Edinburgh Research Explorer

# Real-Time NLOS/LOS Identification for Smartphone-Based Indoor Positioning Systems Using WiFi RTT and RSS
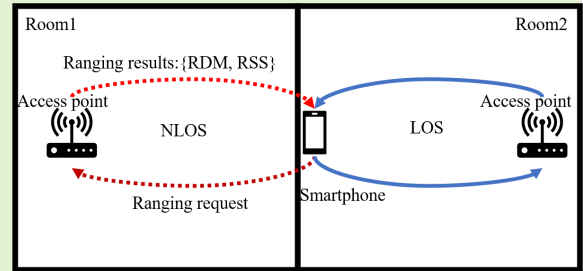
OPEN ACCESS

# Real-time NLOS/LOS Identification for Smartphone-based Indoor Positioning Systems using WiFi RTT and RSS

Yinhuan Dong, Tughrul Arslan, *Senior Member, IEEE*, Yunjie Yang, *Member, IEEE*

**Abstract**—**The accuracy of smartphone-based positioning systems using WiFi usually suffers from ranging errors caused by non-line-of-sight (NLOS) conditions. Previous research usually exploits several distribution features from a long time series (hundreds of samples) of WiFi received signal strength (RSS) or WiFi round-trip time (RTT) to achieve a high identification accuracy. However, the long time series or large sample size attributes to high power and time consumption in data collection for both training and testing. This will also undoubtedly be detrimental to user experience as the waiting time for getting enough samples is quite long. Therefore, this paper proposes three new real-time NLOS/LOS identification methods for smartphone-based indoor positioning systems using WiFi RSS and RTT distance measurement (RDM). Based on our extensive analysis of RSS and RDM dispersion features, three machine learning algorithms were chosen and developed to separate the samples for NLOS/LOS conditions. Experiments show that our best method achieves a discrimination accuracy of over 96% with a sample size of 10. Considering the theoretically shortest WiFi ranging interval of 100ms of the RTT-enabled smartphones, our algorithm is able to provide the shortest latency of 1s to get the testing result among all of the state-of-art methods.**

**Index Terms**—**Real-time NLOS identification, WiFi RTT, WiFi RSS, machine learning, smartphone, positioning.**



## I. INTRODUCTION

LOCATION is vital for numerous applications driven by uncountable mobile users and developers. The global navigation satellite system (GNSS) has been served for years to provide high-precision localization and relevant applications in outdoor scenarios [1]. However, the low penetration of GNSS signal through walls and obstacles sharply decreases the positioning accuracy in the indoor environment [2]. Numerous indoor positioning methods have been proposed to fill the vacancy of providing location-based services (LBS) in indoor scenario these years, such as WiFi [3]–[5], ultra-wideband (UWB) [6]–[8], Radio Frequency Identification Device (RFID) [9]–[11]and, Bluetooth [12]–[14]. Since most smartphones are WiFi-enabled, and WiFi access points (APs) are widely installed in both private and public environments, WiFi-based methods are widely used to provide positioning service to users with smartphones in indoor scenarios.

WiFi-based indoor positioning methods are usually implemented by either the fingerprinting method or range-based method. Fingerprinting method computes the user's position by matching the received signal strength (RSS) from multiple WiFi access points (AP) that are near to the RSS that

The authors are with the School of Engineering, University of Edinburgh, Edinburgh, EH8 9YL, UK (e-mail: yinhuan.dong@ed.ac.uk; tughrul.arslan@ed.ac.uk; y.yang@ed.ac.uk).

is pre-recorded at known locations. The range-based method usually computes the user location by a certain algorithms, such as multilateration and the least square, through the estimated distance between the AP and the smartphone according to the RSS [15]. Especially, the protocol of fine time measurement (FTM) standardized by IEEE 802.11-2016 brought the new technology of round-trip time (RTT), which could provide meter-level positioning accuracy [16]. Promoted by Google, various manufacturers claim that their updated Android-powered smartphones are WiFi-RTT enabled, this includes Google, Xiaomi, LG, Samsung, Sharp, and so on [17]. These WiFi-RTT enabled smartphones can send WiFi ranging requests to nearby APs to get the ranging results (such as RSS and RTT-based distance measurement) in a short period of time ($\geq 200$ ms in this study) without connecting to the APs.

However, both fingerprinting and range-based localization methods are not satisfactory as the WiFi signal is vulnerable to multi-path effects, especially in non-line-of-sight (NLOS) conditions when obstacles block the clear line-of-sight between the transmitter (AP) and the receiver (smartphone). Therefore, NLOS conditions should always be identified first. The work from Xiao *et al.* [15] extracts the multiple features from a group of RSS samples to distinguish between LOS/NLOS conditions. The algorithm could achieve the accuracy of around 95% and over 90% using the testing set

(collected in the same experiment environment as the training set) and validation set (collected in a different environment) in the group size of 1000, respectively. Yu *et al.* [18] proposed a method to reduce the impact of NLOS and multi-path through the combination of real-time ranging model based on WiFi RTT and pedestrian dead reckoning (PDR). Genter *et al.* presented a distance estimation error model with the Gaussian mixture model to calibrate the measuring distance using WiFi RTT. The work from Han *et al.* [19] uses support vector machine (SVM) to classify the NLOS and LOS conditions with the features extracted from a group of WiFi RSS and RTT samples. The accuracy of such method achieves over 92% on the testing set (collected one the same site as the training set) while the group size is 99.

Although the above-mentioned state-of-art methods could achieve good performance under certain circumstances, most of them are not real-time enough or not for real-time positioning. Given a fixed sampling rate, it cost much time and power consuming for smartphone to collect hundreds of samples to identify the NLOS and LOS conditions between the phone and the AP. This is even worse in practice as the ranging-based method usually needs at least three APs to calculate the user's location. Some researchers have worked to reduce the influence of NLOS conditions in tracking the user's position by integrating the PDR method. However, it is not suitable when computing the absolute position. Even though some work could achieve good performance on the testing set collected on the same site as the training set; the robustness has not been validated.

This paper proposes and compares three real-time NLOS/LOS identification algorithms for smartphone-based indoor positioning systems using ranging results of WiFi RSS and RTT-based distance measurement (RDM). By analyzing a series of temporal WiFi samples (ranging samples), we extract several features of the ranging data and employ three machine learning algorithms to distinguish between NLOS and LOS conditions. The main contributions of this work are as follows:

- We propose real-time NLOS and LOS identification methods with the highest identification accuracy but the lowest latency using WiFi RSS and RDM.
- Our study is the first investigation of exploiting the dispersion features of a short series (small number) of WiFi ranging samples for real-time NLOS/LOS use. Three machine learning-based classification techniques are chosen to explore the various features simultaneously.
- Rather than using K-fold cross-validation, we stipulate two testing and validating strategies to evaluate the robustness of the proposed algorithms in different environments.
- Our experiments are based on the data collected by commercial smartphones and WiFi access points in real experiment sites without any pre-setting or reconfiguring infrastructure.

