

SENTINEL: Securing Indoor Localization Against Adversarial Attacks With Capsule Neural Networks

Danish Gufran¹, *Graduate Student Member, IEEE*, Pooja Anandathirtha, *Student Member, IEEE*,
and Sudeep Pasricha¹, *Fellow, IEEE*

Abstract—With the increasing demand for edge device-powered location-based services in indoor environments, Wi-Fi received signal strength (RSS) fingerprinting has become popular, given the unavailability of GPS indoors. However, achieving robust and efficient indoor localization faces several challenges, due to RSS fluctuations from dynamic changes in indoor environments and heterogeneity of edge devices, leading to diminished localization accuracy. While advances in machine learning (ML) have shown promise in mitigating these phenomena, it remains an open problem. Additionally, emerging threats from adversarial attacks on ML-enhanced indoor localization systems, especially those introduced by malicious or rogue access points (APs), can deceive ML models to further increase localization errors. To address these challenges, we present SENTINEL, a novel embedded ML framework utilizing modified capsule neural networks to bolster the resilience of indoor localization solutions against adversarial attacks, device heterogeneity, and dynamic RSS fluctuations. We also introduce *RSSRogueLoc*, a novel dataset capturing the effects of rogue APs from several real-world indoor environments. Experimental evaluations demonstrate that SENTINEL achieves significant improvements, with up to 3.5× reduction in mean error and 3.4× reduction in worst-case error compared to state-of-the-art frameworks using simulated adversarial attacks. SENTINEL also achieves improvements of up to 2.8× in mean error and 2.7× in worst-case error compared to state-of-the-art frameworks when evaluated with the real-world *RSSRogueLoc* dataset.

Index Terms—Adversarial attacks, adversarial training, capsule neural networks, device heterogeneity, evil twin attacks, man-in-the-middle attacks, rogue access points (APs), Wi-Fi received signal strength (RSS) fingerprinting.

I. INTRODUCTION

IN RECENT years, indoor localization has gained attention for its versatile applications across several industries, such as healthcare, asset tracking, smart homes, location-based advertising, and much more [1]. Technology giants, such as Apple, Google, Meta, and Microsoft, are making substantial investments in indoor localization research to improve the

Manuscript received 6 August 2024; accepted 10 August 2024. This work was supported in part by the National Science Foundation under Grant CNS-2132385. This article was presented at the International Conference on Compilers, Architectures, and Synthesis for Embedded Systems (CASES) 2024 and appeared as part of the ESWEEK-TCAD Special Issue. This article was recommended by Associate Editor S. Dailey. (*Corresponding author: Danish Gufran.*)

The authors are with the Department of Electrical and Computer Engineering, Colorado State University, Fort Collins, CO 80523 USA (e-mail: danish.gufran@colostate.edu; pooja.anandathirtha@colostate.edu; sudeep@colostate.edu).

Digital Object Identifier 10.1109/TCAD.2024.3446717

accuracy and reliability of indoor location-based services [2]. However, achieving high-precision indoor localization remains a formidable challenge due to the inherent complexities and dynamic nature of indoor environments.

Traditional navigation systems, such as the global positioning system (GPS), have found widespread adoption in popular tools, such as Google Maps, Apple Maps, and Waze, mainly owing to their commendable localization accuracies in outdoor settings. However, the dependence of GPS on satellite signals and clear sky visibility poses a significant limitation, rendering this approach ineffective for indoor use [3]. In response to this challenge, researchers have shifted their attention to alternate wireless infrastructures that could be a better fit for localization across indoor spaces, such as Wi-Fi, Bluetooth, and ZigBee. Among these alternatives, Wi-Fi-based localization systems utilizing received signal strength (RSS) have gained significant traction [1], [2], [3], [4]. This surge in popularity for this solution is attributed to the ubiquitous availability of Wi-Fi in indoor spaces and the capability of modern edge devices to capture Wi-Fi RSS, making it a viable option for indoor localization [4].

Wi-Fi RSS is obtained by measuring the signal strength of nearby Wi-Fi routers or access points (APs) via edge devices. This captured RSS data can be used to estimate the current indoor location of an edge device. As the edge device moves, it periodically captures new RSS measurements, reflecting the edge device’s mobility. Leveraging this changing RSS data, many techniques have been proposed for accurate indoor localization, with geometric model-based [5] and fingerprinting model-based [4], [6] approaches emerging prominently. Geometric models utilize propagation methods, such as trilateration [7] and triangulation [8] to pinpoint an edge device’s location. However, these solutions are prone to inaccuracies as they are particularly sensitive to RSS fluctuations caused by dynamic changes and complexities within indoor environments. On the other hand, fingerprinting model-based systems eschew propagation methods by creating a database of Wi-Fi RSS patterns (“fingerprints”) of visible Wi-Fi APs collected throughout the indoor space to estimate location. Fingerprinting models have been shown to exhibit greater resilience to RSS fluctuations, demonstrating higher accuracies than geometric methods [4], [9].

Fingerprinting-based localization solutions comprise two distinct phases: 1) an offline phase and 2) an online phase. During the offline phase, Wi-Fi RSS fingerprints are systematically captured across multiple reference points (RPs)



Fig. 1. Impact of rogue APs on three popular ML-based indoor localization solutions [15], [16], [17] from prior work.

85 within a building floorplan. These fingerprints are then often
 86 utilized to train a machine learning (ML) model, enabling
 87 it to capture underlying patterns and features within the
 88 collected RSS fingerprints [10]. Once trained, this ML model
 89 is deployed on the edge device, making it available in the
 90 online phase for real-time indoor location predictions.

91 In the online phase, the RSS fingerprints may exhibit
 92 fluctuations due to diverse factors in the indoor environments.
 93 These factors include signal attenuation, reflections from
 94 objects, human interference, and multipath fading, which can
 95 introduce fluctuations in the collected RSS fingerprints [11].
 96 Furthermore, edge device heterogeneity exacerbates this issue.
 97 Even among edge devices utilizing the same Wi-Fi chipset
 98 (from the same manufacturer), differences in hardware, soft-
 99 ware, antenna configurations, and firmware settings can
 100 introduce fluctuations in RSS fingerprints [11]. As a result,
 101 training an ML model can be challenging as heterogeneous
 102 and noisy RSS can result in poor generalization and result in
 103 inaccurate location predictions. Priors works have shown up to
 104 a 41% reduction in location accuracy due to these factors [12].
 105 Additionally, the often-overlooked factor of adversarial attacks
 106 can not only perturb the RSS fingerprints (thereby introducing
 107 stronger fluctuations) but also compromise the accuracy and
 108 effectiveness of localization with the edge device, emphasizing
 109 the need for more robust and secure localization systems.

