is examined using both a pre-trained model (PT) and a fine-tuned model (FT).

1) *Top-1 accuracy*: As depicted in Fig. 6 and Fig. 7, the mean accuracy and its variability are illustrated by dots and shaded regions, respectively. The performance of *RF-Prox* in terms of Top-1 accuracy is represented by blue lines, with dashed lines highlighting enhancements post fine-tuning. Initially, accuracy for the pre-trained models (PT-W, PT-C) ranges between 86.4%/81.0% and 94.3%/91.0% for categories decreasing from nine to two, underscoring the system's robust generalizability. After applying domain transfer learning, the fine-tuned models (FT-W, FT-C) demonstrate improved accuracies ranging from 94.6%/87.6% to 99.5%/93.0%, indicating superior adaptation across varied categorical scenarios and effectively bridging the source-target discrepancy through transfer learning.

2) *NDCG*: The NDCG performance of *RF-Prox* is portrayed by red lines in Fig. 6 and Fig. 7, with dashed lines denoting enhancements following fine-tuning. To assess the system's distance-awareness capability, a reference point is chosen at random, and additional UEs are positioned linearly at uniform intervals. The pre-trained models (PT-W,

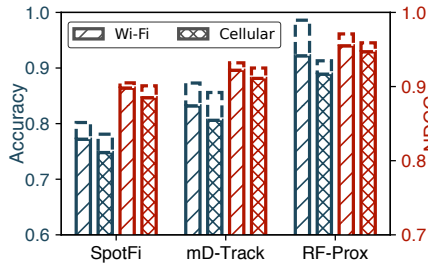


Fig. 8. Comparison with localization methods

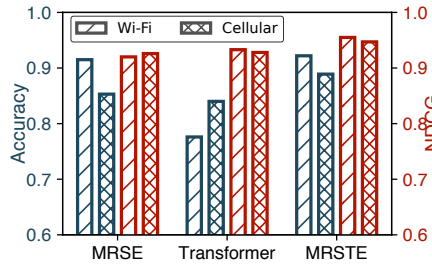


Fig. 9. Effectiveness of each component

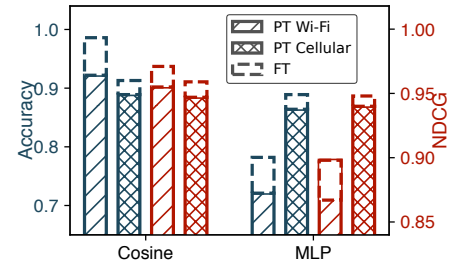


Fig. 10. Comparison of fusion methods

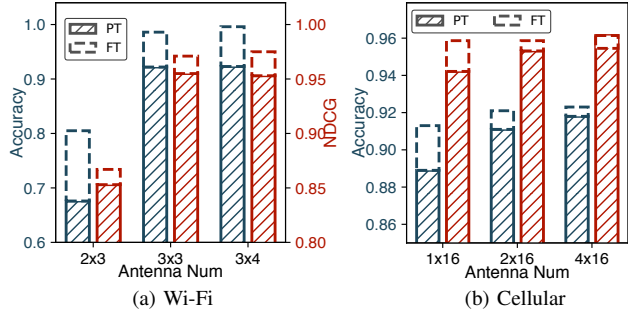


Fig. 11. Impact of antenna number

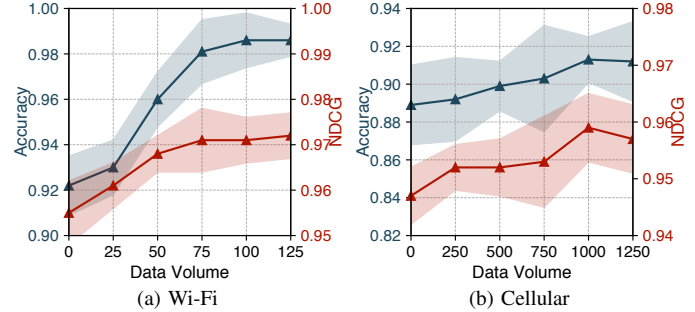


Fig. 12. Impact of fine-tune data volume

PT-C) achieve NDCG scores ranging from 0.932/0.937 to 0.972/0.955 for categories decreasing from nine to two, highlighting the system's generalizability. Post domain transfer learning, the fine-tuned models (FT-W', FT-C') attain NDCG scores between 0.939/0.951 and 0.997/0.965, showcasing the system's potent distance-awareness and the successful mitigation of the source-target gap via transfer learning.

3) *System latency & model parameters*: Utilizing the PyTorch Profiler, we assessed the computational demands, including the floating point operations (FLOPs), model parameters, and inference timing for each component, as documented in Table I. Remarkably, the total number of model parameters is lower than that of typically employed small-scale models (e.g., ResNet-18 [21]), suggesting significant potential for direct deployment on various edge-embedded devices for real-time inference. Additionally, the PMAN's notably smaller parameter count compared to the MRSTE emphasizes efficient fine-tuning and domain adaptation with minimal data volume.

4) *Comparison with localization methods*: To validate the robustness and expressiveness of the high-order spatio-temporal features extracted by MRSTE, we compared its performance to that achieved using multipath AoA and ToF data processed by SpotFi and mD-Track. The pre-trained *RF-Prox* model surpasses both mD-Track and SpotFi in accuracy and NDCG within Wi-Fi-based indoor and cellular-based outdoor scenarios, with respective gains of 9.0%/0.033 and 15.0%/0.057 for indoor, and 8.3%/0.036 and 14.1%/0.062 for outdoor scenarios. These improvements underscore the advanced capabilities of MRSTE in leveraging high-order spatio-temporal features for superior performance. It's important to note that the reported performances are averaged across various common category numbers, maintaining consistency in reporting metrics across the following sections.

TABLE I
SYSTEM LATENCY & NUMBER OF MODEL PARAMETERS

Parameters	Multi-resolution Spatio-Temporal Encoder	Proximity Metric Adaptation Network	Overall
FLOPs (M)	26.12	0.2	26.32
Model Parameters (k)	280.73	18.58	299.31
Inference Time (ms)	13.67	0.13	13.8

C. Component study

In this section, we undertake a component study to evaluate the significance of each module within *RF-Prox*.

1) *MRSTE*: The core of *RF-Prox* comprises two main components: a CNN-based multi-resolution spatial feature extraction module and a transformer-based temporal feature processing module. The spatial features are derived from varying numbers of antennas and subcarriers, while temporal features are extracted from device motion. As illustrated in Fig. 9, *RF-Prox* demonstrates superior performance over both MRSE and Transformer across all metrics in both indoor and outdoor scenarios. This underscores the efficacy of *RF-Prox* in integrating the strengths of both components to enhance spatio-temporal feature extraction, thereby boosting overall performance.

2) *PMAN*: Within the PMAN module, we explore two fusion methods applied to pairs of Channel State Information (CSI): cosine similarity and CSI concatenation followed by Multi-Layer Perceptron (MLP) processing. As depicted in Fig. 10, cosine similarity consistently outperforms the CSI concatenation approach in all evaluated metrics within both environments. This discrepancy is attributed to the fact that concatenating CSIs introduces superfluous sequential data, whereas CSI pairs are inherently unordered and independent. Consequently, the unordered nature of cosine similarity yields superior performance.

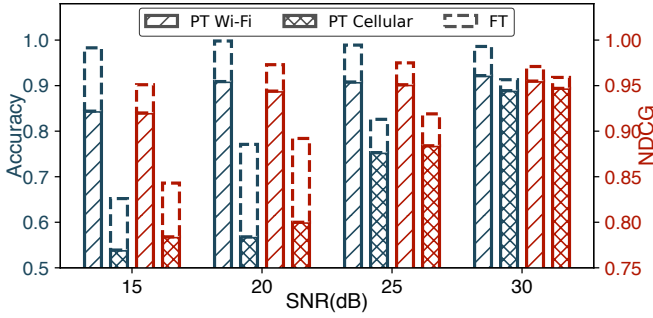


Fig. 13. Performance under different SNR

D. Micro-benchmarks

In this section, we conduct a robustness analysis focusing on antenna number, SNR, and the volume of data used for fine-tuning, to assess their impact on system efficacy.