## II. PROBLEM ANALYSIS AND MOTIVATION

### A. Ranging with WiFi RTT

As previously mentioned, the protocol of FTM is able to calculate the distance between RTT-enabled smartphones and
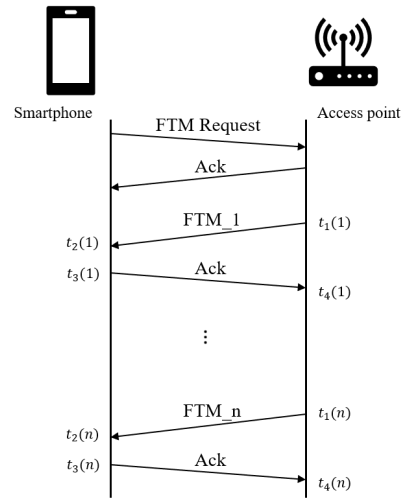


Fig. 1. Overview of FTM protocol (one FTM request gives one burst with $n$ FTMs ($n \leq 31$).

RTT-enabled access points. As shown in Figure 1, the access point sends acknowledgement (ACK) to the smartphone once the FTM request is received (an Initial FTM Request (iFTMR) should be sent first). This gives a single burst that contains multiple FTMs (maximum of 31, excluding the iFTMR). One burst can happen with a burst period ranging from 100ms to every $2^{16}*100$ms (1.8 hours) , which ultimately depends by the master [20]. As the timestamps record the time when the signal is sent and received, the RTT can be calculated by subtracting the timestamp from the AP and the time delay occurs in the smartphone by:

$$\frac{(\sum_{k=1}^{n} t_4(k) - \sum_{k=1}^{n} t_1(k)) - (\sum_{k=1}^{n} t_3(k) - \sum_{k=1}^{n} t_2(k))}{n} \quad (1)$$

accordingly, the distance between the smartphone and access point (RDM) can be estimated by multiplying half the RTT to the velocity of light ($c = 3*10^8 m/s$):

$$RDM = \frac{1}{2} * RTT * c \quad (2)$$

Promoted by Google, the function of WiFi RTT was introduced in Android 9 (API level 28) for more practical use. The FTM request is named as a ranging request in Android system for the RTT-enabled smartphones. The successful ranging request gives the user multiple measurements, such as RTT-based distance (in mm), RSS (in dBm), timestamp (in ms), and so forth. We define a pair of RTT-based measured distance (RDM) and RSS from one successful ranging request as one ranging sample in this paper.

### B. Problem Statement

Considering sending requests and receiving results to an AP in an indoor scenario from a smartphone, the signal path between the smartphone and the AP is usually blocked by some obstacles, defined as NLOS conditions. As illustrated in Figure 2, the signals between the smartphone and the APs in the same room are in LOS condition as there is no interference on the paths from any obstacles. However, the signals between
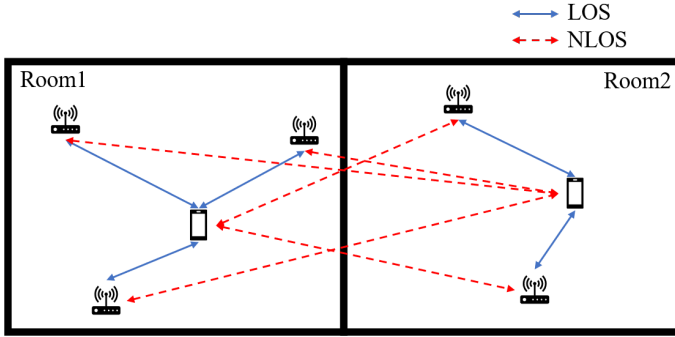
Fig. 2. Illustration of LOS and NLOS cases.

the phone and the APs in different rooms are blocked by a wall, which makes such signals in NLOS conditions.

As previously discussed in Section I, RSS is detrimentally affected by the multi-path effect, especially in the NLOS environment. Similarly, the NLOS condition usually causes an inaccurate distance measurement between the smartphone and AP due to a delay and fluctuation of the signal's travel time. We can observe from the ranging results in Figure 3 that the NLOS causes significant errors in measuring the RSS and RDM, which will degrade the positioning accuracy. Therefore, identifying NLOS signals enable developers and researchers to eliminate such signals or reduce their detrimental effects on further signal processing.

*C. Motivation*

Some previously proposed methods have already extracted and analyzed some features of RSS values from a set of $N$ RSS samples $[RSS_1, RSS_2, ..., RSS_N]$ to identify NLOS signals. For instance, Xiao *et al.* [15] verified that the combination of some features of RSS (such as mean and kurtosis) could achieve a very high identification accuracy while the set size of $N$ is 1000. However, this takes about 100s (considering the shortest sampling interval of 100ms) for an Android-powered smartphone to collect enough samples in practice. As we consider the real-time NLOS identification for smartphones, a small set size should always be considered in both training and testing of the algorithm. Nevertheless, owing to fewer samples used in real-time NLOS/LOS identification, we can infer that most of the distribution features of RSS that were verified to be effective to NLOS/LOS identification in previous studies cannot help distinguish NLOS from LOS for real-time use. Figure 4 illustrates the probability distributions of skewness and kurtosis of RSS of a set of 10 samples. We can observe that the fitted curves of NLOS and LOS samples show almost the same shape and a large portion of overlapping areas, which makes the skewness and kurtosis of RSS fail to identify different conditions.

Different from previous studies, we exploit the dispersion of the RDM and RSS samples in this study. We notice that the ranging samples (especially for RDM) in Figure 3 show high dispersion in different conditions. For example, we calculate the quartile deviation and the number of quartiles of RDM samples with the set size of 10. The distribution of the results

is shown in Figure 5. We can observe that there is usually a shift between the two fitted curves of the samples collected in different conditions, which makes the dispersion features are able to help discriminate NLOS and LOS conditions.

Although our statistical analysis shows that some features may or may not help identify NLOS and LOS conditions, there are still some uncertainties about how the combination of such features could improve the identification accuracy and which combination shows the best performance. As we notice in Figure 3 that RSS usually shows higher sparsity but lower dispersion than RDM, it is speculated that some extracted features of RSS may not be able to help to identify the NLOS and LOS conditions. The features with a low contribution to the improvement of identification accuracy should be eliminated, as more features will increase the complexity of the algorithm. Therefore, we propose to use machine learning techniques detailed in the next section to find the best combination of features to distinguish between NLOS and LOS conditions by evaluating the multiple features simultaneously. All the features and their symbols are listed in Table I.

## III. NLOS/LOS IDENTIFICATION EMPLOYING MACHINE LEARNING

The task here is to decide whether a given set of ranging samples corresponds to NLOS or LOS conditions. Machine learning-based algorithms have been widely adopted by many studies to classify NLOS and LOS samples. Least squares-support vector machine (LS-SVM) is one of the most popular algorithms that has been implemented in various studies proposed by Xiao *et al.* [15], Chitambira *et al.* [32], Han *et al.* [21], and so on. Alternatively, random forest shows great performance in classification problems with low computation complexity [26]. In recent years, deep learning-based NLOS/LOS identification methods (such as [24], [25]) have also attracted attention. Therefore, we design and implement these three algorithms to solve the NLOS/LOS discrimination problem.

