

Sub-Resolution mmWave FMCW Radar-based Touch Localization using Deep Learning

Raghunandan M. Rao, Amit Kachroo, Koushik A. Manjunatha, Morris Hsu, Rohit Kumar

Abstract—Touchscreen-based interaction on display devices are ubiquitous nowadays. However, capacitive touch screens, the core technology that enables its widespread use, are prohibitively expensive to be used in large displays because the cost increases proportionally with the screen area. In this paper, we propose a millimeter wave (mmWave) radar-based solution to achieve sub-resolution error performance using a network of four mmWave radar sensors. Unfortunately, achieving this is non-trivial due to inherent range resolution limitations of mmWave radars, since the target (human hand, finger etc.) is ‘distributed’ in space. We overcome this using a deep learning-based approach, wherein we train a deep convolutional neural network (CNN) on range-FFT (range vs power profile)-based features against ground truth (GT) positions obtained using a capacitive touch screen. To emulate the clutter characteristics encountered in radar-based positioning of human fingers, we use a metallic finger mounted on a metallic robot arm as the target. Using this setup, we demonstrate sub-resolution position error performance. Compared to conventional signal processing (CSP)-based approaches, we achieve a $2 - 3\times$ reduction in positioning error using the CNN. Furthermore, we observe that the inference time performance and CNN model size support real-time integration of our approach on general purpose processor-based computing platforms.

Index Terms—mmWave radar, deep neural network, sub-resolution touch localization, large displays.

I. INTRODUCTION

Modern displays use capacitive touchscreens for enabling touch-based interaction with the device, wherein touch localization is performed by processing the changes in electrical properties of the touchscreen layers across the display [1]. In general, the touchscreen cost scales linearly with the area of the display covered by the touchscreen. As a result, it becomes prohibitively expensive to use capacitive touchscreens in large displays (for instance, display size > 15 inch). Furthermore, since the size of interactive elements (icons, sliders, buttons, etc.) tend to be large on a large screen, the positioning error requirement can often be relaxed from the typical mm-level accuracy to a few cm, without significantly impacting the user experience (UX). This work is focused on enabling accurate touch positioning in the latter regime.

A. Related Work

To reduce the cost while providing accurate positioning performance, there is significant interest in using alternative technologies such as Ultrasound [2], WiFi [3], Radio Frequency [4], [5], mmWave [6], [7] and Ultrawideband

(UWB) [8]–[10] radar. Yun *et al.* [2] designed a device-free system using ultrasound signals to track human finger motion. Their algorithm is based on processing the channel impulse response (CIR) to estimate the absolute distance and the distance displacement using multiple CIRs, resulting in a median tracking error $\delta r_{\text{track},50} = 1$ cm. Wu *et al.* [3] proposed a sub-wavelength finger motion tracking system using one WiFi transmitter and two WiFi receivers. Their approach used a channel quotient-based feature to detect minute changes in the channel state information (CSI) due to finger movement, resulting in $\delta r_{\text{track},90} = 6$ cm. Wang *et al.* [4] propose a Radio Frequency Identification (RFID) sensor worn on the finger to demonstrate precise tracking of air handwriting gestures. The authors demonstrated tracking errors of $\delta r_{\text{track},90} = 9.7$ cm and 10.5 cm in Line-of-Sight (LoS) and Non-LoS (NLoS) conditions respectively. Shang-guan [5] proposed an air handwriting system based on two differentially polarized antennas to track the orientation and position of an RFID-tagged pen, achieving a 90th percentile position error ($\delta r_{\text{pos},90}$) of 11 cm. Xiao *et al.* [7] proposed an RF backscattering-based system to track handwriting traces performed using a stylus in which a RFID tag is embedded. The authors demonstrated $\delta r_{\text{track},50} = 0.49$ cm at a writing speed of 30 cm/s. Wei [6] designed a 60 GHz radio-based passive tracking system to position different writing objects such as pen/pencil/marker to obtain $\delta r_{\text{pos},90} = 0.8/5/15$ cm respectively. Their approach relies on passive backscattering of a single carrier 60 GHz signal, using which the initial location is acquired. The low tracking error is obtained by tracking its phase over time. In [8], the authors propose an Inertial Measurement Unit (IMU)-UWB radar fusion-based tracking approach to implement a stylus-aided handwriting use-case that achieves $\delta r_{\text{pos},50} = 0.49$ cm. However, in the absence of the IMU, the authors report $\delta_{\text{pos},90} = 6$ cm.

Even though the works [2]–[4], [7] report a low tracking error, the position error is high. While this trade-off is acceptable in finger tracking applications such as handwriting recognition, it is unacceptable for on-screen interaction where the performance is dictated by the position error, not by the accuracy of the reconstructed trajectory. On the other hand, works such as [8] that achieve cm-level position accuracy necessitates the use of additional IMU sensors, that adds system/computational complexity, and friction to the UX.

B. Contributions

In this work, we bridge this gap by proposing a mmWave radar sensor network-based positioning framework that uses a

R. M. Rao, K. A. Manjunatha, M. Hsu, and R. Kumar are with Amazon Lab126, Sunnyvale, CA, USA, 94089 (email: {raghmrao, koushiam}@amazon.com, mhsu@lab126.com, rrohk@amazon.com).