1) *Antenna Number*: Antenna configurations are varied as 2×3 , 3×3 , and 3×4 arrays for indoor settings, and 1×16 , 2×16 , or 3×16 for outdoor scenarios. As depicted in Fig. 11, both accuracy and NDCG metrics exhibit an upward trend with the increase in the number of antennas in both scenarios. In Wi-Fi-based environments, a 3×3 array provides satisfactory results, achieving accuracy and NDCG of 92.2% / 0.955 with the pre-trained model, and 98.6% / 0.971 post fine-tuning. Similarly, for cellular-based environments, a 1×16 array also shows commendable performance, with accuracy and NDCG of 88.9% / 0.947 (pre-trained) and 91.3% / 0.959 (fine-tuned) respectively. Although additional antennas could potentially improve performance by providing more channel information for enhanced multipath resolution discrimination, the marginal gains diminish beyond a certain point. From a practical perspective, it is prudent to balance the benefits against the costs of using an excessive number of antennas.

2) *SNR*: The SNR settings are established at 15, 20, 25, 30 dB, reflective of typical conditions in communication environments. As shown in Fig. 13, the system maintains robust performance across various SNR levels in Wi-Fi-based indoor environments. For instance, at an SNR of 25 dB, the system achieves an accuracy of $90.8 \pm 1.8\%$ / $98.9 \pm 0.7\%$ and an NDCG of 0.951 ± 0.005 / 0.975 ± 0.005 for the pre-trained and fine-tuned models, respectively. Conversely, performance in cellular-based outdoor environments degrades significantly at lower SNRs, attributed to the complex and dynamic nature of these settings, particularly when unmanned aerial vehicles (UAVs) are involved. Notably, fine-tuning enhances model resilience in lower SNR environments, evidencing the model's adaptability through domain transfer even under challenging conditions.

3) *Fine-tuning data volume*: For fine-tuning, data volumes are set at 0, 25, 50, 75, 100, 125 for indoor environments and 0, 250, 500, 750, 1000, 1250 for outdoor scenarios, with zero data equivalent to employing the pre-trained model. As illustrated in Fig. 12, the system exhibits substantial performance improvements with minimal fine-tuning data. The performance plateau observed at 125 data points indoors and 1250 outdoors suggests scenario-specific data requirements for optimal performance, influenced by device mobility and environmental complexity. The findings affirm the PMAN's robust adaptability across various settings.

VI. RELATED WORK

This section offers an insightful summary of the research landscape surrounding our work.

Wireless-based Localization Techniques. Cellular networks have support for a wide range of positioning methods. For outdoor scenarios, the Cell ID (CID) method [25] leverages the cellular network's awareness of the user equipment's (UE) serving cell to provide basic location insights, albeit with constrained accuracy. Observed Time Difference Of Arrival (OTDOA) [26], [27] employs multilateration, estimating positions through the Time of Arrival (ToA) from several base stations. This technique, reliant on base station infrastructure, exhibits efficacy predominantly in Line-Of-Sight (LOS) situations. Assisted Global Navigation Satellite System (AGNSS) [28] utilizes satellite signal measurements retrieved by systems such as Galileo (Europe) and GPS (US) with high accuracy (i.e. few meters), but AGNSS can be compromised by extreme weather conditions which disrupt satellite signal communication with ground devices [2].

Indoor localization solutions exploit various channel attributes, such as Angle of Arrival (AoA) [5], Time of Flight (ToF) [6], and their fusion [8], aiming for meter-level accuracy. However, these approaches are prone to significant errors in Non-Line-of-Sight (NLoS) settings. Wireless fingerprinting techniques [10], [11] achieve finer accuracy by matching signal features against a pre-compiled database, but their adaptability is limited. Prior research primarily focused on pinpointing the location of individual devices, whereas *RF-Prox* innovatively facilitates proximity assessments between two indirectly connected devices for the first time.

Deep Learning in Wireless Sensing. Deep learning architectures have been extensively applied across a variety of wireless sensing tasks, including gesture [15], [29]–[33] and gait recognition [34]–[38], respiration monitoring [39]–[43], fall detection [44]–[48], and tracking [8], [49]–[51]. Recent innovations have incorporated advanced deep learning concepts such as adversarial [52] and meta-learning [17]. Contrary to the conventional reliance on time-frequency spectrograms as inputs [53], *RF-Prox* represents a pioneering approach by utilizing end-to-end complex-valued neural networks for wireless sensing applications, further enhanced by a transfer learning framework to excel in domain generalization.

VII. CONCLUSION

This paper introduces *RF-Prox*, a novel system designed for the proximity estimation of non-directly connected devices, marking a significant innovation in this domain. Utilizing a sophisticated Multi-Resolution Spatio-Temporal Encoder (MRSTE), *RF-Prox* is capable of extracting domain-agnostic spatio-temporal features from wireless signals. These features are then processed through the Proximity Metric Adaptation Network (PMAN), which converts the extracted latent representations into a set of proximity metrics specifically tailored to the target domain. We implement and evaluate *RF-Prox* on both Wi-Fi-based indoor environments and cellular-based outdoor scenarios. Our results demonstrate that through the incorporation of a transfer learning mechanism, *RF-Prox* efficiently leverages extensive source domain data to learn generalized representations. Moreover, it exhibits remarkable adaptability to new target domains with minimal fine-tuning.

1
2 data. As the inaugural system of its kind, *RF-Prox* represents a
3 pivotal breakthrough in the proximity estimation landscape for
4 non-directly connected devices, offering substantial potential
5 for future applications and research.

6
7 ACKNOWLEDGMENT

8 This work is supported in part by the NSFC under grant
9 62372265, 62271081.

10
11
12 REFERENCES

13 [1] K. Qian, C. Wu, Z. Yang, Y. Liu, and K. Jamieson, "Widar: Decimeter-level passive tracking via velocity monitoring with commodity wi-fi," in *Proceedings of the ACM MobiHoc*, 2017.

14 [2] G. Crowley and I. Azeem, "Chapter 23 - extreme ionospheric storms and their effects on gps systems," in *Extreme Events in Geospace*, N. Buzulukova, Ed. Elsevier, 2018, pp. 555–586. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/B9780128127001000236>

15 [3] A. Bensky, *Wireless positioning technologies and applications*. Artech House, 2016.

16 [4] J. Xiong and K. Jamieson, "Arraytrack: a fine-grained indoor location system," in *Proceedings of the USENIX NSDI*, 2013.

17 [5] M. Kotaru, K. Joshi, D. Bharadia, and S. Katti, "Spotfi: Decimeter level localization using wifi," in *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication*, 2015, pp. 269–282.

18 [6] D. Vasisht, S. Kumar, and D. Katabi, "Decimeter-level localization with a single wifi access point," in *13th USENIX Symposium on Networked Systems Design and Implementation (NSDI 16)*, 2016, pp. 165–178.

19 [7] E. Soltanaghaei *et al.*, "Multipath triangulation: Decimeter-level wifi localization and orientation with a single unaided receiver," in *Proceedings of the ACM MobiSys*, 2018.

20 [8] Y. Xie, J. Xiong, M. Li, and K. Jamieson, "md-track: Leveraging multi-dimensionality for passive indoor wi-fi tracking," in *Proceedings of the ACM MobiCom*, 2019.

21 [9] G. Chi, Z. Yang, J. Xu, C. Wu, J. Zhang, J. Liang, and Y. Liu, "Wi-drone: wi-fi-based 6-dof tracking for indoor drone flight control," in *Proceedings of the ACM MobiSys*, 2022.

22 [10] Z. Yang, C. Wu, and Y. Liu, "Locating in fingerprint space: wireless indoor localization with little human intervention," in *Proceedings of the ACM MobiCom*, 2012.