110 Adversarial attacks can mislead popular ML models,
 111 including state-of-the-art deep learning (DL) algorithms that
 112 have been shown to be vulnerable to adversarial examples.
 113 Goodfellow et al. [13] verified the discovery by misleading
 114 the popular GoogLeNet [14] model with adversarial examples.
 115 Similarly, ML-based indoor localization systems also face
 116 the threat of adversarial attacks. The presence of malicious
 117 (or rogue) APs in the building floorplan can be used to
 118 create adversarial attacks by mimicking a legitimate AP and
 119 broadcasting erroneous RSS values. In Fig. 1, we illustrate
 120 the detrimental impact of the presence of rogue APs on
 121 three popular ML-based indoor localization solutions based on
 122 K -nearest neighbors (KNNs) [15], Gaussian process classifier
 123 (GPC) [16], and deep neural networks (DNNs) [17]. This
 124 experiment was conducted on an indoor path in a building
 125 measuring 55 m in length containing 55 RPs (1 RP per meter),
 126 with up to 203 visible APs (per RP). The experiment incor-
 127 porated the popular fast gradient sign method (FGSM) [30]
 128 technique to simulate the presence of rogue APs, resulting in
 129 significantly increased indoor localization errors, with average
 130 error increases of $3.33\times$ for KNN, $3.0\times$ for GPC, and $5.71\times$
 131 for DNN, highlighting the negative impact of the rogue APs
 132 on localization accuracy.

To tackle the challenges posed by RSS fluctuations in 133
 dynamic indoor environments, edge device heterogeneity, and 134
 rogue AP attacks, in this work we introduce SENTINEL, a 135
 novel embedded ML framework that employs modified capsule 136
 neural networks tailored specifically for indoor localization 137
 and rogue AP resilience, offering a more practical, secure, 138
 and real-time solution for indoor localization. The major 139
 contributions of our SENTINEL framework are as follows. 140

- 1) We design a novel modified capsule neural network 141
 specifically for the RSS fluctuation challenges in indoor 142
 localization, tailored to a) overcome the spatial invari- 143
 ance problem in prior DL-based indoor localization 144
 efforts and b) enable lightweight deployment on edge 145
 devices. 146
- 2) We study the effects of rogue AP attacks and propose 147
 an adversarial training setup together with the modified 148
 capsule neural network for resilience against adversarial 149
 (rogue) AP attacks for the first time in indoor localiza- 150
 tion. 151
- 3) We introduce a new Wi-Fi RSS fingerprint dataset called 152
RSSRogueLoc [35] that captures AP attacks from rogue 153
 APs in real-world indoor environments for the first time. 154
- 4) We conduct a performance comparison with SENTINEL 155
 against state-of-the-art indoor localization solutions, to 156
 highlight its effectiveness in the presence of diverse 157
 adversarial attacks, edge device heterogeneity, and RSS 158
 fluctuations across diverse indoor building paths. 159

160 II. RELATED WORK

Wi-Fi fingerprinting-based indoor localization has gained 161
 significant recognition, evident in competitions hosted by 162
 industry giants like Microsoft and NIST [2]. Several classical 163
 ML-based solutions, such as ones based on the KNN [15] 164
 and GPC [16] algorithms, have showcased their poten- 165
 tial in addressing RSS fluctuations arising from dynamic 166
 effects in indoor environments. These fluctuations encom- 167
 pass various factors, including human interference, obstacles, 168
 movement of furniture or equipment, variable population den- 169
 sity, signal interference, reflections by objects, and shadowing 170
 effects [19], [40], [41]. 171

Despite the demonstrated promise of these ML solutions, 172
 they often face challenges in maintaining robustness against 173
 fluctuations introduced by edge device heterogeneity. The 174
 heterogeneity issue arises from differences in Wi-Fi chipsets 175
 and noise filtering software employed by different manu- 176
 facturers of edge devices. As these chipsets and software 177
 stacks are crucial for extracting RSS fingerprints [11], [19], the 178
 heterogeneity within them introduces additional complexities 179
 for traditional ML-based indoor localization systems. 180

In response to these challenges, researchers have explored 181
 the use of more powerful DL algorithms for indoor localiza- 182
 tion, including DNNLOC [17], MLPLOC [18], LC-DNN [19], 183
 CNNLOC [21], SANGRIA [22], ANVIL [23], and TIPS [24]. 184
 DNNLOC [17], MLPLOC [18], and LC-DNN [19] employ 185
 DNNs along with improved RSS preprocessing methods 186
 to enhance feature correlation in the RSS fingerprints. 187
 CNNLOC [21] proposes a modified convolutional neural 188

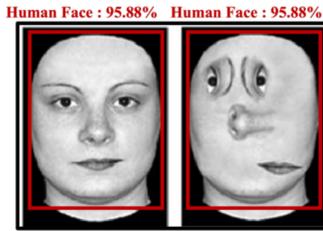


Fig. 2. Spatial invariance problem in DL algorithms. Both cases are classified as valid human faces by a CNN model.

network (CNN), to improve on these efforts by enhancing the model’s ability to capture relevant features in the RSS fingerprints. SANGRIA [22] employs DNN-based autoencoders while ANVIL [23], [42] utilizes attention neural networks, to improve focus on critical input features. TIPS [24] leverages transformer-based encoding of RSS fingerprints for improved resilience against fluctuations introduced by dynamic indoor environments and device heterogeneity. However, these approaches are still significantly impacted by more complex heterogeneity effects in emerging devices and are also susceptible to adversarial attacks, due to the spatial invariance problem in DL algorithms.

Most DL algorithms, particularly CNNs, suffer from the spatial invariance problem where the DL algorithm has a propensity to focus solely on the presence of features in the data while neglecting the precise relative positions of the features [25]. Alterations in the position of each feature can lead to mispredictions by the DL model. This limitation is illustrated in Fig. 2, where the VGGFace algorithm [26], a CNN-based model, struggles to differentiate between the two faces. In the figure on the left, a normal human face is depicted, while the figure on the right presents an abnormal face with jumbled feature positions. The model assigns the same output classification probability to both cases. The concern regarding feature positions is particularly relevant in the context of RSS fingerprints for indoor localization, where positions of certain features represent crucial information and can be specific to a particular RP. When an edge device moves to a different RP, the positions of these features may undergo changes based on the characteristics of the new RP location. Thus, it is imperative to account for the dynamic nature of feature positions when designing practical indoor localization solutions.

To address this limitation and enhance feature extraction, researchers have embraced more recent DL algorithms, such as vision transformers (VITAL) [27], [43] and capsule neural networks (EDGELOC) [28] for indoor localization. VITAL [27], uses vision transformers, introduces positional encoding for each feature, aiming to overcome the spatial invariance limitations posed by CNNs. Similarly, EDGELOC [28] uses a simple capsule neural network derived from [38], treating each captured feature as a vector, considering both magnitude and direction of features. These frameworks show the potential to greatly mitigate the effects of dynamic environments and heterogeneity for indoor localization. However, the introduction of adversarial attacks

especially arising from rogue APs can not only jumble the feature positions but also introduce new malicious features in the data. Such attacks can easily mislead state-of-the-art localization frameworks and compromise user security.