A. Kachroo is with Amazon Web Services (AWS), Santa Clara, CA, 95054 (email: amkachro@amazon.com).

Deep Convolutional Neural Network (CNN) to achieve sub-resolution position accuracy. We build a robot-based testbed for characterizing positioning performance, in which we use a robot-mounted metal finger as the distributed target. We collect data for multiple runs to capture the metal finger’s signature for each radar sensor at different locations of the screen, as well as the corresponding ground truth position using a capacitive touchscreen. The conventional signal processing (CSP)-based approach that uses range calibration coupled with the nonlinear least squares (NLS) algorithm [11] results in $\delta r_{\text{pos},90} = 3.7$ cm. On the other hand, inspired by the LeNet model [12], we design a CNN that outperforms the CSP-based approach and achieves sub-resolution accuracy, with $\delta r_{\text{pos},90} = 1.6$ cm. Finally, we demonstrate that the small model size and CNN inference time makes real-time implementation feasible on general purpose processor-based computing platforms.

II. SYSTEM DESIGN

A. Working Principle

In contemporary consumer electronic devices, touch localization on a capacitive touchscreen displays rely on changes in the electrical properties of carefully arranged material layers when a finger touches the screen. In essence, the location is estimated by determining the ‘touch cell’ where there is maximum variation in the capacitance [13]. In contrast, contactless methods can also be used if accurate distance [14]–[16] and/or angle information [17] of the finger is known relative to sensors *with known positions*. In this work, we design a low-cost touch positioning system that uses multiple Frequency Modulated Continuous Wave (FMCW) mmWave radar sensors to locate the “finger” (target) on the screen. FMCW mmWave radar sensors are attractive for short range sensing applications because they can be operated at low-power and manufactured with low cost. This is in part due to the low bandwidth (~ 1 MHz) of the baseband signal processing chain, despite the large sweep bandwidth (> 1 GHz) [18].

However, the main factor that limits accurate finger localization is the limited range resolution (Δr_{res}) of the radar, given by $\Delta r_{\text{res}} = \frac{c}{2f_{\text{BW}}}$, where c is the free-space velocity of light, and f_{BW} the chirp bandwidth of the FMCW radar [19]. For example, the radar will not be able to distinguish objects that are closer than 3 cm (in the radial direction) for $f_{\text{BW}} = 5$ GHz. As a result, the finger and the rest of the hand often appear as a single target to each radar sensor, thus resulting in a range error that is dictated by Δr_{res} .

B. Experimental Testbed Setup

The main focus of this work is to evaluate the position error of a mmWave radar-based touch solution. To undertake this, we built a testbed whose schematic is shown in Fig. 1a. The setup is based on a 15.6 inch display, on which $N_{\text{rad}} = 4$ radar sensors are placed at the corners of the screen using 3D printed fixtures, as shown in Fig. 1b. The display is equipped with a capacitive touchscreen, which is interfaced with a dedicated computer to obtain the ground truth (position and time) for

TABLE I: Radar Sensor Parameters

| Parameter | Value |
|--|-----------|
| Waveform type | FMCW |
| Chirp Bandwidth (f_{BW}) | 4.874 GHz |
| Range Resolution (Δr) | 3.075 cm |
| Frame Rate (f_r) | 120 Hz |
| Number of IF samples/chirp (N_{IF}) | 64 |
| Number of chirps/frame (N_{ch}) | 8 |
| Number of RX antennas/sensor (N_{rx}) | 3 |
| Number of radar sensors (N_{rad}) | 4 |

each touch event. A metal finger is used as the target which is to be localized by the radar sensor network. The metal finger is mounted on a programmable robot to achieve precise control over its trajectory during the data collection session, as shown in Fig. 1a. The robot is controlled by another dedicated computer, and is programmed to touch the display on a grid of points, as shown in Fig. 2. The spacing between each point on the grid is approximately 1 cm along both vertical and horizontal directions. It is important to note that localization performance in this setup is typically limited by Δr_{res} , since the metal finger (target, analogous to the human finger) and the metallic robotic arm (analogous to the rest of the hand) on which it is mounted *will appear as a single target* to the radar. The radar configuration used is shown in Table I. It is worthwhile to note that this waveform is compliant with the latest FCC final rule [20].

C. Radar Signal Pre-Processing

The signal processing pipeline for generating the feature is shown in Fig. 3. Let f, s, r, c, j and i denote the frame index, IF sample index, range bin, chirp index, RX antenna index, and the sensor index respectively. Each radar sensor transmits the FMCW waveform with parameters shown in Table I. For each frame f , the received waveform is then down-converted to get the intermediate frequency (IF) signal $x_{\text{IF}}(f, s, c, j, i)$, that corresponds to the radar return. From this, the range information corresponding to all scattering objects in the radar’s field of view (FoV) is obtained by computing the beamformed ‘range-FFT’ $x_r(f, r, c, j)$ using the following sequence of steps.