23 [11] D. Li, J. Xu, Z. Yang, Y. Lu, Q. Zhang, and X. Zhang, "Train once, locate anytime for anyone: Adversarial learning based wireless localization," in *Proceedings of the IEEE INFOCOM*, 2021.

24 [12] J. Shi, M. Sha, and X. Peng, "Adapting wireless mesh network configuration from simulation to reality via deep learning based domain adaptation," in *Proceedings of the USENIX NSDI*, 2021.

25 [13] Y. Gao, G. Chi, G. Zhang, and Z. Yang, "Wi-prox: Proximity estimation of non-directly connected devices via sim2real transfer learning," in *GLOBECOM 2023 - 2023 IEEE Global Communications Conference*, 2023, pp. 5629–5634.

26 [14] Z. Yang, Z. Zhou, and Y. Liu, "From rssi to csi: Indoor localization via channel response," *ACM Computing Surveys (CSUR)*, vol. 46, no. 2, pp. 1–32, 2013.

27 [15] 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.

28 [16] R. Song, D. Zhang, Z. Wu, C. Yu, C. Xie, S. Yang, Y. Hu, and Y. Chen, "Rf-url: unsupervised representation learning for rf sensing," in *Proceedings of the ACM MobiCom*, 2022.

29 [17] S. Ding, Z. Chen, T. Zheng, and J. Luo, "Rf-net: A unified meta-learning framework for rf-enabled one-shot human activity recognition," in *Proceedings of the ACM SenSys*, 2020.

30 [18] C. Trabelsi, O. Bilaniuk, Y. Zhang, D. Serdyuk, S. Subramanian, J. F. Santos, S. Mehri, N. Rostamzadeh, Y. Bengio, and C. J. Pal, "Deep complex networks," in *International Conference on Learning Representations*, 2018.

31 [19] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 31, no. 1, 2017.

32 [20] J. Wang, D. Vasisht, and D. Katabi, "Rf-idraw: Virtual touch screen in the air using rf signals," *ACM SIGCOMM Computer Communication Review*, vol. 44, no. 4, pp. 235–246, 2014.

33 [21] 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.

34 [22] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

35 [23] A. Alkhateeb, "DeepMIMO: A generic deep learning dataset for millimeter wave and massive MIMO applications," in *Proc. of Information Theory and Applications Workshop (ITA)*, San Diego, CA, Feb 2019, pp. 1–8.

36 [24] Y. Wang, L. Wang, Y. Li, D. He, and T.-Y. Liu, "A theoretical analysis of ndcg type ranking measures," in *Conference on learning theory*. PMLR, 2013, pp. 25–54.

37 [25] S. M. Razavi, F. Gunnarsson, H. Rydén, Å. Busin, X. Lin, X. Zhang, S. Dwivedi, I. Siomina, and R. Shreevastav, "Positioning in cellular networks: Past, present, future," in *2018 IEEE Wireless Communications and Networking Conference (WCNC)*. IEEE, 2018, pp. 1–6.

38 [26] M. Huang and W. Xu, "Enhanced lte toa/otdoa estimation with first arriving path detection," in *2013 IEEE Wireless Communications and Networking Conference (WCNC)*. IEEE, 2013, pp. 3992–3997.

39 [27] H. Ryden, A. A. Zaidi, S. M. Razavi, F. Gunnarsson, and I. Siomina, "Enhanced time of arrival estimation and quantization for positioning in lte networks," in *2016 IEEE 27th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*. IEEE, 2016, pp. 1–6.

40 [28] S. Madry *et al.*, *Global navigation satellite systems and their applications*. Springer, 2015.

41 [29] N. Yu, W. Wang, A. X. Liu, and L. Kong, "Qgesture: Quantifying gesture distance and direction with wifi signals," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 2, no. 1, pp. 1–23, 2018.

42 [30] H. Abdelnasser, M. Youssef, and K. A. Harras, "Wigest: A ubiquitous wifi-based gesture recognition system," in *2015 IEEE conference on computer communications (INFOCOM)*. IEEE, 2015, pp. 1472–1480.

43 [31] Q. Pu, S. Gupta, S. Gollakota, and S. Patel, "Whole-home gesture recognition using wireless signals," in *Proceedings of the 19th annual international conference on Mobile computing & networking*, 2013, pp. 27–38.

44 [32] S. Tan and J. Yang, "Wifinger: Leveraging commodity wifi for fine-grained finger gesture recognition," in *Proceedings of the 17th ACM international symposium on mobile ad hoc networking and computing*, 2016, pp. 201–210.

45 [33] R. H. Venkatnarayan, G. Page, and M. Shahzad, "Multi-user gesture recognition using wifi," in *Proceedings of the 16th Annual International Conference on Mobile Systems, Applications, and Services*, 2018, pp. 401–413.

46 [34] Y. Zhang, Y. Zheng, G. Zhang, K. Qian, C. Qian, and Z. Yang, "Gaitsense: towards ubiquitous gait-based human identification with wi-fi," *ACM Transactions on Sensor Networks*, 2021.

47 [35] W. Wang, A. X. Liu, and M. Shahzad, "Gait recognition using wifi signals," in *Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, 2016, pp. 363–373.

48 [36] Y. Zeng, P. H. Pathak, and P. Mohapatra, "Wiwho: Wifi-based person identification in smart spaces," in *2016 15th ACM/IEEE International Conference on Information Processing in Sensor Networks (IPSN)*. IEEE, 2016, pp. 1–12.

49 [37] W. Wang, A. X. Liu, M. Shahzad, K. Ling, and S. Lu, "Understanding and modeling of wifi signal based human activity recognition," in *Proceedings of the 21st annual international conference on mobile computing and networking*, 2015, pp. 65–76.

50 [38] C. Wu, F. Zhang, Y. Hu, and K. R. Liu, "Gaitway: Monitoring and recognizing gait speed through the walls," *IEEE Transactions on Mobile Computing*, vol. 20, no. 6, pp. 2186–2199, 2020.

51 [39] H. Abdelnasser, K. A. Harras, and M. Youssef, "Ubibreathe: A ubiquitous non-invasive wifi-based breathing estimator," in *Proceedings of the 16th ACM international symposium on mobile ad hoc networking and computing*, 2015, pp. 277–286.

52 [40] J. Liu, Y. Wang, Y. Chen, J. Yang, X. Chen, and J. Cheng, "Tracking vital signs during sleep leveraging off-the-shelf wifi," in *Proceedings of the 16th ACM international symposium on mobile ad hoc networking and computing*, 2015, pp. 267–276.

53 [41] X. Wang, C. Yang, and S. Mao, "Tensorbeat: Tensor decomposition for monitoring multiperson breathing beats with commodity wifi," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 9, no. 1, pp. 1–27, 2017.

54 [42] C. Wu, Z. Yang, Z. Zhou, X. Liu, Y. Liu, and J. Cao, "Non-invasive detection of moving and stationary human with wifi," *IEEE Journal on Selected Areas in Communications*, vol. 33, no. 11, pp. 2329–2342, 2015.