Adversarial training has emerged as a potential solution to address the challenges from adversarial attacks in ML [29]. Popular solutions typically incorporate a subset of adversarial samples along with the training data to allow robustness in the presence of adversarial attacks during inference. Adversarial samples are generated using several popular adversarial methods out of which the FGSM [30] has been widely employed to simulate the effects of adversarial attacks, owing to its simplicity. ADVLOC [31] and CALLOC [32] are two recent solutions that incorporate adversarial training, aiming to address the effects of adversarial attacks in indoor localization. Both ADVLOC [31] and CALLOC [32] integrate FGSM samples during training for adversarial resilience. CALLOC additionally employs curriculum learning along with attention neural networks to enhance feature extraction between the original and adversarial samples, to improve overall robustness. Nevertheless, both solutions fall short of addressing the multitude of challenges associated with dynamic environments, heterogeneity, and adversarial attacks concurrently. Additionally, these solutions heavily rely on simulated data for measuring the efficacy of the model’s performance against adversarial attacks in the online phase. Their performance in real-world adversarial scenarios has not yet been carefully studied.

After carefully studying the simultaneous challenges of dynamic environments, edge device heterogeneity, adversarial attacks, and lack of real-world adversarial attack data to measure the effectiveness of adversarial resilience in indoor localization, in this work we propose SENTINEL, a novel embedded ML framework that goes beyond state-of-the-art DL solutions to better address the spatial invariance problem and improve robustness using an enhanced capsule neural network with techniques that more comprehensively improve resilience to real-world indoor localization challenges. Another important contribution of our work is the design of a newly curated RSS fingerprint dataset called *RSSRogueLoc* [35] that captures the presence of rogue APs within indoor building paths, to analyze the impact of adversarial attacks on indoor localization frameworks in real-world environments, for the first time.

III. ADVERSARIAL ATTACKS IN INDOOR LOCALIZATION

Adversarial attacks involve deliberately perturbing input data to deceive an underlying ML model [30]. This perturbation typically consists of adding noise to individual data values (datapoints) either by introducing new or malicious features (new datapoints) or disrupting the magnitude and positions of features in the input data. These adversarial perturbations exploit limitations in the manner in which features and patterns are learned by the ML model during training, thereby causing mispredictions with the ML model [30].

In the context of indoor localization, Wi-Fi RSS fingerprints are measured in decibels referenced to 1 mW (dBm) and

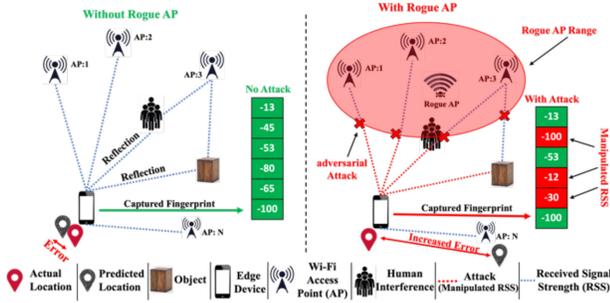


Fig. 3. RSS fluctuations in indoor environment depicting real-world scenarios with and without the presence of rogue APs.

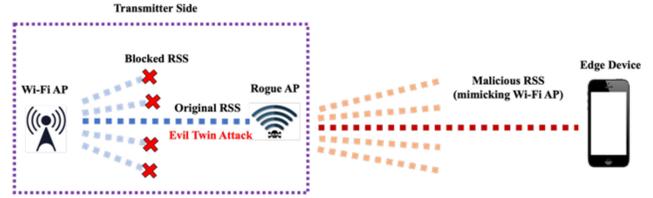


Fig. 4. Evil twin attack during indoor localization.

291 typically range from -100 dBm (weak signal) to 0 dBm
 292 (strong signal). These fingerprints are very susceptible to
 293 fluctuations due to dynamic indoor environments and edge
 294 device heterogeneity, and perturbations due to adversarial
 295 attacks especially in the presence of rogue APs, as shown in
 296 Fig. 3. Rogue APs can perturbate specific or all datapoints
 297 within an RSS fingerprint. This perturbed data may exhibit
 298 features characteristic of a different RP location, leading to
 299 increased prediction errors, as shown in Fig. 3.

300 Rogue APs pose a threat to indoor localization systems
 301 by introducing deliberate perturbations through two distinct
 302 pathways: 1) the transmitter side and 2) the channel side.

303 1) *Transmitter Side*: This attack is executed from the
 304 transmitter side, specifically on the APs deployed in
 305 the indoor environment. The attack targets a legitimate
 306 AP in the environment, attempting to infect it with
 307 malicious data (malware). Once successful, the resulting
 308 rogue AP gains complete control over the legitimate AP,
 309 compromising the security of any operations performed
 310 by the legitimate AP. This poses a significant security
 311 risk, as the rogue AP can now manipulate RSS, leading
 312 to an increase in localization errors. This attack can
 313 compromise the robustness of the indoor localization
 314 solution in that environment.

315 2) *Channel Side*: This attack is executed from the channel
 316 side, specifically within the spatial domain between
 317 a legitimate AP and the edge device. The rogue AP
 318 monitors communication between the legitimate AP
 319 and edge devices and introduces carefully calibrated
 320 interference with the signals traveling through this space.
 321 Once successful, the rogue AP can manipulate the RSS
 322 captured by the edge device, that may mimic the char-
 323 acteristics of a different RP location. This manipulation
 324 compromises the robustness of the indoor localization
 325 solution, as the altered RSS can lead to increase in
 326 localization errors.

327 A. Rogue AP Attack Implementation

328 Rogue APs possess the capability to execute a variety of
 329 attacks. Notably, these attacks can be launched with minimal
 330 information about the target system, rendering them as gray-
 331 box attacks. The nature of gray-box attacks makes rogue APs
 332 an attractive choice for adversaries, as they do not require
 333 comprehensive knowledge of the indoor localization system.

This characteristic transforms rogue AP implementation into a
 334 more plug-and-play system for executing adversarial attacks.
 335 We next describe the two types of rogue AP attacks, illustrat-
 336 ing their underlying methods and potential consequences.
 337

338 1) *Evil Twin Attacks*: This transmitter side rogue AP attack
 339 involves the creation of a malicious wireless network
 340 that mimics a legitimate one. The rogue AP utilizes
 341 malware to infect a legitimate AP, allowing it to gather
 342 critical information, such as the service set identifier
 343 (SSID), media access control identifier (MACID), and
 344 other network parameters [36]. By replicating these
 345 parameters, the rogue AP tricks edge devices into
 346 connecting to it, masquerading as an authentic AP. Fig. 4
 347 demonstrates the implementation of the evil twin attack,
 348 which is explored for the first time in the context of
 349 indoor localization, as part of this work. The rogue AP
 350 initiates the attack by targeting a legitimate Wi-Fi AP,
 351 mimicking its network parameters, and simultaneously
 352 blocking all communications from the legitimate Wi-Fi
 353 AP. Subsequently, the rogue AP broadcasts its own
 354 malicious Wi-Fi network (masquerading the authentic
 355 AP), that can inject malicious features into the RSS
 356 fingerprint collected by the edge device. These malicious
 357 features have the potential to falsify the edge device's
 358 perceived location, making it appear in a different location.
 359 This compromise in location information poses a severe
 360 threat to the entire indoor localization system.