*1) Random forest:* Random forest is an ensemble learning algorithm that trains the model using several classifiers (decision trees) with random sets of features [26]. It makes the final prediction by combining the results from all the classifiers through majority voting. As RF uses the combination of both
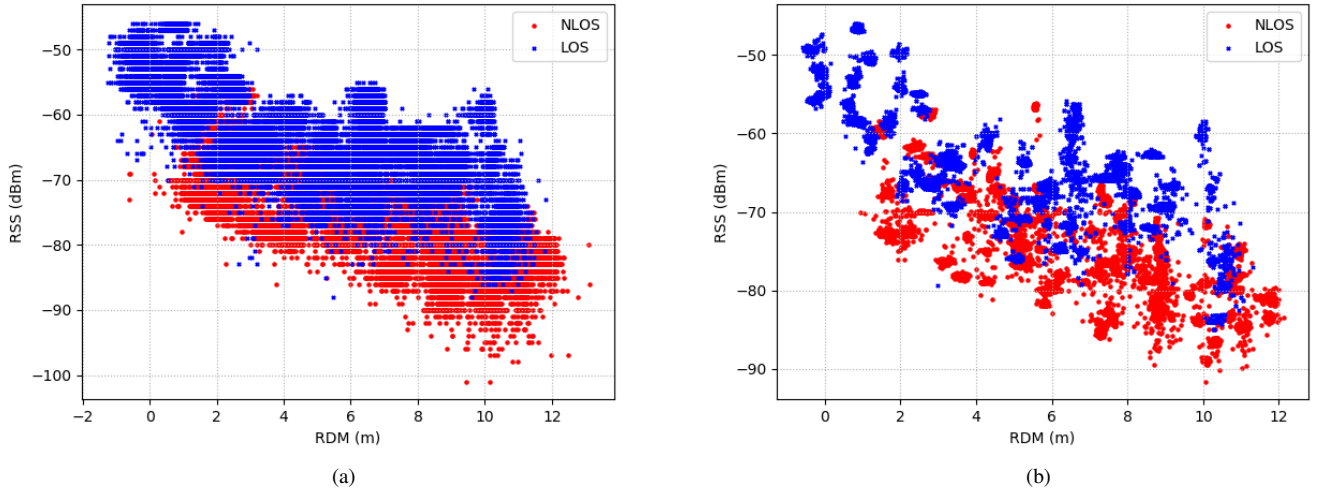
Fig. 3. A large amount of ranging samples collected in a variety of settings, including different smartphones, sampling rate and ground-truth distance (details will be illustrated in Section IV), in NLOS/LOS conditions (a) original ranging samples; (b) grouped ranging samples (group size = 10).
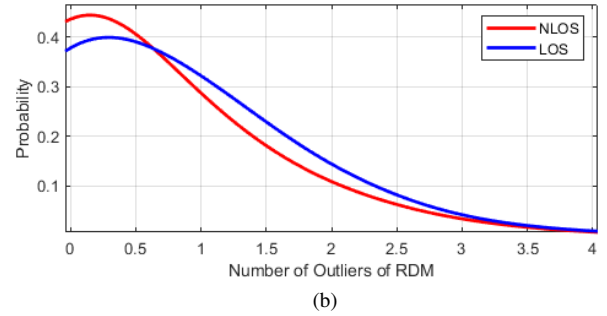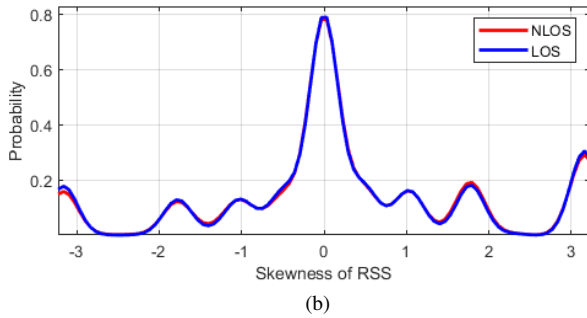


Fig. 4. Probability distribution of the extracted features of RSS (a) Kurtosis of RSS; (b) Skewness of RSS.



Fig. 5. Probability distribution of the extracted features of RDM (a) Quartile Deviation of RDM; (d) Number of Quartiles of RDM.

boosting and bagging, it usually produces a model that is not highly overfitting with high efficiency. As we only have the two labels of NLOS and LOS condition in this study, Classification and Regression Tree (CART) is chosen to solve this binary classification problem (we employ 10 decision trees in this study to avoid high computation complexity [27]). Rather than using information entropy, Gini index is used to evaluate the features and divide the input samples in CART

for faster computation. The Gini index is defined as:

$$\text{Gini}(x) = \sum_{l=1}^{L} p_l (1 - p_l) = 1 - \sum_{l=1}^{L} p_l^2 \tag{3}$$

where $L$ is the number of categories of the dataset, and $p_l$ denotes the probability of the sample's label is $l$. As the dataset $X$ has 2 classes of data in this case, $L$ is set to 2, and hence the Gini index of $X$ according to a given feature $x_i$ could be

computed by:

$$\text{Gini}(X, x_i) = \frac{X_1}{X} \text{Gini}(X_1) + \frac{X_2}{X} \text{Gini}(X_2) \quad (4)$$

The Gini index reflects the uncertainty of the given set of samples. Since the Gini index is the difference between 1 and the sum of the probability squares of category $l$ (as shown in Equation 3 ), the larger Gini index, the higher uncertainty of the samples. Therefore, the optimal partition feature $x_\star$ could be selected by minimizing the Gini index as follows:

$$x_\star = \arg\min \text{Gini}(X, x_i) \quad (5)$$

*2) Least Square Support Vector Machine:* Owing to the easier training process and higher generalization quality of SVM [28], we propose using SVM to identify the NLOS and LOS conditions. Specially, the least square SVM (LS-SVM) [29] is used in this work to avoid the quadratic programming problem of SVM [15].

Given a training set of $N$ samples $\{x, l\}^N$, each sample is composed by a feature matrix $x$ (contains several features mentioned in Section II) and a label $l(x) \in \{0, 1\}$ (in this case, $l(x) = 1$ stands for NLOS condition and $l(x) = 0$ represents LOS condition). The linear classification problem of the training samples could be expressed as the following function:

$$l(x) = sign[w^T \varphi(x) + w_0] \quad (6)$$

where $sign$ is the signum function, $\varphi(x)$ is the eigenvector of the feature matrix $x$, $w$ and $w_0$ are the weighting parameters learned from the training set. As the NLOS and LOS conditions are not linearly separable, we use a Gaussian radial basis function (RBF) [30] to realize better classification performance:

$$k(x, x_i) = \varphi(x)^T \cdot \varphi(x_i) = \exp\left[-\frac{\|x - x_i\|_2^2}{2\sigma^2}\right] \quad (7)$$

where the hyper-parameter $\sigma^2$ is learned from the training data by solving the following optimization problem:

$$\arg\min_{w, w_0, e, \sigma^2} \frac{1}{2}\|w\|^2 + c\frac{1}{2}\sum_{i=1}^{N} e_i^2 \quad (8)$$
$$\text{s.t. } l_i\left[w^T \varphi(x_i) + w_0\right] = 1 - e_i, \forall i$$

where $e$ is the penalty of misclassification, and $c$ is the weighting factor that controls the trade-off between training error and model complexity. This optimization problem is a linear programming problem [29], which could be solved by its Lagrangian dual and Karush-Kuhn-Tucker conditions [31]. Once the parameters of the classifier have been well-trained, the prediction could be computed by:

$$l(x) = \text{sign}\left[\sum_{i=1}^{N} \lambda_i l_i k(x, x_i) + w_0\right] \quad (9)$$
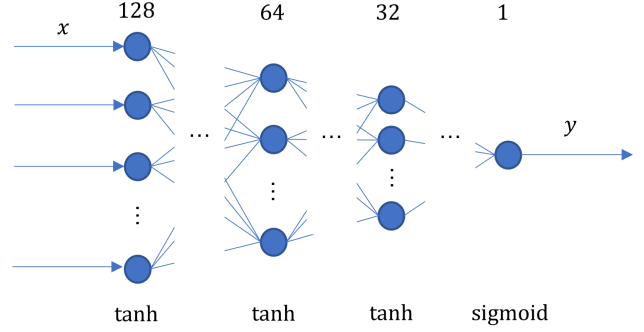
where $\lambda_i$ is the Lagrange multiplier.