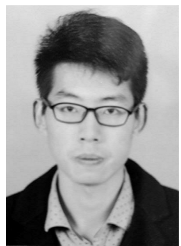
- [43] F. Zhang, C. Wu, B. Wang, M. Wu, D. Bugos, H. Zhang, and K. R. Liu, "Smars: Sleep monitoring via ambient radio signals," *IEEE Transactions on Mobile Computing*, vol. 20, no. 1, pp. 217–231, 2019.
- [44] Z. Yang, Y. Zhang, and Q. Zhang, "Rethinking fall detection with wi-fi," *IEEE Transactions on Mobile Computing*, 2022.
- [45] S. Ji, Y. Xie, and M. Li, "Sifall: Practical online fall detection with rf sensing," in *Proceedings of the 20th ACM Conference on Embedded Networked Sensor Systems*, 2022, pp. 563–577.
- [46] Y. Hu, F. Zhang, C. Wu, B. Wang, and K. R. Liu, "A wifi-based passive fall detection system," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 1723–1727.
- [47] S. Palipana, D. Rojas, P. Agrawal, and D. Pesch, "Falldefi: Ubiquitous fall detection using commodity wi-fi devices," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 1, no. 4, pp. 1–25, 2018.
- [48] Y. Tian, G.-H. Lee, H. He, C.-Y. Hsu, and D. Katabi, "Rf-based fall monitoring using convolutional neural networks," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 2, no. 3, pp. 1–24, 2018.
- [49] F. Adib, Z. Kabelac, and D. Katabi, "Multi-person localization via rf body reflections," in *12th USENIX Symposium on Networked Systems Design and Implementation (NSDI 15)*, 2015, pp. 279–292.
- [50] 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 *Proceedings of the 16th annual international conference on mobile systems, applications, and services*, 2018, pp. 350–361.
- [51] C. Wu, F. Zhang, Y. Fan, and K. R. Liu, "Rf-based inertial measurement," in *Proceedings of the ACM Special Interest Group on Data Communication*, 2019, pp. 117–129.
- [52] W. Jiang, C. Miao, F. Ma, S. Yao, Y. Wang, Y. Yuan, H. Xue, C. Song, X. Ma, D. Koutsonikolas *et al.*, "Towards environment independent device free human activity recognition," in *Proceedings of the ACM MobiCom*, 2018.
- [53] Z. Yang, Y. Zhang, K. Qian, and C. Wu, "Slnet: A spectrogram learning neural network for deep wireless sensing," in *Proceedings of the USENIX NSDI*, 2023.



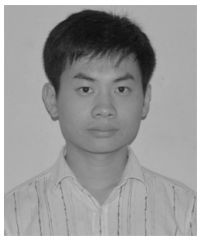
Zhiqing Wei (Member, IEEE) received the B.E. and Ph.D. degrees from the Beijing University of Posts and Telecommunications (BUPT), Beijing, China, in 2010 and 2015, respectively. He is an Associate Professor with BUPT. He has authored one book, three book chapters, and more than 50 papers. His research interests include the performance analysis and optimization of intelligent machine networks. Dr. Wei was granted the Exemplary Reviewer of IEEE Wireless Communications Letters in 2017 and the Best Paper Award of International Conference on Wireless Communications and Signal Processing (WCSP) in 2018. He was the Registration Co-Chair of IEEE/Chinese Institute of Communications (CIC) International Conference on Computer and Communications (ICCC) 2018 and the Publication Co-Chair of IEEE/CIC ICC 2019 and IEEE/CIC ICC 2020.



Yuchong Gao (Student Member, IEEE) received his B.E. degree in School of Information and Communication Engineering from Beijing University of Posts and Telecommunications in 2024. He is currently working towards the Master degree in School of Software, Tsinghua University. His research interests include wireless sensing and mobile computing.



Guoxuan Chi (Member, IEEE) received his B.E. degree in School of Information and Communication Engineering from Beijing University of Posts and Telecommunications in 2019, and the Ph.D. degree in the School of Software, Tsinghua University in 2024. He is currently a research assistant in Tsinghua University. His research interests include wireless sensing and mobile computing.



Zheng Yang (Fellow, IEEE) received the B.E. degree in computer science from Tsinghua University, Beijing, China, in 2006, and the Ph.D. degree in computer science from the Hong Kong University of Science and Technology, Hong Kong, in 2010. He is currently an associate professor with Tsinghua University. His main research interests include Internet of Things and mobile computing. He is the PI of National Natural Science Fund for Excellent Young Scientist. He was the recipient of the State Natural Science Award (second class).



Shijie Cheng received his B.E. degree in Department of Intelligence and Computing from Tianjin University in 2023. He is currently working towards the Master degree in School of Software, Tsinghua University. His research interest is wireless sensing.

The Summary of Difference from the Conference

Version in IEEE GLOBECOM 2023

Preliminary conference paper:

Yuchong Gao, Guoxuan Chi, Guidong Zhang, and Zheng Yang. “Wi-Prox: Proximity Estimation of Non-Directly Connected Devices via Sim2Real Transfer Learning”, IEEE Global Communications Conference (GLOBECOM) 2023.

Main differences:

- (1) We expand the system to support full-scene proximity estimation across different wireless signal types, thereby improving universality, scaling, and facilitating Integrated Sensing and Communication (ISAC) in the 6G era.
 - a) We have modified some expressions in [Section 1](#) about background and motivation:

“Despite the advances in positioning technologies such as the Global Positioning System (GPS), which offers satisfactory outdoor positioning accuracy, its performance is significantly impaired under extreme weather conditions due to interference or obstruction of satellite signals with ground devices. Moreover, rapid urbanization has led to an increased indoor time for individuals, escalating the demand for indoor positioning services across commercial, governmental, and communication sectors. This surge in demand places higher requirements on the coverage of positioning services, highlighting the limitations of GPS in indoor environments and emphasizing the need for integrated solutions suitable for both outdoor and indoor wireless sensing in scenarios where GPS is unavailable.”

- b) We modify the System overview to expand origin simulation-to-reality transfer learning to a more general transfer learning between source and target domains and introduce temporal data into the system in [Section 2](#):

“ RF-Prox is a device proximity estimation system based on the wireless signal, which integrates two pivotal components: the Multi-Resolution Spatio-Temporal Encoder (MRSTE) and the Proximity Metric Adaptation Network (PMAN). As depicted in Fig.2, the workflow of RF-Prox initiates by capturing the Channel State Information (CSI) from wireless links between two distinct mobile devices. This CSI data is then processed through the MRSTE module to extract relevant features. The MRSTE module employs a complex-valued neural network, incorporating residual convolution blocks to derive multi-resolution latent representations from the CSI’s real and imaginary components. Following this, a complex-to-real transformation is performed, converting these complex representations into real-valued spatial features for subsequent analysis. The MRSTE concludes with a transformer-based temporal module, which compresses and extracts latent temporal information.”

“ After processing through MRSTE, the domain-agnostic spatio-temporal features are amalgamated and fed into the PMAN module’s fully connected layers for a comprehensive analysis and comparison. The proximity metric between two devices is then determined using cosine similarity, after transformation by a carefully designed proximity mapping function.

RF-Prox adopts a transfer learning framework, initiating with a model pre-trained in a source domain, which can then be directly applied and fine-tuned within a target domain as necessary. During the pre-training phase, a substantial dataset from the source domain \mathcal{D}_S is utilized to enhance the model’s ability to generalize and extract domain-independent features. After pre-training, fine-tuning the model requires only a minimal dataset from the target domain \mathcal{D}_T . This method enables RF-Prox to be efficiently adapted to new environments, facilitating its practical deployment in varied scenarios.”

- c) We add a related work section in [Section 6](#) to better articulate previous work in the wireless signal-based spatial perception series including localization and tracking, and to facilitate readers’ understanding of the unique contributions of our work:

“ This section offers an insightful summary of the research landscape surrounding our work.

Wireless-based Localization Techniques. Cellular networks have support for a wide range of positioning methods. For outdoor scenarios, the Cell ID (CID) method leverages the cellular network’s awareness of the user equipment’s (UE) serving cell to provide basic location insights, albeit with constrained accuracy. Observed Time Difference Of Arrival (OTDOA) employs multilateration, estimating positions through the Time of Arrival (ToA) from several base stations. This technique, reliant on base station infrastructure, exhibits efficacy predominantly in Line-Of-Sight (LOS) situations. Assisted Global Navigation Satellite System (AGNSS) utilizes satellite signal measurements retrieved by systems such as Galileo (Europe) and GPS (US) with high accuracy (i.e. few meters), but AGNSS can be compromised by extreme weather conditions which disrupt satellite signal communication with ground devices. Indoor localization solutions exploit various channel attributes, such as Angle of Arrival (AoA), Time of Flight (ToF), and their fusion, aiming for meter-level accuracy. However, these approaches are prone to significant errors in Non-Line-of-Sight (NLoS) settings. Wireless fingerprinting techniques achieve finer accuracy by matching signal features against a pre-compiled database, but their adaptability is limited. Prior research primarily focused on pinpointing the location of individual devices, whereas RF-Prox innovatively facilitates proximity assessments between two indirectly connected devices for the first time.