$$x_{\text{IF},\text{zp}}(f, s, c, j, i) = \begin{cases} x_{\text{IF}}(f, s, c, j, i) & \text{for } s < N_{\text{IF}} \\ 0 & \text{for } N_{\text{IF}} \leq s \leq N_{\text{os}}N_{\text{IF}} - 1 \end{cases}, \quad (1)$$

$$x_w(f, s, c, j, i) = x_{\text{IF},\text{zp}}(f, s, c, j, i)w_{\text{IF}}(s), \quad (2)$$

$$x_r(f, r, c, j, i) = \sum_{s=0}^{N_{\text{os}}N_{\text{IF}}-1} x_w(f, s, c, j, i)e^{\frac{-j2\pi sr}{N_{\text{os}}N_{\text{IF}}}}, \quad (3)$$

for $r = 0, 1, \dots, \left(\frac{N_{\text{os}}N_{\text{IF}}}{2} - 1\right)$. Here, zero-padding is performed in (1) to shrink the effective range-bin width to $\Delta r_{\text{os}} = \Delta r/N_{\text{os}}^1$, where $N_{\text{os}} = 8$ is the oversampling factor. The zero-padded signal is then used to compute the range-FFT in (3) after a windowing operation $w_{\text{IF}}(\cdot)$. The purpose of the latter is to trade-off the range-FFT sidelobe level with

¹Note that while this shrinks the range bin width, the range resolution (i.e. minimum radial distance between two targets such that they appear as two distinct targets) is unchanged. Oversampling in the range domain minimizes the contribution of range quantization error in the system.

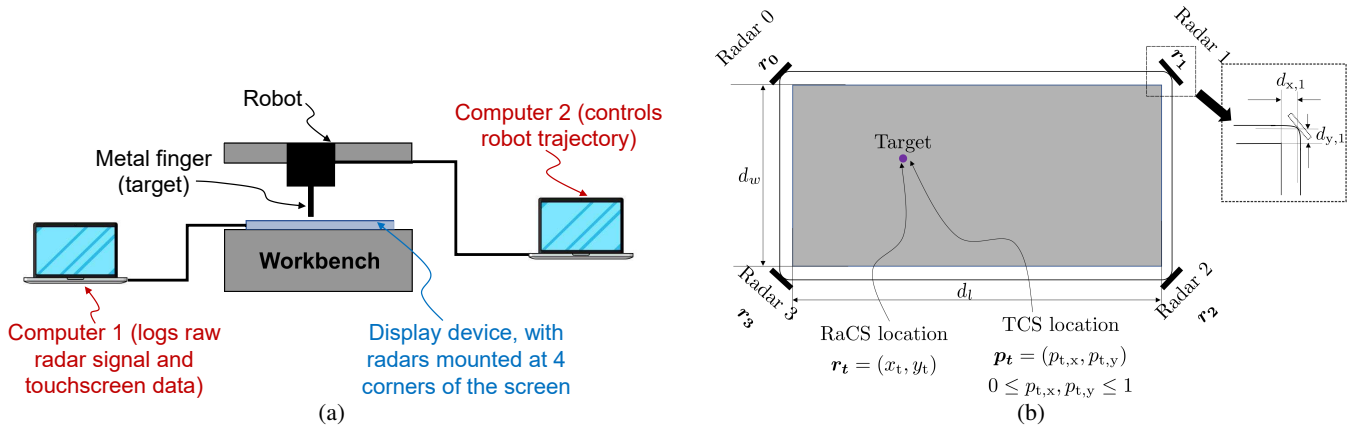


Fig. 1: Schematic of (a) the robot-based data collection setup, and (b) mmWave radar network-based positioning setup, and description of the **R**adar **C**oordinate **S**ystem (RaCS) and **T**ouchscreen **C**oordinate **S**ystem (TCS).

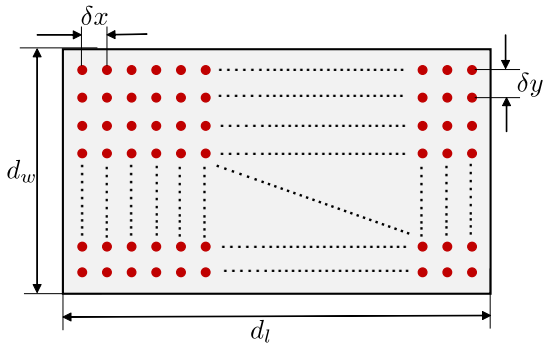


Fig. 2: Schematic of the touch locations on the display. In our setup, $d_l = 34.3$ cm, $d_w = 17.8$ cm, and $\delta x = \delta y = 1$ cm.

the mainlobe width. It is worthwhile to note that the IF signal contains only the in-phase component and hence, is a real-valued signal. Thus, the range-FFT is symmetric about $r = N_{\text{IF}}N_{\text{os}}/2$. Clutter removal is then used to eliminate the scattered returns from static objects using an IIR moving target indication (MTI) filter to get the post-MTI range-FFT signal $x_{\text{MTI}}(f, r, c, j, i)$, given by

$$\begin{aligned} x_c(f, r, c, j, i) &= \beta x_r(f, r, c, j, i) + (1 - \beta)x_c(f, r, c, j, i), \\ x_{\text{MTI}}(f, r, c, j, i) &= x_r(f, r, c, j, i) - x_c(f, r, c, j, i), \end{aligned} \quad (4)$$

where $0 < \beta < 1$ is the IIR filter response parameter and $x_c(f, \dots)$ is the clutter estimate for the f^{th} frame. Finally, to keep the feature dimension manageable for real-time implementation, averaging across chirps and *boresight beamforming* are performed to get the beamformed range-FFT feature $x_{\text{bf}}(f, r, i)$ using²

$$x_{\text{bf}}(f, r, i) = \frac{1}{N_{\text{rx}}N_{\text{ch}}} \sum_{j=0}^{N_{\text{rx}}-1} \sum_{c=0}^{N_{\text{ch}}-1} x_{\text{MTI}}(f, r, c, j, i). \quad (5)$$

Note that for uniform linear/planar arrays, boresight beamforming is equivalent to signal averaging across RX antennas.