Fig. 6. Illustration of the network structure.

TABLE II
DETAILS OF THE NEURAL NETWORK

| | |
|---|---|
| Hidden layers | 128-64-32-1 |
| Activation function | tanh-tanh-tanh-sigmoid |
| Batch size | 32 |
| Learning rate | 0.001 |
| Epoch† | 100 |
| Early stopping patience | 10 |
| Optimizer | Adam |
| Loss function | binary cross-entropy |

† This is a pre-set epoch number. As early stopping strategy is applied, the network usually stops training earlier.

*3) Deep Neural Network:* A deep neural network model usually consists of multiple stacked deep hidden layers, which tries to learn potential features using weights and biases in each layer. Due to its remarkable learning ability and flexible network structure, DNN has been widely used to solve classification and regression problems.

In this paper, we design a 4-layer fully connected DNN structure to solve the NLOS identification problem. We use the extracted features $x$ as the input data. As the structure of the network is shown in Figure 6, from left to right, each layer contains 128, 64, 32 and 1 hidden nodes, respectively. We use tanh as the activation function of the first three layers, whereas sigmoid is used in the last layer to ensure the output $y$ is in the range from 0 to 1. The output $y$ is then mapped into the output label by:

$$l(x) = \begin{cases} 1, & y > 0.5 \\ 0, & else \end{cases} \quad (10)$$

where $l(x) = 1$ stands for NLOS condition and $l(x) = 0$ represents LOS condition. The details of the network structure and training settings are listed in Table II.

## IV. EXPERIMENT SETTINGS

This section illustrates the experiment sites, equipment and methodology that we used to collect the ranging samples in this study.

### A. Experiment Sites

The experiments were conducted in two different real-world sites, including an office and a student accommodation.
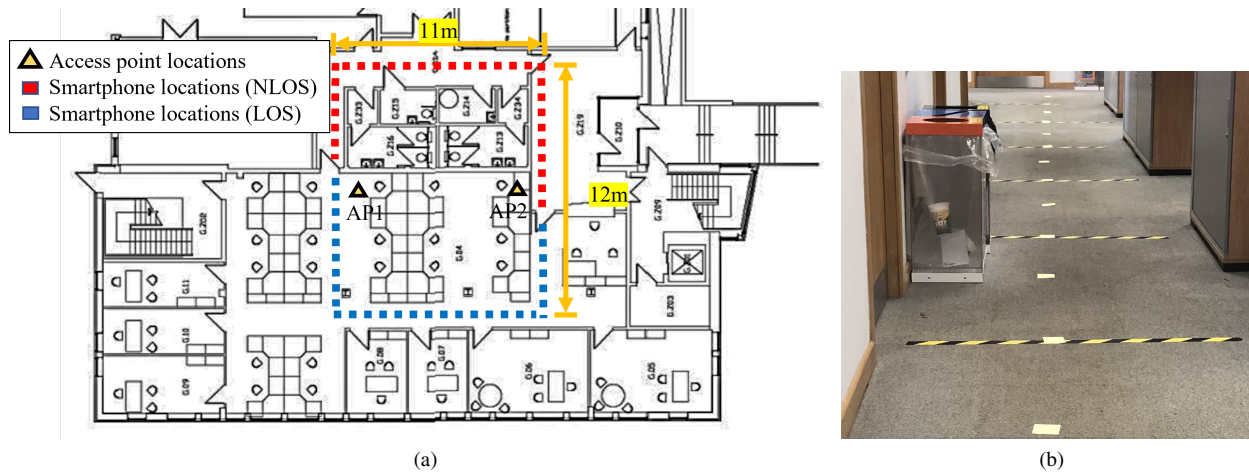
(a)        (b)

Fig. 7. Floor plan of the first floor of the Scottish Microelectronics Centre. The data collected here is for training and testing the performance of the proposed algorithms.
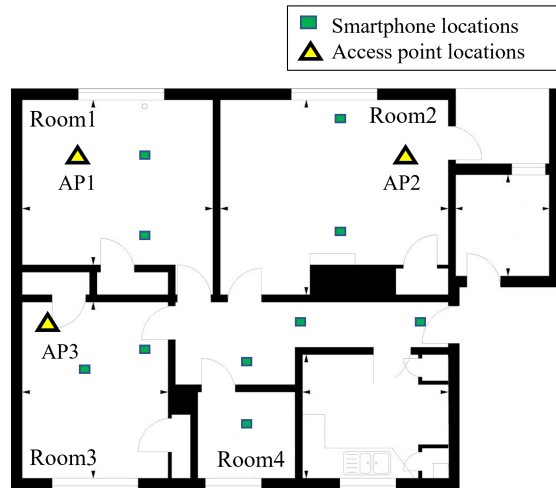


Fig. 8. Floor plan of the student accommodation. The data collected here is for validating the generalization ability of the proposed algorithms.



Fig. 9. Smartphones and access point used in this experiment

Figure 7a shows the office site on the ground floor of the Scottish Microelectronics Centre. It is a complex indoor environment with wooden doors and concrete constructed walls (reinforced with metal rebars), as well as different obstacles. The volunteers were asked to collect the ranging samples from the RTT-enabled access point following the path. The paths are composed of multiple test points at the interval of 1m as shown in Figure 7b. The NLOS samples were collected from AP2 on the path marked in red, while the LOS samples were collected from AP1 following the path marked in blue. The two paths are separated by the wall and other obstacles in the environment. With these settings of test points and access point locations, the ground-truth distances between the smartphone and access point vary from 0.5m to 12m approximately. As the access points were set on two tripods at the height of over 1.7m at the center of the office site, the clear line-of-site was not affected by the desks and chairs surrounding the blue path.

Another experiment site shown in Figure 8 is a student accommodation mainly constructed by concrete and plaster.

The three access points were set on the tripods in different rooms. Data were collected on the test points marked in green. At each point, volunteers were asked to collect samples from all the APs and their conditions. For example, the samples collected in room 1 from AP1 are in LOS condition, and the samples collected from AP2 and AP3 are in NLOS condition.