Deep Learning in Wireless Sensing. Deep learning architectures have been extensively applied across a variety of wireless sensing tasks, including gesture and gait recognition, respiration monitoring, fall detection, and tracking. Recent innovations have incorporated advanced deep learning concepts such as adversarial and meta-learning. Contrary to the conventional reliance on time-frequency spectrograms as inputs, RF-Prox represents a pioneering approach by utilizing end-to-end complex-valued neural networks for wireless sensing applications, further enhanced by a transfer learning framework to excel in domain generalization.”

- (2) In our module design, we have seamlessly integrated transformer-based temporal analysis with CNN-based spatial feature extraction, offering improved adaptability to device mobility and bolstered support for analyzing temporal data in practical

applications.

- a) We introduce the temporal data in [Section 2](#) as described above.
- b) We explain in detail how we integrated the transformer into our network in [Section 3.5](#):

“ Through the complex-to-real transformation module C2R, a time-series of real-valued spatial features are extracted. In order to extract the temporal features carried by the device as it moves, we use the transformer encoder module with the attention mechanism. We first embed the series of \mathbf{X}_S into a high-dimensional representation $\mathbf{X}_E = \text{FC}(\mathbf{X}_S)$ with fully connected layers $\text{FC}(\cdot)$. To introduce temporal information, we use absolute positional encoding for \mathbf{X}_E to get embedded positional information \mathbf{P} , with each element as follows:

$$\begin{aligned} \mathbf{P}(\text{pos}, 2i) &= \sin\left(\text{pos}/10000^{2i/d}\right), \\ \mathbf{P}(\text{pos}, 2i+1) &= \cos\left(\text{pos}/10000^{2i/d}\right), \end{aligned}$$

where pos represents the sequence element's ordinal position, and i refers to the dimension index within the embedding space. Then, we formulate the transformer's encoded input as $\mathbf{X}_{PE} = \mathbf{X}_E + \mathbf{P}$.

In order to better learn the relationships between the elements inside the sequence, we use the attention mechanism to linearly transform the input \mathbf{X}_{PE} into queries, keys and values:

$$\mathbf{Q} = \mathbf{X}_{PE} \mathbf{W}_Q,$$

$$\mathbf{K} = \mathbf{X}_{PE} \mathbf{W}_K,$$

$$\mathbf{V} = \mathbf{X}_{PE} \mathbf{W}_V,$$

where \mathbf{W}_Q , \mathbf{W}_K and \mathbf{W}_V are the weight matrices used for linear transformation. Attention features \mathbf{X}_A are obtained as

$$\begin{aligned} \mathbf{X}_A &= \text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \\ &= \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right) \mathbf{V}, \end{aligned}$$

where d_k is the dimension of keys.”

“ After residual connection and layer normalization, the spatio-temporal features \mathbf{X}_{ST} in the CSI can be finally extracted as

$$\mathbf{X}_{ST} = \text{LayerNorm}(\mathbf{X}_A + \mathbf{X}_{PE}),$$

which encapsulates the spatio-temporal relationship between the device and the access point. Thus, the spatio-temporal features corresponding to two different terminal devices can be further utilized to determine their proximity relationship. ”

- c) We modified the description of the MRSTE model in the contribution to emphasize the integration of the timing model in [Section 1](#):

“ Our proposed Multi-Resolution Spatio-Temporal Encoder is a pioneering attempt at applying complex-valued neural networks to wireless sensing. The multi-resolution design and transformer-based temporal processing model have unique advantages and can also be integrated into other types of wireless sensing applications. ”

(3) We have improved some of the modules to increase the performance of the models.

- a) We replaced the previous ordered MLP-based feature mapping module with an unordered cosine similarity module in [Section 5.3](#):

“ Within the PMAN module, we explore two fusion methods applied to pairs of Channel State Information (CSI): cosine similarity and CSI concatenation followed by Multi-Layer Perceptron (MLP) processing. As depicted in Fig.10, cosine similarity consistently outperforms the CSI concatenation approach in all evaluated metrics within both environments. This discrepancy is attributed to the fact that concatenating CSIs introduces superfluous sequential data, whereas CSI pairs are inherently unordered and independent. Consequently, the unordered nature of cosine similarity yields superior performance. ”

- b) We extend the transfer learning of simulation-to-reality to a more pervasive transfer task under source and target domains in [Section 4.2](#):

“ Our goal is to develop a full-scenario proximity detection system for various wireless signals. However, in some complex scenarios, data acquisition becomes very difficult and consumes a lot of manpower and resources. In order to solve this problem, we propose a pre-training and fine-tuning strategy to fill the performance gap from the source domain to the target domain, which greatly enhances the model’s domain adaptation capability.

Pre-training in source domain. We build both indoor environments and outdoor scenarios on MATLAB for integrated sensing and communication based on the ray tracing model and collected labeled CSI data under thousands of different deployment cases. Pre-training with a large amount of collected data from source domain helps the MRSTE module to learn domain generalized spatio-temporal features without overfitting specific structures.

During the pre-training process, all the parameters in RF-Prox are jointly optimized. Suppose our pre-training model is M_P . Denote the source domain dataset as \mathcal{D}_S , and the set of parameter in MRSTE and PMAN as Θ_M and Θ_P respectively, the optimization (a.k.a backpropagation) process can be written as

$$\{\Theta_M, \Theta_P\} = \arg \min_{\{\Theta_M, \Theta_P\}} \sum_{(\mathbf{H}, \mathbf{d}) \sim \mathcal{D}_S} \frac{1}{|\mathcal{D}_S|} L_P(M_P(\mathbf{H}; \Theta_M, \Theta_P), \mathbf{d}).$$

Fine-tuning for real-world application. The pre-trained model based on the source domain data has a certain generalization capability for domain transfer. To achieve better transfer performance from source domain to target domain, a small amount of target domain data is collected to fine-tune the pre-trained model. During the fine-tuning process, the parameters of the MRSTE are frozen and only the parameters of the PMAN are optimized. the optimization process can be written as”

“

$$\Theta_P = \arg \min_{\Theta_P} \sum_{(\mathbf{H}, \mathbf{d}) \sim \mathcal{D}_T} \frac{1}{|\mathcal{D}_T|} L_P(M_F(\mathbf{H}; \Theta_M, \Theta_P), \mathbf{d}),$$

where M_F refers to the fine-tuning model, and \mathcal{D}_T indicates the target domain dataset.

Upon completion of the pre-training and fine-tuning processes, we get a pre-trained model M_P with high generalizability, and a fine-tuned model M_P which adapts to a specific target domain environment. ”

- (4) To validate the above enhancements, we have comprehensively revised all experiments, incorporated cellular-based outdoor scenarios, redesigned and retrained our models, ensuring they are fine-tuned for the new case. We’ve also conducted a full evaluation and optimization of these additions.
- a) We tried different optimizers, learning rates, and learning rate decay strategies to tune the model for the best performance in [Section 5](#).
 - b) We used the newly proposed MRSTE model to test the overall performance in 8 classification situations in the newly proposed cellular-based outdoor scenarios and the previously mentioned Wi-Fi-based indoor environments in [Section 5.2](#):

“In this section, we investigate the comprehensive efficacy of RF-Prox across both Wi-Fi-enabled indoor environments and cellular-based outdoor scenarios, denoted as W and C, respectively. The performance of the system is examined using both a pre-trained model (PT) and a fine-tuned model (FT).