²Since the feature is obtained through linear operations on the raw IF signal, averaging across chirps and boresight beamforming can equivalently be performed on the IF signal prior to range-fft as well.

D. Ground Truth

For each touch event, the capacitive touchscreen-based ground truth (GT) information is composed of the relative location $(p_{t,x}, p_{t,y})$ and touch timestamp (t_{GT}) , such that $0 \leq p_{t,x} \leq 1$ and $0 \leq p_{t,y} \leq 1$. The relative locations are converted to locations in the radar coordinate system $\mathbf{r}_{\text{GT}} = (x_{t,\text{GT}}, y_{t,\text{GT}})$ using knowledge of the reference radar location w.r.t. the touch area. We use the sign convention shown in Fig 1b, using which the GT coordinates are calculated using

$$\mathbf{r}_{t,\text{GT}} = (r_{t,x}d_l + d_{x,0}, -r_{t,y}d_w - d_{y,0}). \quad (6)$$

III. CONVENTIONAL SIGNAL PROCESSING (CSP)-BASED POSITIONING

To improve the accuracy of the conventional signal processing-based estimates, we use range estimates from the beamformed signal, as well as the per-RX signals. Firstly, the post-MTI range-FFT from each RX antenna is averaged across chirps using

$$x_{\text{MTI},c}(f, r, j, i) = \frac{1}{N_{\text{ch}}} \sum_{c=0}^{N_{\text{ch}}-1} x_{\text{MTI}}(f, r, c, j, i). \quad (7)$$

Then, range estimates corresponding to the per-RX ($\hat{r}_{ij}(f)$) as well as beamformed signals ($\hat{r}_{\text{bf},i}(f)$) are estimated using

$$\hat{r}_{ij}(f) = \Delta r_{\text{os}} \cdot \arg \max_r |x_{\text{MTI},c}(f, r, j, i)|^2, \quad (8)$$

$$\hat{r}_{\text{bf},i}(f) = \Delta r_{\text{os}} \cdot \arg \max_r |x_{\text{bf}}(f, r, i)|^2. \quad (9)$$

To have reliable ranging performance in the presence of low SNR conditions and strong clutter regions (e.g. portion of the hand excluding the finger such as shoulder, torso, palm etc.), we invalidate the range estimate when there is no consensus among the different per-RX range estimates. The range estimate from the i^{th} sensor is computed using

$$\hat{r}_i(f) = \begin{cases} \hat{r}_{\text{bf},i}(f) - r_{\text{cal},i} & \text{if } |\hat{r}_{ij}(f) - \hat{r}_{ik}(f)| \leq \Delta r_{\text{th}} \\ \text{nan} & \text{otherwise.} \end{cases}, \quad (10)$$

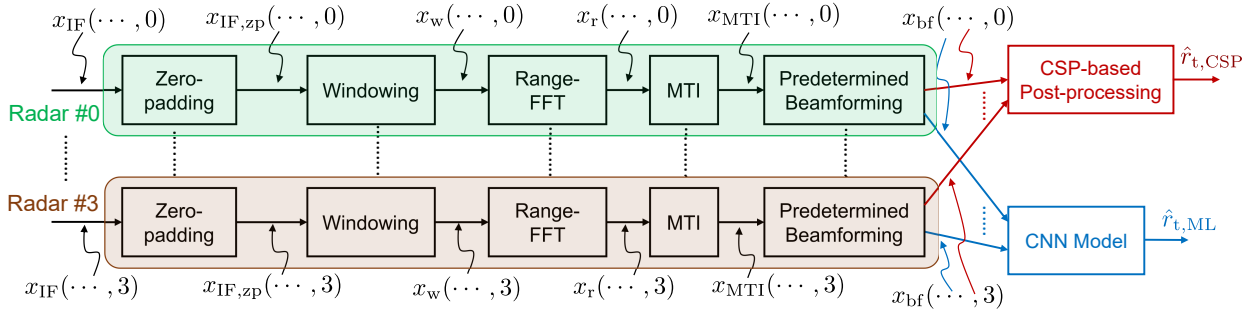


Fig. 3: Flowchart of the radar signal processing pipeline implemented on each radar sensor.