361 2) *Man-in-the-Middle Attacks*: This channel side rogue
 362 AP attack employs ARP (address resolution protocol)
 363 spoofing techniques to intercept communication between
 364 the legitimate Wi-Fi AP and the edge devices [37].
 365 Operating within the spatial domain between the AP
 366 and the edge device, the rogue AP positions itself as an
 367 intermediary, intercepting signals transmitted between
 368 the legitimate AP and the edge device. Unlike direct
 369 communication, the man-in-the-middle attack allows
 370 the rogue AP to inspect, modify, or block the signals
 371 before relaying them to their intended destination. This
 372 interception provides the adversary with the capability to
 373 alter RSS values in real-time, introducing discrepancies
 374 in the RSS features captured by the edge device. Fig. 5
 375 demonstrates the implementation of the man-in-the-
 376 middle attack for indoor localization.

377 B. Adversarial Attack Methods

378 Adversarial perturbations, introduced by malicious entities,
 379 pose a threat to ML models, particularly in privacy-sensitive
 380 domains like indoor localization. We identify and focus on

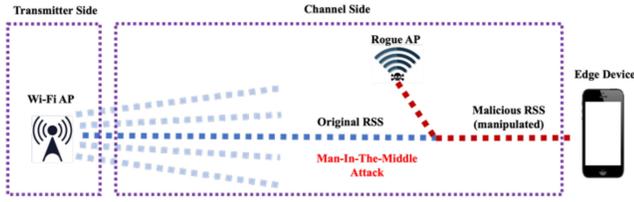


Fig. 5. Man-in-the-middle attack during indoor localization.

three popular adversarial methods in this work: 1) FGSM [30]; 2) projected gradient descent (PGD) [33]; and 3) momentum iterative method (MIM) [34]. Given the gray box nature of adversarial attacks (evil twin and man-in-the-middle attacks, discussed above), adversaries exploit minimal information about the localization framework. These methods introduce carefully calibrated perturbations into the RSS fingerprints using the ML model’s loss function, making them a practical choice for studying the nuanced effects of adversarial attacks in indoor localization.

1) *FGSM*: FGSM leverages the gradient information of the ML model’s loss function with respect to the input data. This method perturbs the original input data by adding a small, controlled perturbations in the direction of the gradient sign. This intentional perturbation systematically alters both the magnitude and positions of features within the input data. Consequently, this perturbation can mislead the ML model by indicating features at a different RP location, thereby increasing errors in location predictions

$$\eta = \epsilon * \text{sign}(\nabla J(\theta, X, Y)) \quad (1)$$

$$X_{\text{Adv}} = X + \eta. \quad (2)$$

In the equations above, η represents the perturbations, θ represents the parameters of the ML model, and X and Y denote the RSS fingerprint and RP class, respectively. The hyperparameter ϵ controls the magnitude of the perturbation and $(\nabla J(\theta, X, Y))$ denotes the loss function of the ML model. X_{Adv} is the perturbed RSS data.

2) *PGD Method*: PGD extends the concepts of FGSM by offering a more sophisticated approach in generating adversarial examples. PGD modifies FGSM by eliminating the sign function in (1) and clipping the perturbations between X and ϵ . While FGSM introduces perturbations in a single step, PGD refines the perturbation over multiple iterations $\{X_{\text{Adv}}(0), X_{\text{Adv}}(1), \dots, X_{\text{Adv}}(N), X_{\text{Adv}}(N+1)\}$.

$$X_{\text{Adv}}(0) = X \quad (3)$$

$$\eta = \text{Clip}_{X, \epsilon} \left\{ \epsilon * \frac{\nabla J(\theta, X, Y)}{L|\nabla J(\theta, X, Y)|_2} \right\} \quad (4)$$

$$X_{\text{Adv}}(N+1) = X_{\text{Adv}}(N) + \eta. \quad (5)$$

In (3), X denotes the original input data and $X_{\text{Adv}}(0)$ denotes the perturbed adversarial sample at the initial iteration (0). Equation (4) computes perturbations η using a clipped function applied to the gradients of the loss function $\nabla J(\theta, X, Y)$ and $L|\nabla J(\theta, X, Y)|_2$ represents

the squared L2 norm (ridge regularization) of the gradients of the loss function. This normalization step ensures that the perturbation is scaled appropriately, maintaining stability in generating the adversarial sample, while being clipped between X and ϵ (magnitude of the perturbation). These perturbations are added to $X_{\text{Adv}}(N)$ iteratively, as shown in (5). This iterative refinement process enhances the potency of adversarial samples by introducing a more calibrated manipulation in feature magnitude and positions within the RSS fingerprint data, leading to more potent adversarial samples compared to FGSM.

3) *MIM*: MIM further refines the adversarial samples from PGD, by incorporating momentum into the perturbation generation process to enhance the efficiency of the perturbation search

$$X_{\text{Adv}}(N+1) = \text{Clip}_{X, \epsilon} \{ \alpha * X_{\text{Adv}}(N) + \eta \}. \quad (6)$$

The perturbation η is calculated using (4), similar to the PGD approach. In (6), α is applied as momentum to the $X_{\text{Adv}}(N)$ of the previous iteration, while being clipped between X and ϵ (magnitude of the perturbation). By incorporating momentum into the perturbation generation process, MIM effectively manipulates RSS features and positions, leading to adversarial samples that induce more significant errors in the localization process, compared to FGSM and PGD. This enhanced perturbation poses substantial challenges to the robustness of indoor localization solutions.

C. Adversarial Attack Formulation for ML Indoor Localization

In formulating adversarial attacks for indoor localization systems, we employ the three distinctive methods discussed above: FGSM, PGD, and MIM. Our objective is to generate adversarial data by introducing perturbations that modify the features embedded within an RSS fingerprint. To generate potential real-world adversarial data effectively, we leverage two key parameters.

1) *Perturbation Strength (ϵ)*: This crucial hyperparameter is used in FGSM, PGD, and MIM methods to introduce perturbations to the RSS fingerprints. In generating adversarial samples for indoor localization, we systematically adjust the ϵ value to encompass various perturbation strengths applicable in real-world scenarios. We vary ϵ from 0.1 to 0.5 to reflect a practical perturbation scenario tailored for indoor localization [39]. This range is considered acceptable because it strikes a balance between being subtle enough to evade detection and significant enough to effectively test the system’s robustness. Smaller values of ϵ (closer to 0.1) represent minor perturbations that are less likely to be noticed but might not challenge the system’s defenses effectively, while larger values (up to 0.5) represent more noticeable perturbations that can more rigorously test the model’s resilience.