### B. Devices and Software

In this experiment, four smartphones include three Google Pixel 2 and one Google Pixel 2XL were used to collect the ranging samples. Google WiFi access points were utilized as the transmitter in the measurements. Some core specifications are listed in Table III. The devices used in this work are shown in Figure 9. In this study, we used the Android application WifiRttScan developed by Google to send ranging requests and

TABLE IV
NUMBER OF SAMPLES COLLECTED IN DIFFERENT EXPERIMENT SITES

|  | Office site (training&testing) | Student accommodation (validation) |
|---|---|---|
| NLOS | 116550 | 10592 |
| LOS | 116027 | 11893 |
| Sum | 232577 | 22485 |

TABLE V
SUBSETS OF DIFFERENT COMBINATIONS OF FEATURES

| Features | $C_1$ | $C_2$ | $C_3$ | $C_4$ | $C_5$ | $C_6$ | $C_7$ | $C_8$ |
|---|---|---|---|---|---|---|---|---|
| $\mu$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| $\sigma$ |  | ✓ |  | ✓ |  |  |  | ✓ |
| $S$ |  |  |  |  |  | ✓ |  | ✓ |
| $K$ |  |  |  |  |  |  | ✓ | ✓ |
| $Q$ |  |  | ✓ | ✓ |  |  |  | ✓ |
| $\lambda$ |  |  | ✓ | ✓ |  |  |  | ✓ |
| $R$ |  |  |  |  | ✓ |  |  | ✓ |

collect the ranging results. The collected data were processed and analyzed using Python 3.8 on a desktop with an Intel i7-9700 CPU (3.00 GHz) and 32GB installed memory.

### C. Data Collection

As previously mentioned, the lowest latency of sending the ranging request from one RTT-enabled smartphone is 100ms in theory. However, it is recommended by Google that the sampling interval should not be shorter than 200ms to avoid collision and other software problems. Therefore, the lowest latency is set to 200ms in this work, which gives 10 ranging samples every two seconds if all requests are successful (one sample per 200ms). Although there is no clear evidence that illustrates the sampling rate would affect the samples' quality, we set the ranging period at 200ms, 250ms, 333ms, and 500ms for the Google Pixel 2XL and other three Google Pixel 2, respectively, to reduce the uncertainty. The smartphones were always kept waking up (the application was always running in the foreground), face up and oriented parallel to the ground during the data collection to avoid the effects that may caused by some reasons, such as the gesture of holding the phone, orientation of the antenna, or other software problems. Besides, the volunteers were asked to collect the samples statically. This means that the volunteers who held the smartphone can not move once the data collection starts. At least 600 samples were collected by each device at each point.

## V. EVALUATION

This section evaluates the proposed algorithms. We first introduce the construction of the two datasets. This is followed by the training and testing results of the machine learning-based real-time NLOS/LOS identification algorithms. Finally, we select the best features and validate the reliability of the trained models.

### A. Datasets

As listed in Table IV, we construct two datasets from the data collected in different environments. The first dataset contains the samples collected in the environment of the office site that few people were walking around. This set is used for training and testing the machine learning-based NLOS/LOS identification algorithm. The second set of samples collected in the student accommodation is used for validating the generalization ability of our methods in a different experiment site. It is essential that the algorithms could work in different scenarios to avoid repeated training process which is usually labour intensive and time consuming.

### B. Training

We train and test the proposed method trough different combinations of the features. The combinations of features are represented by $C_i$ $(i = 1, 2, \ldots, 8)$. The subsets of different combinations of features are shown in Table V. As the features are extracted from either RDM or RSS signals, we design four different schemes to assemble the features from different signals:

- $C_i^{RDM}(i = 1, 2, \ldots, 8)$: Each subset of different features in this scheme uses only RDM features, including $\mu_{RDM}$, $\sigma_{RDM}$, $Q_{RDM}$, $\lambda_{RDM}$, $R_{RDM}$, $\mathcal{S}_{RDM}$, $\mathcal{K}_{RDM}$.
- $C_i^{RSS}(i = 1, 2, \ldots, 8)$: Each subset of different features in this scheme uses only RDM features, including $\mu_{RSS}$, $\sigma_{RSS}$, $Q_{RSS}$, $\lambda_{RSS}$, $R_{RSS}$, $\mathcal{S}_{RSS}$, $\mathcal{K}_{RSS}$.
- $C_i^{FTM}(i = 1, 2, \ldots, 8)$: This scheme uses all the features from RDM and RSS (also called FTM features).
- $C_i^{SEL}(i = 1, 2, \ldots, 8)$: The selected features are used in this scheme (also called SEL features). This contains all the features from RDM samples and only the mean value of RSS: $\mu_{RDM}$, $\mu_{RSS}$,$\sigma_{RDM}$, $Q_{RDM}$, $\lambda_{RDM}$, $R_{RDM}$, $\mathcal{S}_{RDM}$, $\mathcal{K}_{RDM}$.

### C. Testing

We evaluate the proposed algorithms using three metrics, i.e. fail rate (the algorithm fail in detecting NLOS), false alarm rate (the algorithm identifies NLOS while the samples are from LOS), and overall false detection rate (the sum of fail rate and false alarm rate). They are denoted by $P_N$, $P_L$ and $P_O$, respectively.

*1) Using a single source of the signal:* We first evaluate the test results using the features of a single source of the signal (either RSS features or RDM features). The NLOS identification errors are shown in Figure 10. We can observe from Figure 10a that for each machine learning-based NLOS identification algorithm, the extracted features from the RSS signal are not able to reduce the NLOS identification error significantly. Compared to the benchmark false detection rates of 0.2399, 0.1614, and 0.1568 using DNN, RF, and SVM with only RSS mean values ($C_1^{RSS}$), there is a little improvement when the extracted features are used. The lowest errors of different algorithms are 0.1808 ($C_5^{RSS}$), 0.1604 ($C_3^{RSS}$) and 0.1568($C_1^{RSS}$) , respectively.

On the contrary, we can observe from Figure 10b that some of the extracted RDM features can help to improve the identification rate significantly. Compared to the benchmark overall false detection rates of using DNN, RF, and SVM with
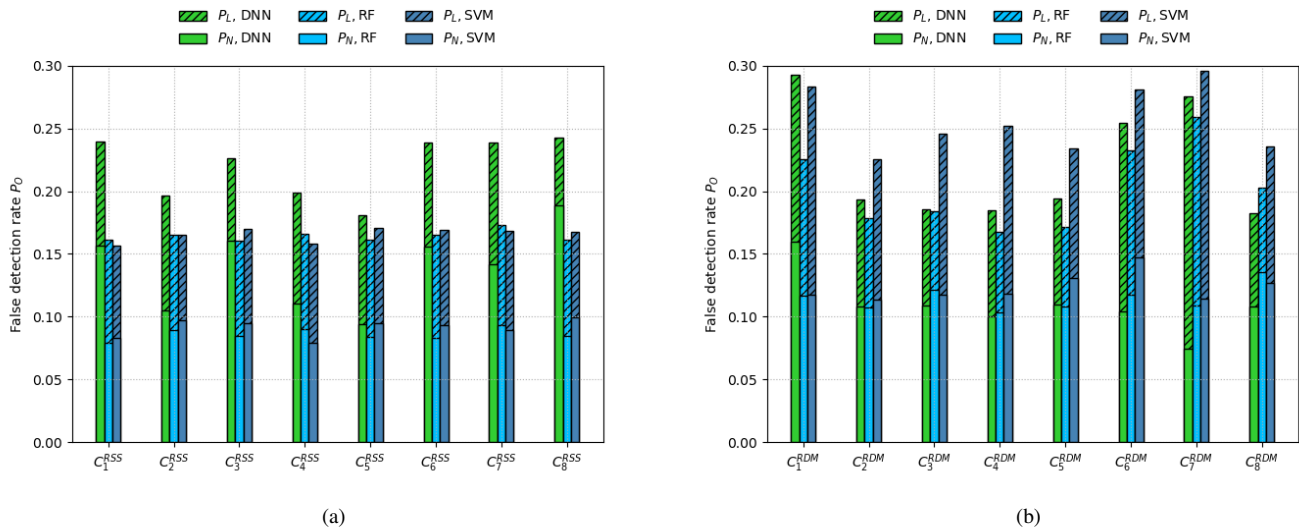
Fig. 10. Testing results of fail rate ($P_N$), false alarm rate ($P_L$) and false detection rate ($P_O$) of the machine learning-based NLOS identification algorithms using different features (a) RSS features; (b) RDM features.