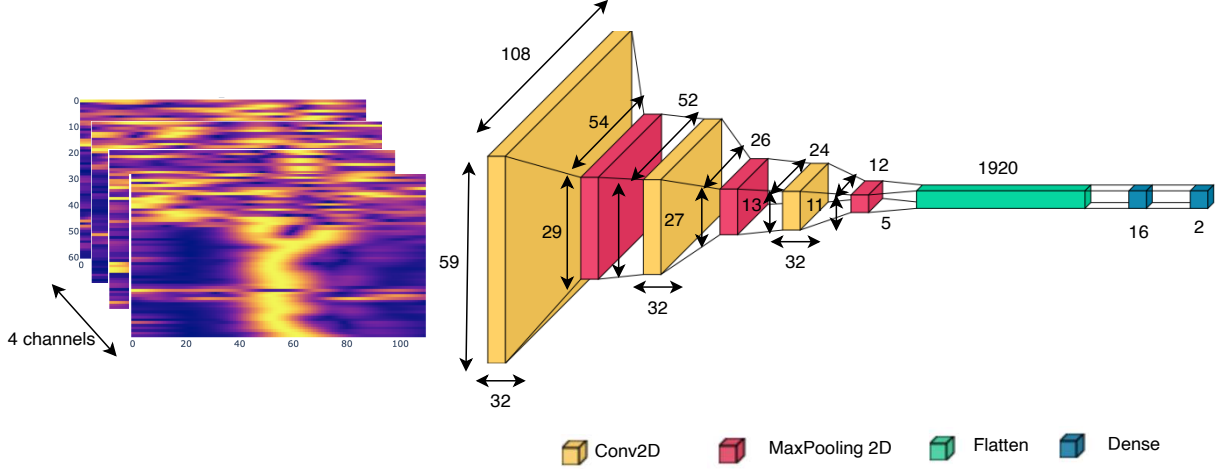


Fig. 4: The four-channel input ($\bar{x}_{bf}(\mathcal{F}_i(n), \mathcal{R})$) consists of heatmaps from the four mmWave radars. The Deep Convolutional Neural Network consisting of three Convolution2D layers and Maxpooling2D and one Dense layer. The final dense layer outputs the 2D position estimate ($\hat{r}_{t,ML}$) of the touch event on the display.

where Δr_{th} is the range consensus tolerance, and $r_{cal,i}$ is the range offset for the i^{th} sensor, which is obtained using a one-time calibration of the localization environment. Let f_n be the radar frame index corresponding to the n^{th} touch event. Then, the range estimates are averaged over a window of $N_w = 5$ frames, to mitigate the unavailability of range estimates, resulting in a range estimate $\bar{r}_i(f_n) = \frac{1}{N_{val}} \sum_{m=0}^{N_w-1} \hat{r}_i(f_n - m) \mathbb{1}[\hat{r}_i(f_n - m) \neq \text{nan}]$, where $N_{val} = |\{m | \hat{r}_i(f_n - m) \neq \text{nan}\}|$.

Finally, the target's position estimate ($\hat{r}_{t,CSP}(f)$) is obtained by solving the nonlinear least squares (NLS) problem [11]

$$\hat{r}_{t,CSP}(f_n) = \arg \min_{\mathbf{r}} \sum_{i=0}^{N_s-1} (\|\mathbf{r}_i - \mathbf{r}\|_2 - \bar{r}_i(f_n))^2. \quad (11)$$

IV. DEEP NEURAL NETWORK-BASED POSITIONING

A. Feature Generation

The datastream from each radar sensor and the capacitive touchscreen are collected independently without explicit synchronization. The relatively high sampling rate of the radar (120 Hz) and the touchscreen (90 Hz) w.r.t. the finger movement speed eliminates the need for explicit sensor

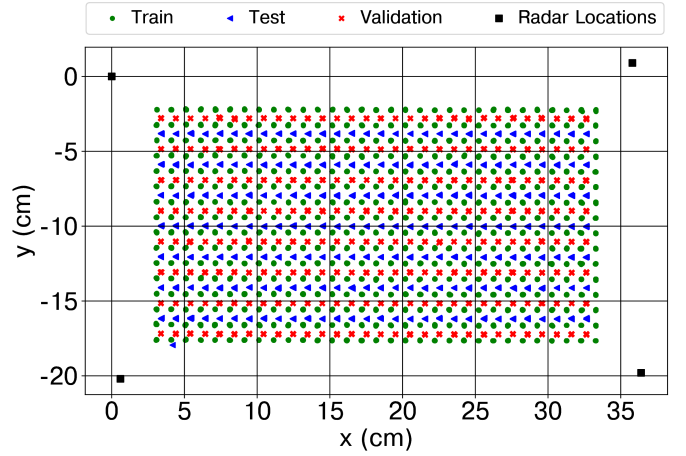


Fig. 5: Partitioning of touch points into train/validation/test datasets, and their locations relative to the radar sensors. In our testbed, the radar locations are $\mathbf{r}_0 = (0, 0)$ cm, $\mathbf{r}_1 = (35.8, 0.9)$ cm, $\mathbf{r}_2 = (36.4, -19.8)$ cm, and $\mathbf{r}_3 = (0.6, -20.2)$ cm.

synchronization. The radar frame indices corresponding to each touch event are found using the GT touch time t_{GT} . Suppose the i^{th} radar frame index corresponding to the touch

event is $f_{GT,i}$. For each touch event, the feature contains the beamformed range-FFT for all radar sensors for the frames $\mathcal{F}_i = \{f_{GT,i-25}, f_{GT,i-24}, \dots, f_{GT,i+25}\}$ for the range bin indices $\mathcal{R} = \{0, 1, \dots, R_{\max}\}$, where $R_{\max} = \left\lceil \frac{N_{\text{os}} \sqrt{d_l^2 + d_w^2}}{\Delta r} \right\rceil$ is the maximum possible target distance on the display. In our setup, $R_{\max} = 110$. The feature for the n^{th} touch event is the tensor $\bar{\mathbf{x}}_{bf}(\mathcal{F}_i(n), \mathcal{R}) \in \mathbb{C}^{61 \times 110 \times 4}$, comprising of four heatmaps, each corresponding to one of the radar sensors.