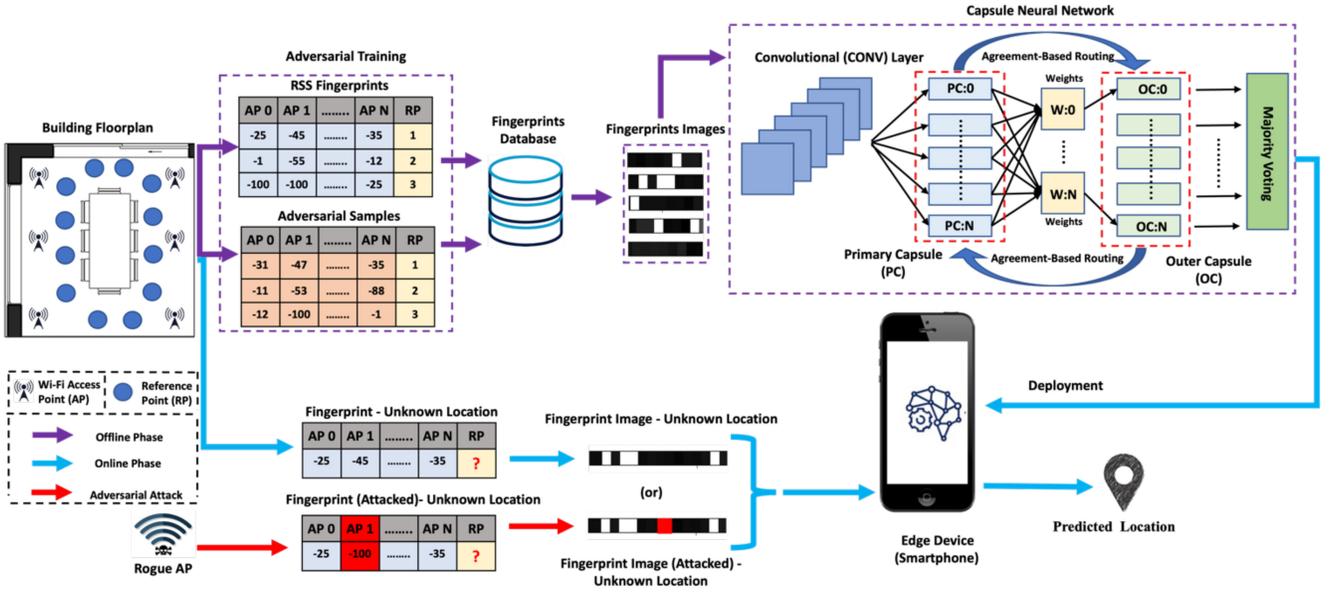


Fig. 6. Overview of the SENTINEL framework, including the offline (training) phase and online (inference) phase.

479 2) *Compromised APs* (φ): This parameter represents the
 480 quantity of legitimate APs that are subject to com-
 481 promise by the rogue AP within the indoor system.
 482 In a typical scenario, rogue APs selectively attack a
 483 subset of legitimate APs. We utilize φ as a parameter to
 484 investigate the impact of the quantity of compromised
 485 APs on indoor localization performance. φ is
 486 set to range from 0 to 100, indicating the percentage of
 487 attacked APs, thus covering the spectrum from 0% to
 488 100% of compromised APs. These attacked APs then
 489 introduce perturbations defined by the parameter ϵ .

490 IV. SENTINEL FRAMEWORK: OVERVIEW

491 The SENTINEL framework consists of three key compo-
 492 nents: 1) adversarial training; 2) fingerprint image generation;
 493 and 3) and the capsule neural network, as shown in Fig. 6.
 494 The framework initiates in an offline phase, where RSS
 495 fingerprints are captured across different RPs within the
 496 building floorplan. Multiple fingerprints are collected per RP
 497 to effectively capture data variability. These fingerprints are
 498 labeled and stored in an RSS fingerprint database, forming
 499 the offline training data for the SENTINEL framework. To
 500 fortify the framework against adversarial attacks, we employ
 501 an adversarial training mechanism (discussed in Section IV-A),
 502 which introduces adversarial samples derived from the RSS
 503 fingerprint database. Post-adversarial training, we transform
 504 both original (from RSS fingerprint database) and adversarial
 505 fingerprints into fingerprint images using the fingerprint image
 506 generation mechanism (discussed in Section IV-B), resulting
 507 in grayscale images. These grayscale images encapsulate
 508 crucial information about the indoor floorplan. The grayscale
 509 images then serve as input to the capsule neural network
 510 modified for the task at hand and carefully designed to address
 511 the spatial invariance problem in DL. The capsule neural
 512 network comprises five subcomponents: 1) convolutional layer

(CONV); 2) primary capsule (PC) layer; 3) outer capsule (OC) 513
 layer; 4) an agreement-based routing algorithm; and 5) the 514
 majority voting layer (all discussed in Section IV-C). 515

The domain-specific capsule neural network, once trained, 516
 is deployed on edge devices for predictions during the online 517
 phase. In the online phase, the edge devices (with the pre- 518
 trained ML model), scan for available RSS fingerprints at an 519
 unknown RP location. These received fingerprints are inher- 520
 ently susceptible to RSS fluctuations and potential adversarial 521
 attacks (introduced by rogue APs). 522

523 A. Adversarial Training Mechanism

The SENTINEL framework enhances its resilience against 524
 adversarial attacks by implementing an adversarial training 525
 mechanism. This approach fortifies our capsule neural net- 526
 work by exposing it to a diverse mixture of adversarial and clean 527
 RSS examples during the training process. The fundamental 528
 concept behind adversarial training is to modify the loss 529
 function by incorporating adversarial examples, thereby ren- 530
 dering the capsule neural network resistant to adversarial 531
 attacks 532

$$\nabla J(\theta, X, Y) = \nabla J(\theta, X, Y) + \nabla J(\theta, X + \eta, Y). \quad (7) \quad 533$$

In (7), η represents the perturbation introduced into the input 534
 data using different adversarial methods, such as FGSM, PGD, 535
 and MIM, calculated using the gradients of the loss function 536
 [(1), (3), and (5)] with respect to the input data. In Section V, 537
 we evaluate the performance of various adversarial training 538
 methods to assess SENTINEL's efficacy in defending against 539
 adversarial attacks in the online phase. 540

541 B. Fingerprint Image Generation

Post creation of the RSS fingerprint database (with clean 542
 + adversarial samples), the fingerprints are transformed into 543
 grayscale images to encapsulate crucial information about the 544

indoor floorplan. Initially, the RSS fingerprints are arranged into matrices or tensors, with shape of (H, W) , where H represents the height (typically 1), and W signifies the width, representing the number of visible APs within the indoor environment. Each element in this tensor corresponds to the RSS measured by a specific AP at a particular RP. To convert these RSS fingerprint tensors into grayscale images, a mapping process is applied. This mapping function translates the RSS values into pixel intensities, ensuring that higher RSS values are represented with brighter pixels and lower RSS values with darker pixels. The resulting grayscale images have a shape of (N, H, W, C) , where N denotes the RPs, H represents the height (usually 1), W signifies the width (number of visible APs), and C represents the number of channels (typically 1 for grayscale). This conversion preserves the spatial information of RSS across the indoor space, facilitating effective localization.