only RDM mean values ($C_1^{RDM}$), the subsets of $C_3^{RDM}$ and $C_4^{RDM}$ can usually provide lower errors. For example, the lowest error of using DNN with features in $C_4^{RDM}$ is 0.1847, which is 37% lower than the benchmark error of 0.2928. This illustrates that the dispersion features of standard deviation and quartile deviation of the RDM signal effectively improve identification accuracy.

Although there is some improvement in the detection rate with some extracted features, the NLOS identification error is still high when a single source of the signal is used. Therefore, we propose to use both signals and their features to provide a better identification rate. The evaluation is as follows.

*2) Using hybrid sources of signals:* We evaluate the proposed machine learning-based NLOS identification algorithms with extracted features from RSS and RDM signals. We can observe from Figure 11a that the detection errors are sharply reduced when both RDM and RSS features are involved in detecting the NLOS and LOS cases. The lowest errors of using DNN, RF and SVM are 0.05956 ($C_3^{FTM}$), 0.0304 ($C_1^{FTM}$) and 0.0492 ($C_1^{FTM}$), which are at least 70% lower than using either RDM or RSS features. Nevertheless, in exchange for the extremely low identification errors, the computation complexity is higher. As we use both RSS and RDM features, each subset of $C_i^{FTM}$ contains doubled features than either $C_i^{RSS}$ or $C_i^{RDM}$. To maintain the low detection error and reduce the computation complexity (number of features), we evaluate the machine learning algorithms on some selected features next.

The analysis of our previous experiments on NLOS identification using a single source of the signal has shown that the RSS features cannot help reduce the error significantly. In contrast, some of the RDM features (such as standard deviation and quartile deviation) could improve the detection accuracy. Hence, rather than using all features from both signals, we keep only RDM features and the mean value of RSS in each subset of $C_i^{SEL}$. As the results are shown in Figure 11b, the detection errors of three machine learning-based algorithms

using SEL features maintain at a similar level as FTM features in general; however, the number of features nearly halved. We can also observe from the figure that the standard deviation and quartile deviation of RDM in subset $C_2^{SEL}$, $C_3^{SEL}$ and $C_4^{SEL}$ can usually provide better results than other features.

**Summary**: Our testing results reveal that the RSS features provide a slight improvement in identifying NLOS. In contrast, some RDM features, such as standard deviation and quartile deviation, can help to improve the detecting accuracy significantly. However, the identification accuracy of using either RSS or RDM features is still not high enough (about 80%, approximately). Hence, we applied the machine learning algorithms to the FTM features (the joint of RSS and RDM features). The detection accuracy can be improved significantly to higher than 90% when RSS and RDM features are used. Nevertheless, more features may provide better detection accuracy but also contribute to higher computation complexity. Only useful features should be kept. As our previous experiments have shown that most of RSS features cannot help improve the detection accuracy significantly, we designed SEL features by selecting only the mean value of RSS signals and the RDM features. The testing results demonstrate that the SEL features can well maintain the detection accuracy at a similar level as the FTM features but halve the number of features. In the next subsection, we will illustrate how we validate the algorithms in a different site (than from where we collected the training samples) to verify the robustness against environmental change.

### D. Validation

It is essential that the designed algorithms are able to work in different sites without repeated training to avoid the time-consuming site survey. In order to select the best algorithm, we validate the accuracy of the algorithms in a different site of student accommodation (Figure 8) than from where we collected the training samples (Figure 7a).
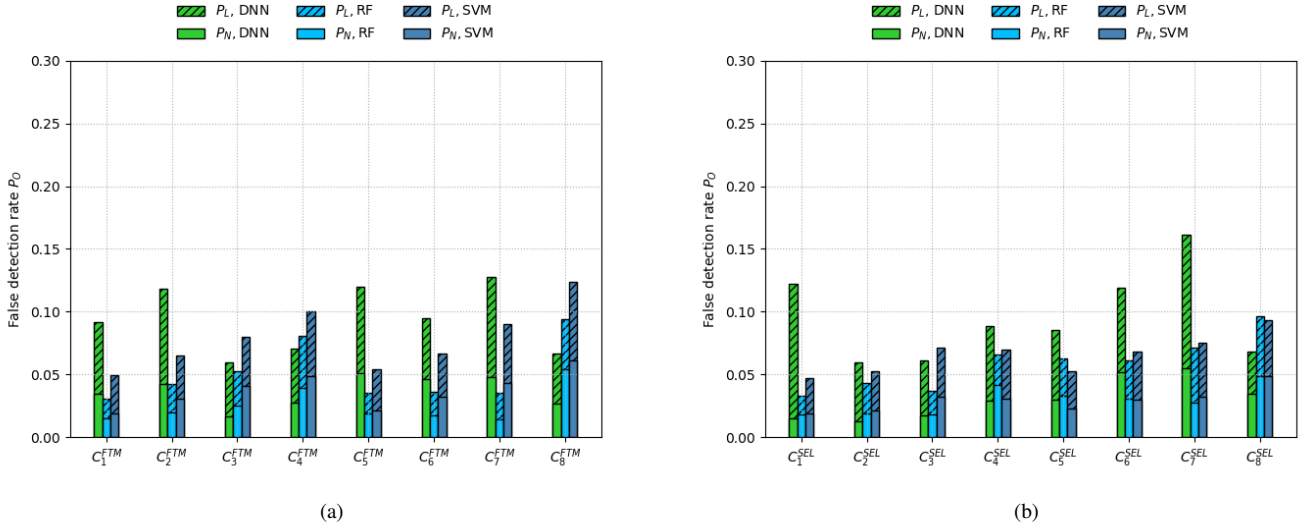
Fig. 11. Testing results of fail rate ($P_N$), false alarm rate ($P_L$) and false detection rate ($P_O$) of the machine learning-based NLOS identification algorithms using different features (a) FTM (RDM and RSS) features ; (b) SEL features (selected features of RDM and RSS).
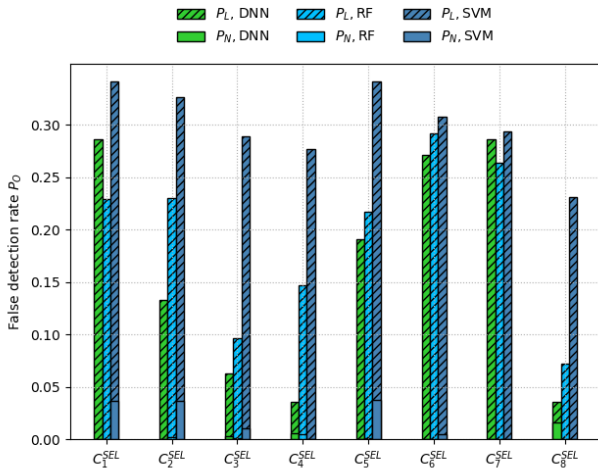


Fig. 12. Validating results of fail rate ($P_N$) and false alarm rate ($P_L$) and false detection rate ($P_O$) of the machine learning-based NLOS identification algorithms using different features (a) RDM and RSS; (b) RDM and RSS (selected features).