B. Machine Learning Model Architecture

The proposed CNN-based architecture used in this work is shown in Fig. 4, which is similar to the LeNet-5 architecture [12] except it uses 3 layers of convolution with 3 max-pooling layers rather than 3 layers of convolution with 2 average pooling in the original architecture. We use three cascaded 2D-Convolutional + 2D-Maxpooling layers with filter sizes of 32, 32, and 32 respectively, each with a kernel size of 3×3 with ReLu activation. The maxpooling layer after convolution is used not only to reduce the feature size but also reduce over-fitting. This improves the generalization and also reduces the memory requirements to host the model on the device. After the convolution and pooling operations, the output is flattened, followed by a dense layer with 16 units and then the final 2 dense unit to generate the 2D position coordinates $\hat{\mathbf{r}}_{t,ML} = (\hat{x}_{t,ML}, \hat{y}_{t,ML})$.

V. EXPERIMENTAL RESULTS

A. Data Collection

The effective touch area on the 15.6 inch display is a rectangular area of length $d_l = 34.3$ cm and width $d_w = 17.8$ cm. In a single session, the data is collected across a point grid with an arbitrary offset from Radar 0 (\mathbf{r}_0), such that consecutive touch points are separated by 1 cm along both the axes, as shown in Fig. 2. Data is collected in two stages:

- 1) *Training Dataset*: In this stage, we collected data for 50 sessions. In each session, the robot touched the screen across a 31×16 grid. After accumulating data across all sessions, the training data has the dimension (24799, 61, 110, 4).
- 2) *Validation and Test Datasets*: In this stage, we collected data for 15 sessions. In each session, the robot touched the screen across a 30×15 grid. The grid pattern was designed such that the validation/test touchpoints have a position offset of $\delta \mathbf{r}' = (0.5 \text{ cm}, 0.5 \text{ cm})$ relative to the training touchpoints. This offset is introduced to test the generalization performance of the machine learning model to unseen data. Finally, touchpoints in the odd/even rows are allocated to the validation/test datasets, as shown in Fig. 5. Hence, the validation and dataset dimensions are (3600, 61, 110, 4) and (3150, 61, 110, 4) respectively.

B. Range Calibration

For the conventional signal processing approach, the range calibration offset for each radar sensor ($r_{\text{cal},i}$) is estimated using the training dataset. For the n^{th} touchpoint in the

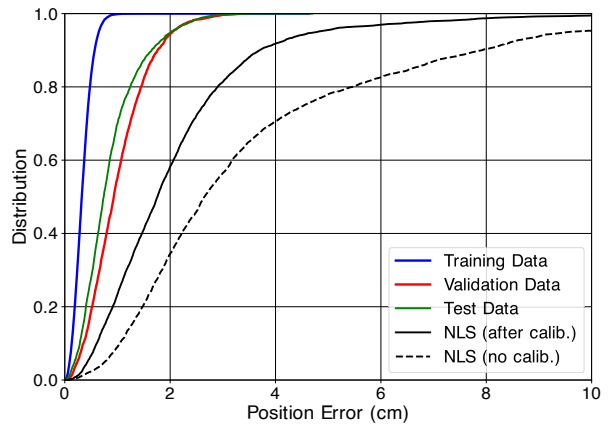


Fig. 6: Distribution of position error for the Conventional Signal Processing and ML-based approaches.

TABLE II: Position Error Performance Comparison on the Test Dataset

| Performance Metric | CNN-based | CSP-based |
|---------------------------|-----------|-----------|
| Median Pointwise RMSE | 0.84 cm | 1.92 cm |
| 90%ile Pointwise RMSE | 1.6 cm | 3.85 cm |
| Median Error (all points) | 0.82 cm | 1.75 cm |
| 90%ile Error (all points) | 1.62 cm | 3.7 cm |

training set, the range estimate corresponding to each radar sensor is computed using (10), and is compared to the GT distance $r_{t,GT,i}(n) = \|\mathbf{r}_{t,GT}(n) - \mathbf{r}_i\|_2$. Suppose the corresponding range error is $\Delta r_i(n) = r_{t,GT,i}(n) - \bar{r}_i(f_n)$. Then, the range calibration offset is estimated by computing the empirical average of the range errors for each sensor, i.e. $r_{\text{cal},i} = \frac{1}{N_{\text{train}}} \sum_{n=0}^{N_{\text{train}}-1} \Delta r_i(n)$, where $N_{\text{train}} = 24799$. Note that this is the least-squares (LS) solution of the range-bias estimation problem under the model $\hat{r}(n) = r(n) + \epsilon$ for $n = 0, 1, \dots, N_{\text{train}} - 1$, where ϵ is the range bias.

C. Position Error Performance Comparison

Fig. 6 shows the marginal distribution of position error (marginalized across the entire test dataset) for (a) training/validation/test datasets (CNN-based approach), and (b) the test dataset (CSP-based approach). First, we observe that there is more than a $2\times$ improvement in the validation/test position error performance in the median and 90th percentile value, when using our CNN-based approach. Furthermore, we observe that these values are well within the range resolution of the radar waveform ($\Delta r_{\text{res}} = 3.075$ cm). On the other hand, we observe that while range calibration significantly improves the position error performance, the range-calibrated CSP-based approach achieves a 90th percentile position error of 3.7 cm, which is $\sim 20\%$ higher than Δr_{res} . In alignment with our understanding of the physical limitations imposed by the radar waveform, the CSP-based method is unable to achieve sub-resolution accuracy. The key performance statistics are summarized in Table II.