C. Capsule Neural Network Architecture

The capsule neural network is a pivotal component of the SENTINEL framework, comprising five subcomponents: 1) the convolutional (CONV) layer; 2) PC layer; 3) OC layer; 4) an agreement-based routing algorithm; and 5) a majority voting layer. The enhanced capsule neural network in SENTINEL possesses several key differences from EDGELOC [28] which uses a simple capsule neural network: 1) unlike [28], SENTINEL integrates a majority voting layer to enhance prediction output; 2) unlike [28], SENTINEL is tailored specifically for processing grayscale fingerprint images; 3) [28] targets device heterogeneity only, whereas SENTINEL optimizes hyperparameters differently to simultaneously target mitigation of dynamic environment induced RSS fluctuations, device heterogeneity, and adversarial attacks; and 4) SENTINEL is pruned in the number of capsules (both PC and OC layers) and neurons within each capsule, resulting in a more lightweight deployment on resource-constrained edge devices than [28] while maintaining accuracy. We compare SENTINEL against EDGELOC [28] in Section V. In the rest of this section, we describe the various components of our SENTINEL capsule neural network.

1) *Convolutional (CONV) Layer*: The CONV layer captures spatial features within the grayscale fingerprint images. This layer employs convolutional filter kernels to extract distinctive patterns and features from the input images. Let us denote the grayscale RSS fingerprint image as IM , which has dimensions (N, H, W, C) . The convolutional layer consists of multiple filters kernels, denoted as F , which are applied to IM . The F slide across the entire IM , performing element-wise multiplications and summations, generating feature maps that highlight spatial features within the IM

$$\text{CON}(p, q) = \sum_{i=0}^H \sum_{j=0}^W \text{IM}(p-i, q-j) * F(i, j). \quad (8)$$

In the equation above, $\text{CON}(p, q)$ denotes the feature at position (p, q) in the CONV feature map and $F(i, j)$ represents the corresponding element of the filter kernel.

$IM(p-i, q-j)$ represents the pixel value of IM at position $(p-i, q-j)$. The summation is performed over the height (H) and width (W) of F . During training, the network learns the optimal values of F through backpropagation. This process enables the CONV layer to automatically detect and extract relevant spatial features from the input RSS fingerprint images, providing meaningful representations that contribute to the overall accuracy of the localization process.

2) *PC Layer*: The PC layer receives the spatial features extracted by the CONV layer and serves as the next processing stage in the capsule neural network. A capsule is defined as a group of neurons, where each capsule within the PC layer generates a vector, referred to as the ‘‘activity vector.’’ This vector captures both the magnitude (presence) and position of each feature in the RSS fingerprint. Unlike traditional neural networks (such as MLPs and CNNs) where neurons in subsequent layers are densely connected to all neurons in the preceding layer, the PC layer consists of capsules, where each capsule corresponds to a specific spatial feature detected by the CONV layer. The activity vector (u_{ij}) for capsule i is obtained through a series of computations

$$S_i = \sum_j V_{ij} * \text{CON}_j \quad (9)$$

$$u_{ij} = \text{Squash}(S_i) = \frac{\|S_i\|^2}{1 + \|S_i\|^2} * \frac{S_i}{\|S_i\|}. \quad (10)$$

In (9), S_i represents the input for each capsule i , which is calculated as the weighted sum of outputs from the CONV layer using weight tensors (V_{ij}). These weight tensors determine the contribution of each feature from the CONV layer, enabling the PC layer to selectively focus on relevant spatial features. Subsequently, S_i is squashed using a nonlinear activation function known as the squash function. The squash function transforms S_i into activity vectors u_{ij} , which represent the magnitude and position of the detected spatial features within the RSS fingerprint. This enables the PC layer to encode spatial relationships between features, enhancing the network’s ability to capture meaningful representations of the indoor environment.

3) *OC Layer*: The OC layer performs classifications based on the activity vectors (u_{ij}), received from the PC layer. Each capsule in the OC layer corresponds to an RP class which determines the probability of the input fingerprint image belonging to that class. The classification process in the OC layer involves computing the agreement score between the u_{ij} and the weights tensors (W_{ij}) associated with each capsule in the OC layer

$$a_i = u_{ij} * W_{ij} \quad (11)$$

$$P_i = \text{Softmax}(a_i). \quad (12)$$

In (11), a_i represents the agreement score for capsule i . The W_{ij} contains the weight tensors associated with the connections between the PC and OC layers, determining the importance of each spatial feature for the

classification of the corresponding RP class. In (12), P_i denotes the predicted RP of capsule i after applying the Softmax function to a_i from (11). This function assigns probabilities to each RP class based on a_i , facilitating the classification process.

- 4) *Agreement-Based Routing Algorithm*: The agreement-based routing algorithm plays a crucial role in refining the weight tensors (W_{ij}) between the PC and OC layers. After the OC layer receives activity vectors (u_{ij}) from the PC layer, the agreement scores (a_i) are computed using (11), representing the agreement between the u_{ij} and W_{ij} associated with each capsule in the OC layer. The goal of the routing algorithm is to iteratively adjust these weight tensors based on the a_i achieved. The routing process involves several iterative steps, where a_i are used to update the W_{ij} in a way that maximizes agreement between the a_i and the predicted RP classes. This iterative refinement enhances the network's ability to accurately classify input fingerprint images.
- 5) *Majority Voting Layer*: The majority voting layer is the final component of the proposed capsule neural network. This layer aggregates the predictions (P_i) generated by the OC layer for each capsule. The majority voting mechanism aims to determine the final prediction by selecting the RP class with the highest number of aligned predictions from the capsules in the OC layer

$$\text{Prediction} = \text{Argmax}(P_0, P_1, \dots, P_n). \quad (13)$$

In (13), n represents the total number of RP classes. The Argmax function selects the RP class with the highest probability as the final prediction. By ensuring that a majority of capsules agree on the final class, the majority voting layer reduces the impact of erroneous predictions from individual capsules.

V. EXPERIMENTS

A. Experimental Setup

In this section, we describe our experimental setup, designed to evaluate the performance of our proposed SENTINEL framework in real-world scenarios. Our objective is to conduct comprehensive comparisons with state-of-the-art indoor localization frameworks, including CNNLOC [21], VITAL [27], EDGELOC [28], ADVLOC [31], and CALLOC [32], using simulated (FGSM, PGD, and MIM) and real-world *RSSRogueLoc* [35] data. Data was collected during regular working hours, incorporating both dynamic and static occupants to reflect realistic conditions. Table I shows an overview of the real devices utilized in our experiments.