As the results are shown in Figure 12, for each machine learning-based NLOS identification algorithm, SEL features of $C_7^{SEL}$ can always provide the lowest detection error. This illustrates that the combination of all extracted features can effectively help to improve identification accuracy. As the results of the top three feature subsets with the lowest error are listed in Table VI, besides using all features of $C_7^{SEL}$, the three machine learning algorithms can still achieve good results when $C_3^{SEL}$ and $C_4^{SEL}$ features are used. Although there is a slight increase in the error when some features are eliminated, $C_3^{SEL}$ and $C_4^{SEL}$ features can maintain the detection error at a similar level of using $C_7^{SEL}$ but almost halve the number of features. This again proves that the standard deviation and quartile deviation of the RTT-based distance measurements can help to detect the NLOS signals.

We can also observe from the results that the DNN-based

NLOS identification algorithm provides 96.42% accuracy, which outperforms the others in terms of the algorithms trained by all selected features. The RF-based algorithm can also provide good accuracy of 92.76%. However, SVM shows the worst accuracy of 76.90%. When some of the features are eliminated, and only the dispersion features are kept ($C_4^{SEL}$ and $C_5^{SEL}$), although there is a slight decrease in the detection accuracy, DNN-based solution can still maintain the best performance. This result indicates that the DNN-based algorithm is more robust with the selected dispersion features than RF and SVM in the NLOS identification problem.

In addition, the training time of each model is listed in Table VI. There is a trade-off between the detection accuracy and the training time using DNN. Although the DNN solution can always provide the best accuracy, the higher computation complexity leads to a much longer training time than others. In contrast, the RF solution shows extremely low training costs with moderate detection accuracy.

### E. Summary

In summary, the testing and validation results can be concluded as follows:

- Using a single signal source signal of either RSS or RDM cannot provide good performance in real-time NLOS identification for smartphone-based indoor positioning systems. While the hybrid of the two signals and their features can sharply reduce the identification error.
- The extracted RDM features of RTT signals are more helpful to improve the NLOS identification accuracy than RSS features. Specifically, the dispersion features, such as standard deviation and quartile deviation of the RTT-based distance measurements, can improve the identification accuracy whether which other features are added or eliminated. This study also follows our previous analysis in Section II that RSS usually shows higher sparsity but lower dispersion than RDM, which makes

TABLE VI

COMPARISON OF DIFFERENT MACHINE LEARNING-BASED ALGORITHMS USING THE TOP 3 FEATURE SUBSETS

| Algorithm | Top 3 feature subsets (s) | Number of features | Correct detection rate $P_C$ (1-$P_O$) | Training time (s) |
|---|---|---|---|---|
| DNN | $C_7^{SEL}$, $C_4^{SEL}$, $C_3^{SEL}$ | 8, 5, 4 | 96.42%, 96.42%, 93.73% | 52.33, 48.05, 44.45 |
| Random Forest | $C_7^{SEL}$, $C_3^{SEL}$, $C_4^{SEL}$ | 8, 4, 5 | 92.76%, 90.33%, 85.34% | 0.90, 0.96, 0.93 |
| LS-SVM | $C_7^{SEL}$, $C_4^{SEL}$, $C_3^{SEL}$ | 8, 4, 5 | 76.90%, 72.36%, 71.12% | 0.91, 0.93, 0.93 |

the RDM features more suitable to identify NLOS and LOS cases.

- The DNN-based NLOS identification algorithm can provide the best detection accuracy in both training and validating data with the selected features, which shows good generalization ability to the environment change. However, the DNN solution also provides the highest computation complexity, which takes the longest time to train the model.
- Compared to DNN, RF solution can provide moderate detection accuracy and robustness but significantly training costs. A good balance between the detection accuracy and the training cost is essential, as repeated training may be needed in some complex indoor environments with a large volume of data, such as airports and mega shopping malls.

## VI. CONCLUSION AND FUTURE WORK

This paper proposed three real-time NLOS/LOS identification algorithms for smartphone-based indoor positioning systems using WiFi ranging. Previous research shows that some distribution features (such as kurtosis and skewness of RSS) from a long series of samples can effectively identify NLOS/LOS for non-real-time use. However, according to our analysis of multiple extracted features from a large amount of ranging samples, we infer that such features cannot help improve the identification accuracy when the sample size shrinks for real-time use. Rather than using distribution features in the literature, we explored the effect of dispersion features of ranging samples on real-time identification. Three machine learning algorithms have been adopted to investigate the impact of different feature combinations on real-time NLOS/LOS discrimination accuracy. The testing results illustrate that the dispersion features of RDM combined with the mean value of RSS could provide the best real-time discrimination accuracy. In addition, the validation results show that the proposed DNN-based algorithm has the highest generalization ability to environmental change but also the highest training complexity (longest training time). In comparison, the Random Forest-based solution can provide moderate accuracy with a much lower training complexity (shorter training time).

As this study focuses on real-time NLOS/LOS identification for positioning use, the algorithms are designed to distinguish the ranging samples collected in different conditions while the user is stationary. Future work will investigate the NLOS/LOS identification in tracking and navigating (considering the user motion) in indoor scenarios.

## REFERENCES

[1] B. Hofmann-Wellenhof, H. Lichtenegger, and J. Collins, *Global Positioning System: Theory and Practice*, vol. 1. Wien, Austria: Springer-Verlag, 1993.

[2] H. M. Hussien, Y. N. Shiferaw, and N. B. Teshale, "Survey on Indoor Positioning Techniques and Systems," in *International Conference on Information and Communication Technology for Development for Africa*, 2017, pp. 46-55: Springer.

[3] P. Bahl and V. N. Padmanabhan, "RADAR: An in-building RF-based user location and tracking system," in *Proc. 19th Annu. Joint Conf. IEEE Comput. Commun. Soc. (INFOCOM)*, vol. 2, Mar. 2000, pp. 775–784.

[4] M. Youssef and A. Agrawala, "The Horus WLAN location determination system," in *Proc. 3rd Int. Conf. Mobile Syst., Appl., Services*, 2005, pp. 205–218.

[5] W. Sun, M. Xue, H. Yu, H. Tang, and A. Lin, "Augmentation of fingerprints for indoor WiFi localization based on Gaussian process regression," *IEEE Trans. Veh. Technol.*, vol. 67, no. 11, pp. 10896–10905, Nov. 2018.