Fig. 7 shows the heatmaps of RMSE position error for the CNN and CSP-based methods, for different touch regions on the display. These heatmaps are computed for the test

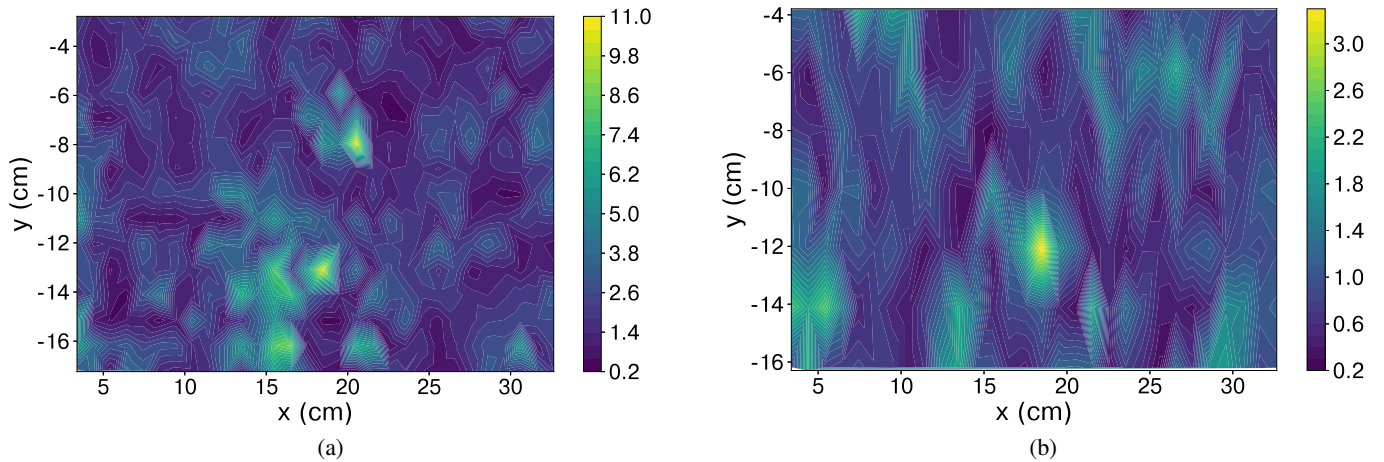


Fig. 7: Comparison of the root mean square (RMSE) position error heatmap (colorbar unit is cm) when using (a) conventional signal processing (after range calibration), and (b) our proposed CNN Model (from Fig. 4).

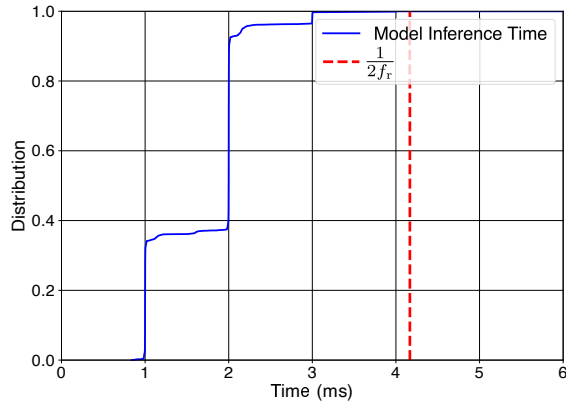


Fig. 8: Distribution of the CNN inference time ($t_{\text{CNN,inf}}$, shown in blue) versus the half radar frame repetition interval ($1/2f_r$, shown in red).

dataset, visualized in Fig. 5. We observe that the CNN-based approach results in more than a $3\times$ improvement in the worst-case (maximum) RMSE, compared to that of the CSP-based approach. In general, there is more than a $2\times$ improvement in RMSE position error for the CNN-based approach relative to that of CSP, when comparing the point-wise median and 90th percentile RMSE, as shown in Table II.

D. Feasibility of Real-Time Implementation

Inference execution time and model size are two important aspects of the CNN model that determine whether it is suitable for real-time implementation. Our model has $\sim 9 \times 10^4$ parameters, with a total size of ~ 350 KB. Thus, the memory required to store the model is quite small, and can be fitted on any standard system-on-chip (SoC).

For integrating any ML-based algorithm into a real-time localization system, it is important that the inference time ($t_{\text{CNN,inf}}$) be smaller than the radar frame repetition interval ($1/f_r$). To evaluate feasibility, we used a computer with an Intel i7-1185G7 processor, 16 GB RAM, and no GPU.

Fig. 8 shows the distribution of $t_{\text{CNN,inf}}$ characterized on the test dataset. We observe that the median as well as the 90th percentile inference time is 2 ms, which is considerably smaller than the radar frame interval (8.33 ms, in our system). Thus, the small model size and inference time indicates that our proposed CNN-based approach is well-suited for real-time implementation on general purpose processor-based platforms.