To ensure a comprehensive evaluation across diverse environmental conditions, we select building floorplans with varying factors, such as path length, the number of visible APs, and environmental noise characteristics, as shown in Fig. 7. Our data collection strategy is designed to facilitate thorough training and testing of the SENTINEL framework. For each building floorplan, we allocate five fingerprints per RP for training and one fingerprint per RP, per device, and per building, for testing. Acknowledging the substantial effort

TABLE I
DEVICES USED TO COLLECT RSS FINGERPRINTS

Device Name	Wi-Fi Chipset	Acronym	Year
BLU Vivo 8	MediaTek Helio P10	BLU	2017
Google Pixel 6a	Google Tensor G1	GOOGLE	2022
HTC U11	Qualcomm Snapdragon 835	HTC	2017
Motorola Z2	Qualcomm Snapdragon 835	MOTO	2017
Nokia 7.1	Qualcomm Snapdragon 636	NOKIA	2018
OnePlus Nord 200	Qualcomm Snapdragon 480	ONEPLUS	2021
Xiaomi Redmi 10A	MediaTek Helio G88	REDMI	2022
Samsung A14	Samsung Exynos 850	SAMSUNG	2023

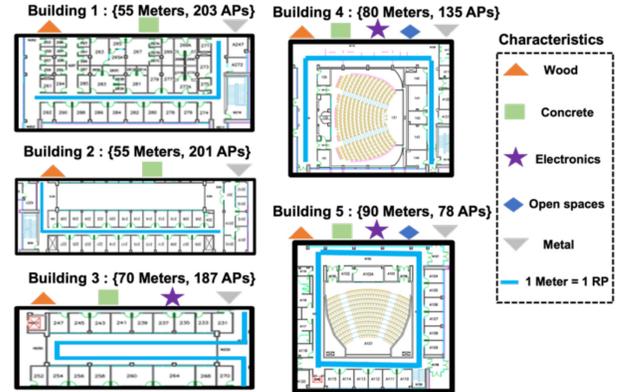


Fig. 7. Building floorplan layouts with varying path length, visible APs, and characteristics.

required to gather a large volume of offline training data, we restrict the collection of offline data to a single device. To facilitate this, we designate the MOTO device as the primary training device. All devices in Table I are used in the online phase during testing.

The SENTINEL framework is configured with specific architectural hyperparameters. The CONV layer is equipped with 32 filters and the PC layer comprises eight capsules with each capsule containing a dimension of 32 neurons. Furthermore, the OC layer contains capsules equal to the number of RP classes with a dimension of 32 neurons each, trained over 300 epochs using the Adam optimizer (learning rate = 0.001) and the sparse categorical cross-entropy loss function. The capsule neural network architecture results in a total of 2 117 687 trainable parameters, with a compact model size of 8.07 MB, facilitating low overhead deployment on most resource-constrained edge devices. Additionally, the SENTINEL framework incorporates an adversarial training mechanism aimed at enhancing its resilience against potential adversarial attacks. Adversarial samples are generated using the FGSM, PGD, and MIM approaches with ϵ set to 0.1 and φ set to 100% (for training only). Each variant of our trained capsule neural network, augmented with adversarial samples, is denoted with a suffix. For instance, the model trained without adversarial samples is referred to as SENTINEL-NONE, while models trained with FGSM, PGD, and MIM samples are labeled SENTINEL-FGSM, SENTINEL-PGD, and SENTINEL-MIM, respectively.

735 *B. Effects of Adversarial Training on Heterogeneity*

736 In this section, we evaluate the performance of the SENTINEL
 737 framework under various adversarial training scenarios (FGSM,
 738 PGD, and MIM), separately. In Fig. 8, we present heatmaps
 739 depicting the performance of the four SENTINEL variants:
 740 1) SENTINEL-NONE; 2) SENTINEL-FGSM; 3) SENTINEL-
 741 PGD; and 4) SENTINEL-MIM. These models are individually
 742 trained on data collected exclusively from a single device
 743 (MOTO) and incorporate their respective adversarial training
 744 techniques. SENTINEL-NONE is trained without including any
 745 adversarial samples, providing a comparison of the effects of
 746 including adversarial training to the SENTINEL framework.
 747 Evaluation of these model variants are conducted using data
 748 acquired from all eight available devices across the five building
 749 floorplans, without any adversarial interference.

750 In Fig. 8, the x -axis of each heatmap represents the testing
 751 devices, while the y -axis corresponds to the different buildings
 752 used for evaluation. Each cell within the heatmap indicates
 753 the average prediction error (in meters) across all RPs for
 754 a specific combination of test device and building floorplan.
 755 We observe differences in prediction errors across all the
 756 SENTINEL variants, due to the differences in adversarial
 757 training methods used. We note an increase in prediction errors
 758 when going from buildings 1–5, which can be attributed to
 759 increasing environmental dynamic causing higher variations in
 760 the selected building paths. For instance, building 1 exhibited
 761 low environmental noise, likely due to fewer people moving
 762 along the path during the testing. It also had relatively shorter
 763 path lengths, which overall resulted in lower prediction errors.
 764 In contrast, building 5 experienced higher environmental noise
 765 due to significantly more people moving along the path during
 766 the testing phase, and longer path lengths, leading to higher
 767 prediction errors. SENTINEL-FGSM consistently exhibits
 768 the lowest prediction errors, followed by SENTINEL-PGD,
 769 SENTINEL-NONE and SENTINEL-MIM. This trend suggests
 770 that while more advanced adversarial training methods like
 771 PGD and MIM may offer refined perturbations, they also
 772 introduce complexity and potential instability during training,
 773 leading to overfitting. The overfitting occurs because the adver-
 774 sarial samples generated by PGD and MIM involve multiple
 775 iterations of perturbations, making them more complex and
 776 causing feature mismatches between RP classes. As a result,
 777 the model may become overly specialized to these adversarial
 778 examples, reducing its ability to generalize well to unseen,
 779 real-world data. SENTINEL-FGSM however, stands out due
 780 to its balance between perturbation effectiveness and model
 781 stability. Its noniterative nature allows for smaller, controlled
 782 perturbations, reducing the chances of a feature mismatch
 783 between legitimate and FGSM samples.

784 To illustrate the impact of device heterogeneity and assess
 785 the performance of the SENTINEL variants, we present
 786 Fig. 9 more clearly. Here, the x -axis represents the testing
 787 devices, and the y -axis denotes the prediction error in meters.
 788 Each bar represents the average prediction error per device
 789 across all building floorplans, with error bars included to
 790 indicate the range of errors observed per testing device,
 791 with the lower whisker representing the best case and the

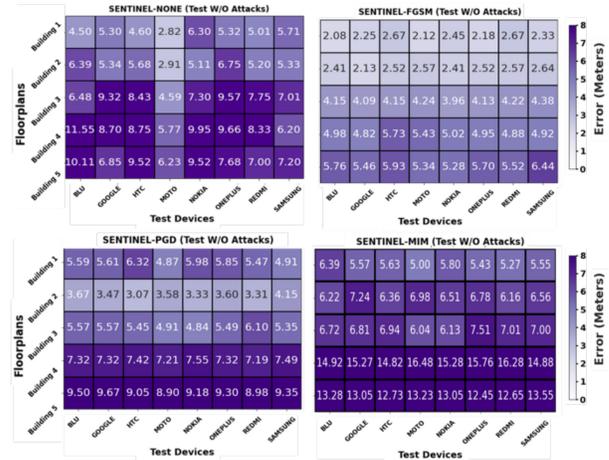


Fig. 8. Performance of the SENTINEL variants across different devices and building floorplans.