1) *Top-1 accuracy*: As depicted in Fig.6 and Fig.7, the mean accuracy and its variability are illustrated by dots and shaded regions, respectively. The performance of RF-Prox in terms of Top-1 accuracy is represented by blue lines, with dashed lines highlighting enhancements post fine-tuning. Initially, accuracy for the pre-trained models (PT-W, PT-C) ranges between 86.4%/81.0% and 94.3%/91.0% for categories decreasing from nine to two, underscoring the system’s robust generalizability. After applying domain transfer learning, the fine-tuned models (FT-W, FT-C) demonstrate improved accuracies ranging from 94.6%/87.6% to 99.5%/93.0%, indicating superior adaptation across varied categorical scenarios and effectively bridging the source-target discrepancy through transfer learning.

2) *NDCG*: The NDCG performance of RF-Prox is portrayed by red lines in Fig.6 and Fig.7, with dashed lines denoting enhancements following fine-tuning. To assess the system’s distance-awareness capability, a reference point is chosen at random, and additional UEs are positioned linearly at uniform intervals. The pre-trained models (PT-W, PT-C) achieve NDCG scores ranging from 0.932/0.937 to 0.972/0.955 for categories decreasing from nine to two, highlighting the system’s generalizability. Post domain transfer learning, the fine-tuned models (FT-W, FT-C) attain NDCG scores between 0.939/0.951 and 0.997/0.965, showcasing the system’s potent distance-awareness and the successful mitigation of the source-target gap via transfer learning.”

- c) We tested the newly proposed MRSTE model’s parameters such as parameter size and inference experiments to prove its ability to perform real-time inference in [Section 5.2](#):

“3) System latency & model parameters: Utilizing the PyTorch Profiler, we assessed the computational demands, including the floating point operations (FLOPs), model parameters, and inference timing for each component, as documented in Table ???. Remarkably, the total number of model parameters is lower than that of typically employed small-scale models (e.g., ResNet-18), suggesting significant potential for direct deployment on various edge-embedded devices for real-time inference. Additionally, the PMAN’s notably smaller parameter count compared to the MRSTE emphasizes efficient fine-tuning and domain adaptation with minimal data volume.”

TABLE I: System latency & number of model parameters

Parameters	Multi-resolution	Proximity Metric	Overall
	Spatio-Temporal Encoder	Adaptation Network	
FLOPs (M)	26.12	0.2	26.32
Model Parameters (k)	280.73	18.58	299.31
Inference Time (ms)	13.67	0.13	13.8

- d) We tested the performance improvement of the newly proposed model compared to the traditional state-of-the-art SpotFi and mD-Track in both indoor and outdoor scenarios in [Section 5.2](#):

“4) Comparison with localization methods: To validate the robustness and expressiveness of the high-order spatio-temporal features extracted by MRSTE, we compared its performance to that achieved using multipath AoA and ToF data processed by SpotFi and mD-Track. The pre-trained RF-Prox model surpasses both mD-Track and SpotFi in accuracy and NDCG within Wi-Fi-based indoor and cellular-based outdoor scenarios, with respective gains of 9.0%/0.033 and 15.0%/0.057 for indoor, and 8.3%/0.036 and 14.1%/0.062 for outdoor scenarios. These improvements underscore the advanced capabilities of MRSTE in leveraging high-order spatio-temporal features for superior performance. It’s important to note that the reported performances are averaged across various common category numbers, maintaining consistency in reporting metrics across the following sections.”

- e) By conducting component studies on the newly proposed MRSTE and PMAN, we proved the excellent spatio-temporal feature extraction capability of the MRSTE module and the excellent performance of cosine similarity in PMAN in [Section 5.3](#):

“1) MRSTE: The core of RF-Prox comprises two main components: a CNN-based multi-resolution spatial feature extraction module and a transformer-based temporal feature processing module. The spatial features are derived from varying numbers of antennas and subcarriers, while temporal features are extracted from device motion. As illustrated in Fig.9, RF-Prox demonstrates superior performance over both MRSE and Transformer across all metrics in both indoor and outdoor scenarios. This underscores the efficacy of RF-Prox in integrating the strengths of both components to enhance spatio-temporal feature extraction, thereby boosting overall performance.

2) PMAN: Within the PMAN module, we explore two fusion methods applied to pairs of Channel State Information (CSI): cosine similarity and CSI concatenation followed by Multi-Layer Perceptron (MLP) processing. As depicted in Fig.10, cosine similarity consistently outperforms the CSI concatenation approach in all evaluated metrics within both environments. This discrepancy is attributed to the fact that concatenating CSIs introduces superfluous sequential data, whereas CSI pairs are inherently unordered and independent. Consequently, the unordered nature of cosine similarity yields superior performance.”

f) we conduct a robustness analysis focusing on antenna number, SNR, and the volume of data used for fine-tuning, to assess their impact on system efficacy in [Section 5.4](#):

“1) Antenna Number: Antenna configurations are varied as 2×3 , 3×3 , and 3×4 arrays for indoor settings, and 1×16 , 2×16 , or 3×16 for outdoor scenarios. As depicted in Fig.11, both accuracy and NDCG metrics exhibit an upward trend with the increase in the number of antennas in both scenarios. In Wi-Fi-based environments, a 3×3 array provides satisfactory results, achieving accuracy and NDCG of 92.2% / 0.955 with the pre-trained model, and 98.6% / 0.971 post fine-tuning. Similarly, for cellular-based environments, a 1×16 array also shows commendable performance, with accuracy and NDCG of 88.9% / 0.947 (pre-trained) and 91.3% / 0.959 (fine-tuned) respectively. Although additional antennas could potentially improve performance by providing more channel information for enhanced multipath resolution discrimination, the marginal gains diminish beyond a certain point. From a practical perspective, it is prudent to balance the benefits against the costs of using an excessive number of antennas.

2) SNR: The SNR settings are established at 15, 20, 25, 30 dB, reflective of typical conditions in communication environments. As shown in Fig.13, the system maintains robust performance across various SNR levels in Wi-Fi-based indoor environments. For instance, at an SNR of 25 dB, the system achieves an accuracy of $90.8 \pm 1.8\%$ / $98.9 \pm 0.7\%$ and an NDCG of 0.951 ± 0.005 / 0.975 ± 0.005 for the pre-trained and fine-tuned models, respectively. Conversely, performance in cellular-based outdoor environments degrades significantly at lower SNRs, attributed to the complex and dynamic nature of these settings, particularly when unmanned aerial vehicles (UAVs) are involved. Notably, fine-tuning enhances model resilience in lower SNR environments, evidencing the model's adaptability through domain transfer even under challenging conditions.”

3) Fine-tuning data volume: For fine-tuning, data volumes are set at 0, 25, 50, 75, 100, 125 for indoor environments and 0, 250, 500, 750, 1000, 1250 for outdoor scenarios, with zero data equivalent to employing the pre-trained model. As illustrated in Fig.12, the system exhibits substantial performance improvements with minimal fine-tuning data. The performance plateau observed at 125 data points indoors and 1250 outdoors suggests scenario-specific data requirements for optimal performance, influenced by device mobility and environmental complexity. The findings affirm the PMAN's robust adaptability across various settings.

g) To foster transparency and collaboration, we have made all related code publicly available¹.

REMARK: Please note that all of the figures and tables below are **NEW** and different from the original paper.

We believe that the above aspects significantly set this paper apart from the previous conference paper.