VI. CONCLUSION

In this paper, we proposed a mmWave FMCW radar-based touch localization system, wherein a deep neural network was trained to accurately localize a robot-mounted metal finger. We demonstrated that the CNN-based approach achieved sub-resolution position error, and significantly outperformed conventional signal processing-based algorithms. Finally, we discussed the feasibility of implementing our proposed approach in real-time. The small (a) CNN model size, and (b) inference time on general purpose computing platforms (relative to the radar frame interval), point towards a very strong feasibility for implementation on a real-time localization system.

In this work, we have focused on accurate localization of robot-mounted targets. In general, extending this work to design localization systems for small targets such as accurate touch localization of human finger, and enabling handwriting on non touchscreen displays are worthwhile to enable low-cost technologies for human-screen interaction.

REFERENCES

- [1] C.-L. Lin, C.-S. Li, Y.-M. Chang, T.-C. Lin, J.-F. Chen, and U.-C. Lin, "Pressure Sensitive Stylus and Algorithm for Touchscreen Panel," *Journal of Display Technology*, vol. 9, no. 1, pp. 17–23, 2013.
- [2] S. Yun, Y.-C. Chen, H. Zheng, L. Qiu, and W. Mao, "Strata: Fine-Grained Acoustic-based Device-Free Tracking," in *Proceedings of the 15th annual international conference on mobile systems, applications, and services*, 2017, pp. 15–28.
- [3] D. Wu, R. Gao, Y. Zeng, J. Liu, L. Wang, T. Gu, and D. Zhang, "FingerDraw: Sub-Wavelength Level Finger Motion Tracking with WiFi Signals," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 4, no. 1, pp. 1–27, 2020.

- [4] 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.
- [5] L. Shangquan and K. Jamieson, "Leveraging Electromagnetic Polarization in a Two-Antenna Whiteboard in the Air," in *Proceedings of the 12th International on Conference on emerging Networking EXperiments and Technologies*, 2016, pp. 443–456.
- [6] T. Wei and X. Zhang, "mTrack: High-Precision Passive Tracking using Millimeter Wave Radios," in *Proceedings of the 21st Annual International Conference on Mobile Computing and Networking*, 2015, pp. 117–129.
- [7] N. Xiao, P. Yang, X.-Y. Li, Y. Zhang, Y. Yan, and H. Zhou, "MilliBack: Real-Time Plug-n-Play Millimeter Level Tracking using Wireless Backscattering," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 3, no. 3, pp. 1–23, 2019.
- [8] Y. Cao, A. Dhekne, and M. Ammar, "ITrackU: Tracking a Pen-like Instrument via UWB-IMU Fusion," in *Proceedings of the 19th Annual International Conference on Mobile Systems, Applications, and Services*, 2021, pp. 453–466.
- [9] N. Hendy, H. M. Fayek, and A. Al-Hourani, "Deep Learning Approaches for Air-Writing Using Single UWB Radar," *IEEE Sensors Journal*, vol. 22, no. 12, pp. 11 989–12 001, 2022.
- [10] F. Khan, S. K. Leem, and S. H. Cho, "In-Air Continuous Writing Using UWB Impulse Radar Sensors," *IEEE Access*, vol. 8, pp. 99 302–99 311, 2020.
- [11] R. Zekavat and R. M. Buehrer, *Handbook of Position Location: Theory, Practice and Advances*. Wiley-IEEE Press, 2019.
- [12] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based Learning Applied to Document Recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [13] Z. Shen, S. Li, X. Zhao, and J. Zou, "CT-Auth: Capacitive Touchscreen-Based Continuous Authentication on Smartphones," *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–16, 2023.
- [14] J. Yan, C. C. J. M. Tiberius, G. J. M. Janssen, P. J. G. Teunissen, and G. Bellusci, "Review of Range-based Positioning Algorithms," *IEEE Aerospace and Electronic Systems Magazine*, vol. 28, no. 8, pp. 2–27, 2013.
- [15] R. M. Rao, A. V. Padaki, B. L. Ng, Y. Yang, M.-S. Kang, and V. Marojevic, "ToA-Based Localization of Far-Away Targets: Equi-DOP Surfaces, Asymptotic Bounds, and Dimension Adaptation," *IEEE Transactions on Vehicular Technology*, vol. 70, no. 10, pp. 11 089–11 094, 2021.
- [16] R. M. Rao and D.-R. Emenonye, "Iterative RNDOP-Optimal Anchor Placement for Beyond Convex Hull ToA-Based Localization: Performance Bounds and Heuristic Algorithms," *IEEE Transactions on Vehicular Technology*, vol. 73, no. 5, pp. 7287–7303, 2024.
- [17] L. Badriasl and K. Dogancay, "Three-Dimensional Target Motion Analysis using Azimuth/Elevation Angles," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 50, no. 4, pp. 3178–3194, 2014.
- [18] A. Santra and S. Hazra, *Deep Learning Applications of Short-Range Radars*. Artech House, 2020.
- [19] F. Uysal, "Phase-Coded FMCW Automotive Radar: System Design and Interference Mitigation," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 1, pp. 270–281, 2020.
- [20] FCC, "FCC Empowers Short-Range Radars in the 60 GHz Band," *Federal Communications Commission, Final Rule*, July 2023. [Online]. Available: <https://www.govinfo.gov/content/pkg/FR-2023-07-24/pdf/2023-15367.pdf>