[6] B. Alavi and K. Pahlavan, B, "Studying the effect of bandwidth on performance of UWB positioning systems," in *Proc. IEEE Wireless Commun. Netw. Conf.*, Las Vegas, NV, Apr. 2006, vol. 2, pp. 884–889.

[7] B. Hanssens, D. Plets, E. Tanghe, C. Oestges, D. P. Gaillot, M. Liénard, T. Li, H. Steendam, L. Martens, and W. Joseph, "An indoor variancebased localization technique utilizing the UWB estimation of geometrical propagation parameters," *IEEE Trans. Antennas Propag.*, vol. 66, no. 5, pp. 2522–2533, May 2018.

[8] K. Yu, K. Wen, Y. Li, S. Zhang, and K. Zhang, "A novel NLOS mitigation algorithm for UWB localization in harsh indoor environments," *IEEE Trans. Veh. Technol.*, vol. 68, no. 1, pp. 686–699, Jan. 2019.

[9] A. R. J. Ruiz, F. S. Granja, J. C. P. Honorato, and J. I. G. Rosas, "Accurate pedestrian indoor navigation by tightly coupling foot-mounted IMU and RFID measurements," *IEEE Trans. Instrum. Meas.*, vol. 61, no. 1, pp. 178–189, Jan. 2012.

[10] C.-H. Huang, L.-H. Lee, C. C. Ho, L.-L. Wu, and Z.-H. Lai, "Real-time RFID indoor positioning system based on Kalman-filter drift removal and Heron-bilateration location estimation," *IEEE Trans. Instrum. Meas.*, vol. 64, no. 3, pp. 728–739, Mar. 2015.

[11] G. Alvarez-Narciandi, J. Laviada, M. R. Pino and F. Las-Heras, "3D location system based on attitude estimation with RFID technology," *2017 IEEE International Conference on RFID Technology & Application (RFID-TA)*, Warsaw, 2017, pp. 80-82, doi: 10.1109/RFID-TA.2017.8098881.

[12] R. Faragher and R. Harle, "Location Fingerprinting With Bluetooth Low Energy Beacons," in *IEEE Journal on Selected Areas in Communications*, vol. 33, no. 11, pp. 2418-2428, Nov. 2015, doi: 10.1109/JSAC.2015.2430281.

[13] N. Yu, X. Zhan, S. Zhao, Y. Wu and R. Feng, "A Precise Dead Reckoning Algorithm Based on Bluetooth and Multiple Sensors," in *IEEE Internet of Things Journal*, vol. 5, no. 1, pp. 336-351, Feb. 2018, doi: 10.1109/JIOT.2017.2784386.

[14] L. Bai, F. Ciravegna, R. Bond and M. Mulvenna, "A Low Cost Indoor Positioning System Using Bluetooth Low Energy," in *IEEE Access*, vol. 8, pp. 136858-136871, 2020, doi: 10.1109/ACCESS.2020.3012342.

[15] Z. Xiao, H. Wen, A. Markham, N. Trigoni, P. Blunsom, and J. Frolik, "Non-line-of-sight identification and mitigation using received signal strength," *IEEE Trans. Wireless Commun.*, vol. 14, no. 3, pp. 1689–1702, Mar. 2015.

[16] "IEEE Standard for Information technology—Telecommunications and information exchange between systems Local and metropolitan area networks—Specific requirements - Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications," in *IEEE Std 802.11-2016 (Revision of IEEE Std 802.11-2012)* , vol., no., pp.1-3534, 14 Dec. 2016, doi: 10.1109/IEEESTD.2016.7786995.

[17] Google. (2020). *Wi-Fi location: ranging with RTT*. Accessed: Jan. 18, 2021.[Online]. Available: https://developer.android.com/guide/topics/connectivity/wifi-rtt.

[18] Yu, Y.; Chen, R.; Chen, L.; Guo, G.; Ye, F.; Liu, Z. "A robust dead reckoning algorithm based on Wi-Fi FTM and multiple sensors". *Remote Sens.* 2019, 11, 504.

[19] C. Gentner, M. Ulmschneider, I. Kuehner, and A. Dammann, "WiFiRTT indoor positioning," in *2020 IEEE/ION Position, Location and Navigation Symposium (PLANS)*, 2020, pp. 1029–1035.

[20] K. Stanton and C. Aldana, "Addition of p802.11-MC fine timing measurement (FTM) to p802.1as-rev: Tradeoffs and proposals," Rep., Mar. 2015.

[21] Han, K.; Yu, S.M.; Kim, S. "Smartphone-based Indoor Localization Using Wi-Fi Fine Timing Measurement." In *Proceedings of the 2019 International Conference on Indoor Positioning and Indoor Navigation (IPIN)*, Pisa, Italy, 30 September–3 October 2019; pp. 1–5.

[22] F. Xiao, Z. Guo, H. Zhu, X. Xie, and R. Wang, "AmpN: Real-time LOS/NLOS identification with WiFi," in *Proc. IEEE Int. Conf. Commun. (ICC)*, Paris, France, May 2017, pp. 1–7.

[23] ] B. Chitambira, S. Armour, S. Wales, and M. Beach, "NLOS identification and mitigation for geolocation using least-squares support vector machines," in *Proc. IEEE Wireless Commun. Netw. Conf.*, 2017, pp. 1–6.

[24] J. Choi, W. Lee, J. Lee, J. Lee, and S. Kim, "Deep learning based NLOS identification with commodity WLAN devices," *IEEE Trans. Veh. Technol.*, vol. 67, no. 4, pp. 3295–3303, Apr. 2018.

[25] ] C. Jiang, J. Shen, S. Chen, Y. Chen, D. Liu, and Y. Bo, "UWB NLOS/LOS classification using deep learning method," *IEEE Commun. Lett.*, vol. 24,no. 10, pp. 2226–2230, Oct. 2020.

[26] L. Breiman, "Random forest," *Mach. Learn.*, vol. 45, no. 1, pp. 5–32, 2001.

[27] T. M. Oshiro, P. S. Perez, and J. A. Baranauskas, "How many trees in a random forest?" in *Proceedings of the 8th International Conference on Machine Learning and Data Mining in Pattern Recognition*, 2012, pp. 154–168.

[28] R. Caruana and A. Niculescu-Mizil, "An empirical comparison of supervised learning algorithms," in *Proc. 23rd ICML*, Pittsburgh, PA, USA, 2006, pp. 161–168.

[29] J. A. K. Suykens and J. Vandewalle, "Least squares support vector machine classifiers," *Neural Process. Lett.*, vol. 9, no. 3, pp. 293–300, Jun. 1999.

[30] M. Buhmann, *Radial Basis Functions: Theory and Implementations*, 1st ed. Cambridge, U.K.: Cambridge Univ. Press, 2003.

[31] H. W. Kuhn and A. W. Tucker, "Nonlinear programming," in *Proc. 2nd Berkeley Symp. Math. Statist. Probab.*, Berkeley, CA, USA, 1951, pp. 481–492.

[32] ] B. Chitambira, S. Armour, S. Wales, and M. Beach, "NLOS identification and mitigation for geolocation using least-squares support vector machines," in *Proc. IEEE Wireless Commun. Netw. Conf.*, 2017, pp. 1–6.