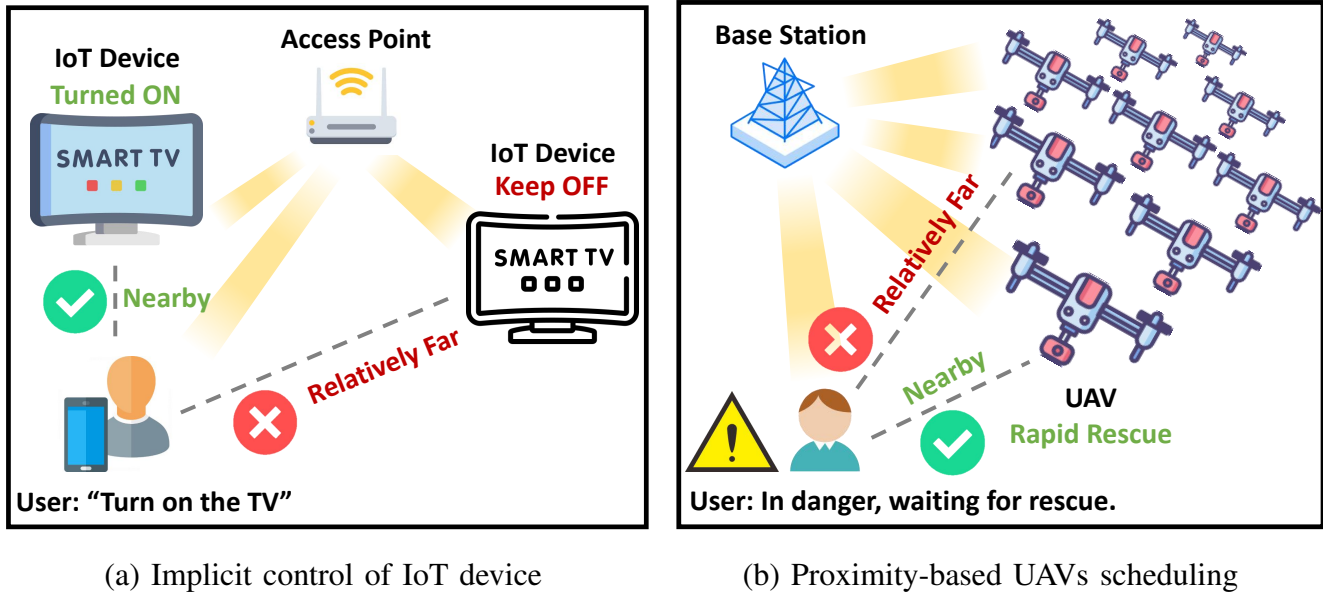


Fig. 1: Illustration of two application scenarios of *RF-Prox*.

¹Our project is available [here](#).

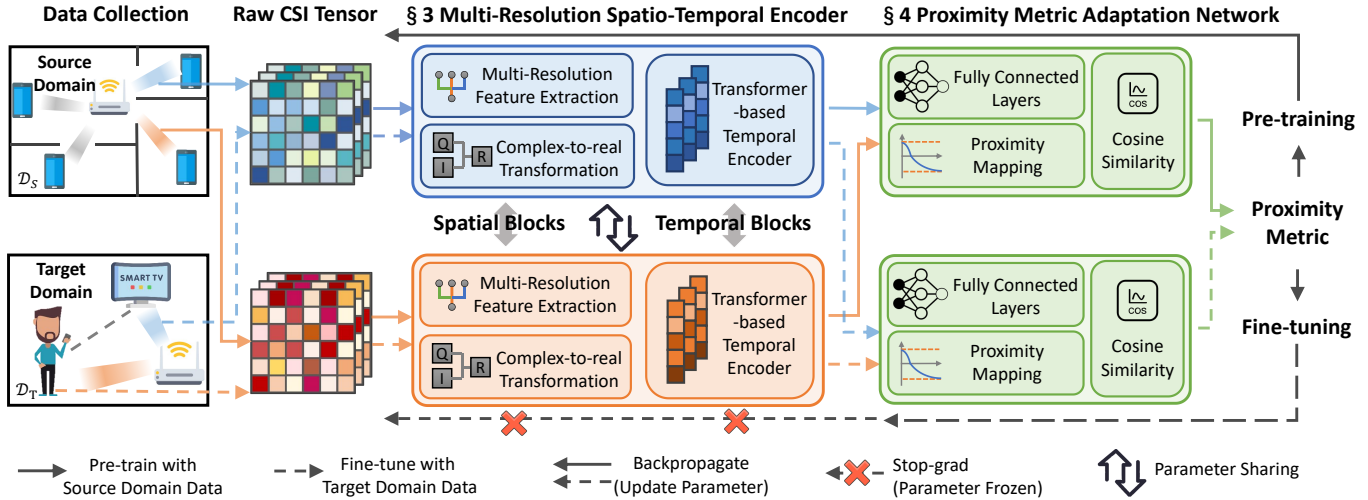


Fig. 2: An overview of the *RF-Prox*, where solid and dashed lines represent data collection from source domain and target domain, respectively, with blue and orange used to distinguish devices.

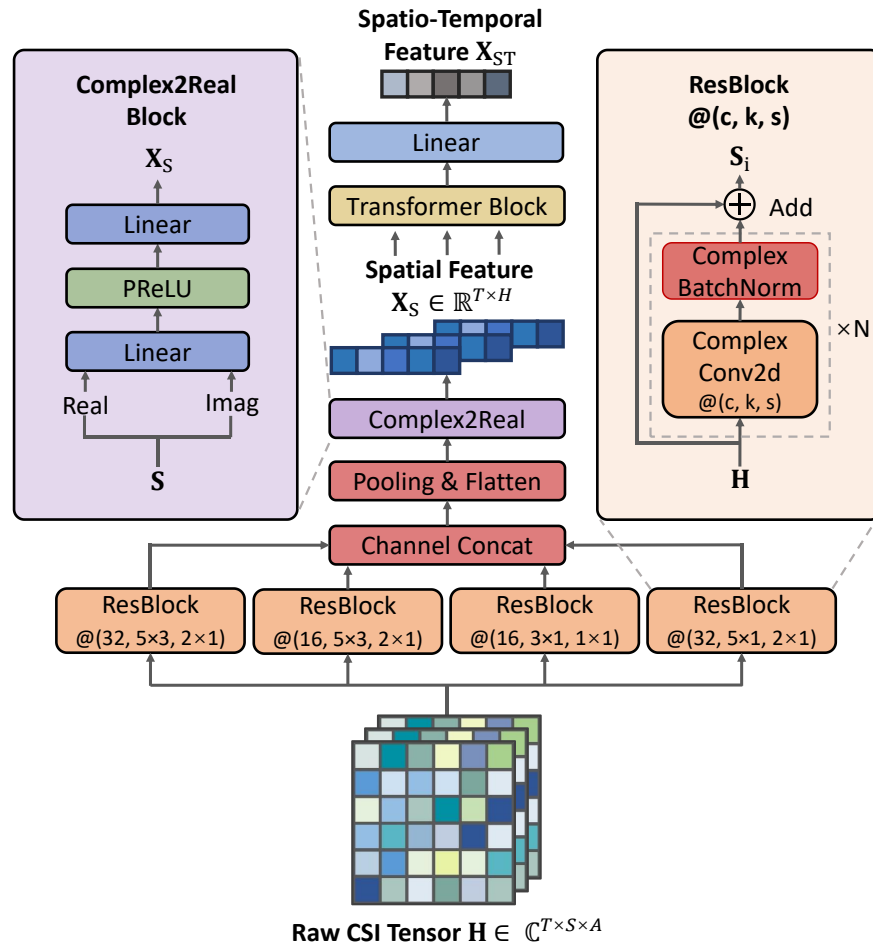


Fig. 3: Illustration of Multi-Resolution Spatio-Temporal Encoder.

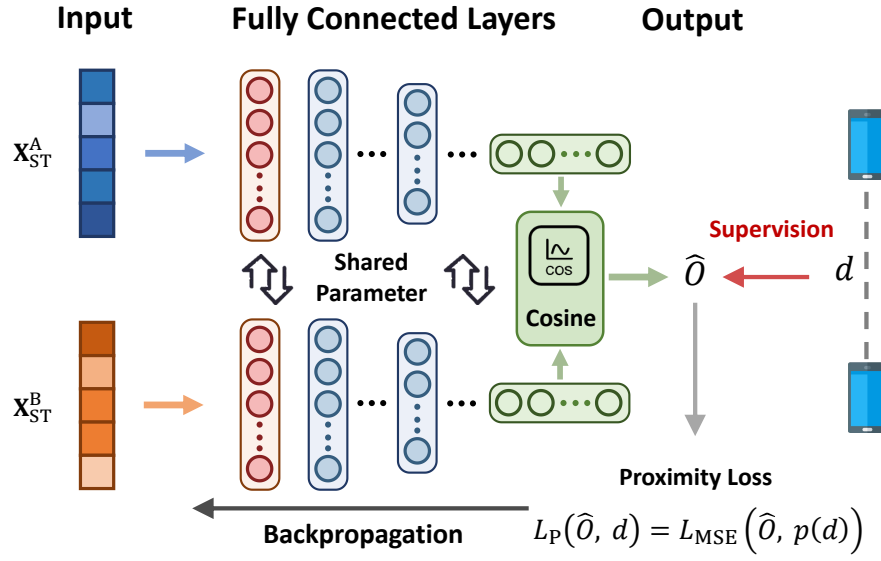


Fig. 4: An overview of the Proximity Metric Adaptation Network.

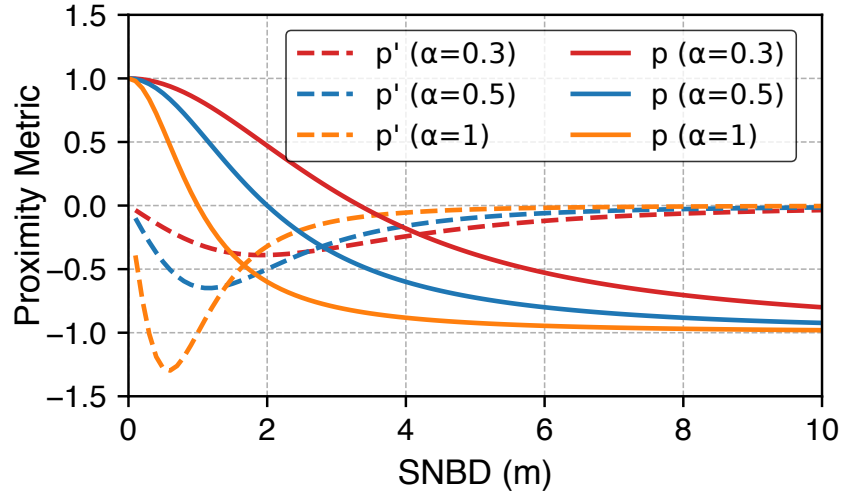


Fig. 5: Illustration of the proximity mapping function varied by α .

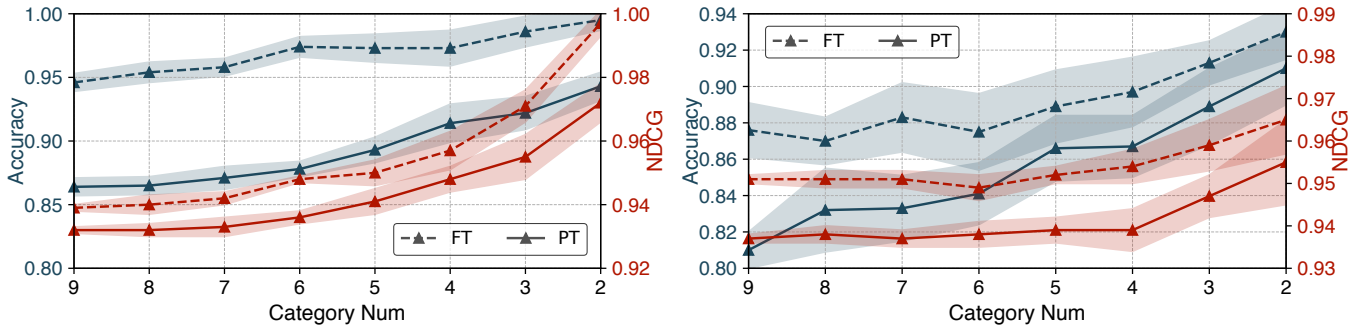


Fig. 6: Top-1 Accuracy and NDCG for Wi-Fi sce-
nario

Fig. 7: Top-1 Accuracy and NDCG for cellular sce-
nario

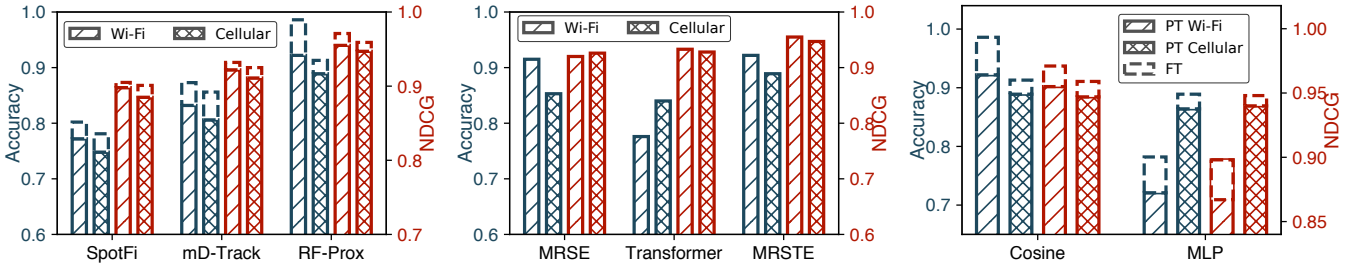


Fig. 8: Comparison with localization methods. Fig. 9: Effectiveness of each component. Fig. 10: Comparison of fusion methods.

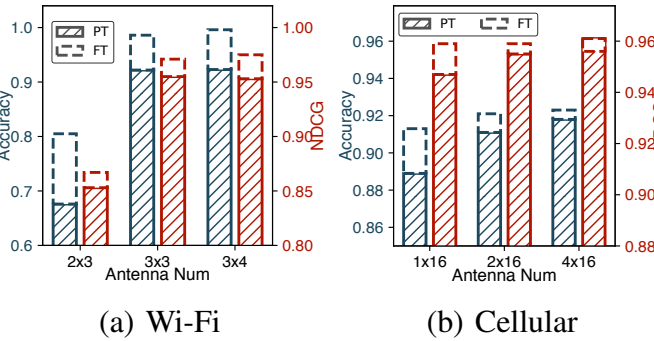


Fig. 11: Impact of antenna number

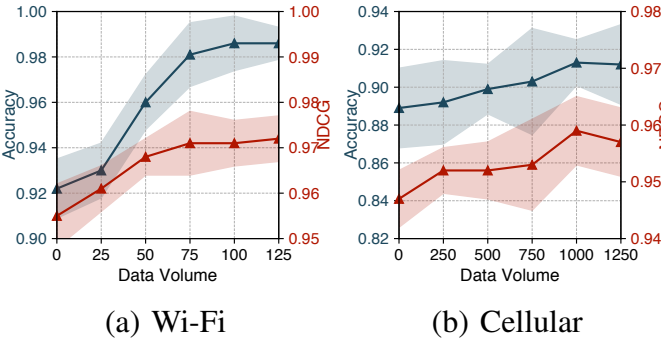


Fig. 12: Impact of fine-tune data volume

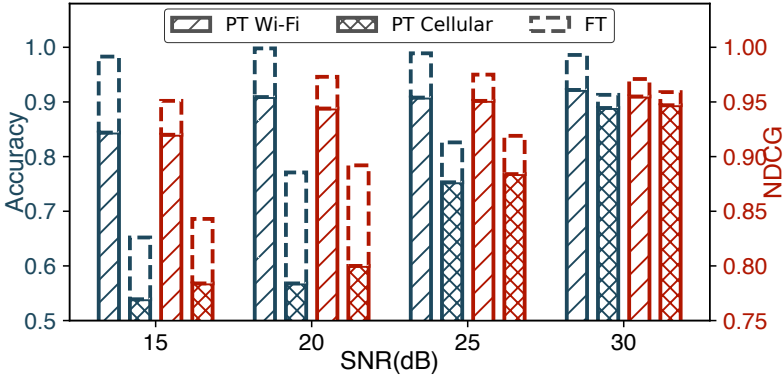


Fig. 13: Performance under different SNR