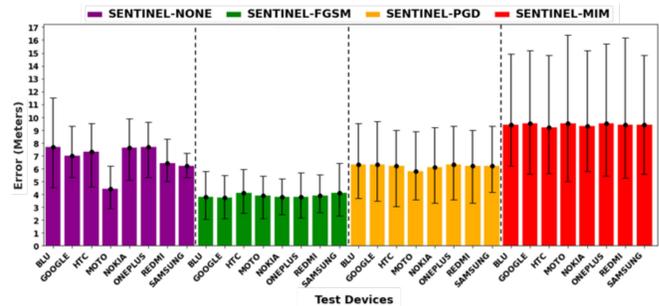


Fig. 9. Performance summary for SENTINEL variants.

792 upper whisker representing the worst-case location error. In
 793 Fig. 9, we observe that the average error per testing device
 794 remains consistent for each SENTINEL variant. However,
 795 the SENTINEL-NONE variant exhibits the least consistency
 796 in prediction errors across the testing devices, with some
 797 devices showing higher errors while others show lower
 798 errors. This suggests lower resilience to heterogeneity for
 799 the SENTINEL-NONE variant. Conversely, other SENTINEL
 800 variants show consistent prediction errors regardless of the
 801 training or testing devices used, indicating better heterogeneity
 802 resilience. Furthermore, incorporating adversarial training not
 803 only strengthens the robustness of the SENTINEL variants
 804 against adversarial attacks but also improves their resilience
 805 to heterogeneity. By subjecting the models to adversarial per-
 806 turbations during training, the variants learn more generalized
 807 features, making them less sensitive to fluctuations from the
 808 testing devices. Particularly noteworthy is the performance of
 809 SENTINEL-FGSM, with up to $1.48\times-2.43\times$ lower average
 810 and worst-case errors compared to the rest of the SENTINEL
 811 variants.

812 *C. Evaluating the Impact of Varying Compromised APs (φ)*

813 In this section, we investigate the impact of varying φ
 814 in the testing phase, using different adversarial attack meth-
 815 ods (FGSM, PGD, and MIM), on the performance of the
 816 SENTINEL variants. To maintain consistency, we set the

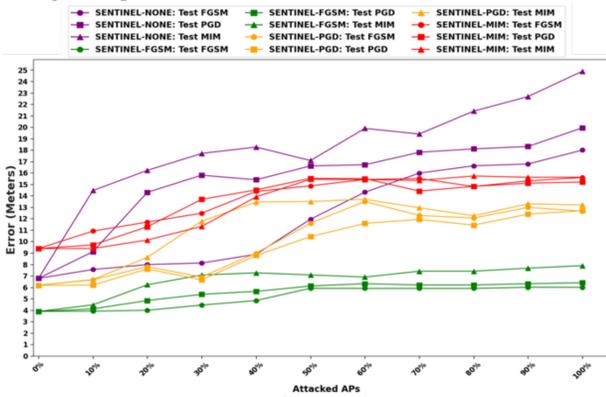


Fig. 10. Performance of the four SENTINEL variants on simulated adversarial attacks through varying φ .

817 attack ϵ to 0.1, indicating 10% added perturbations per φ . In
 818 Fig. 10, the x -axis represents φ , ranging from 0 (no attacked
 819 APs) to 100 (all visible APs being attacked). The y -axis
 820 denotes prediction errors measured in meters and the line plots
 821 illustrate the performance of each SENTINEL variant under
 822 the three adversarial attack methods. In Fig. 10, each marker
 823 indicates the average prediction error across all testing devices
 824 and building floorplans.

825 We observe that as φ increases, the prediction errors for
 826 all SENTINEL variants also increase. However, there is a
 827 stabilization point observed at approximately $\varphi = 50\%$ for
 828 most variants methods (except SENTINEL-NONE, which
 829 lacks adversarial training), suggesting that the performance
 830 of the SENTINEL variants remains relatively unaffected
 831 when a significant portion of APs are compromised. This
 832 stabilization point indicates that the SENTINEL variants are
 833 resilient to attacks involving large numbers of compromised
 834 APs. Additionally, most variants demonstrate resilience against
 835 various adversarial attack methods (except SENTINEL-
 836 NONE), as evidenced by the almost flat line in prediction
 837 errors. Specifically, when subjected to the FGSM attack,
 838 the SENTINEL-FGSM model exhibits 1.90 \times , 2.35 \times , and
 839 2.64 \times lower average errors compared to the SENTINEL-PGD,
 840 SENTINEL-NONE, and SENTINEL-MIM models, respec-
 841 tively. Similarly, under the PGD attack, the SENTINEL-FGSM
 842 model demonstrates 1.69 \times , 2.75 \times , and 2.40 \times lower average
 843 errors compared to the SENTINEL-PGD, SENTINEL-NONE,
 844 and SENTINEL-MIM models, respectively. Lastly, when influ-
 845 enced by the MIM attack, the SENTINEL-FGSM model shows
 846 1.67 \times , 2.71 \times , and 2.15 \times lower average errors compared to
 847 the SENTINEL-PGD, SENTINEL-NONE, and SENTINEL-
 848 MIM models, respectively.

849 D. Evaluating the Impact of Varying Perturbations (ϵ)

850 In this section, we explore the impact of varying levels
 851 of perturbation strength (ϵ) in the testing phase on
 852 the performance of all SENTINEL variants. Our objective
 853 is to investigate how the prediction performance of each
 854 SENTINEL variant is affected by changes in ϵ , ranging from
 855 0 (indicating no attack) to 0.5 (representing a 50% increase
 856 in added perturbations). In Fig. 11, the x -axis represents the

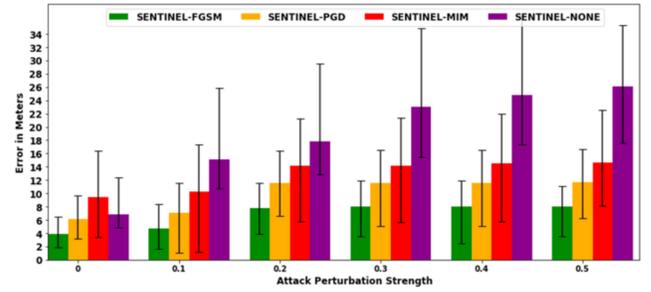


Fig. 11. Performance of the three SENTINEL variants on simulated adversarial attacks through varying ϵ .

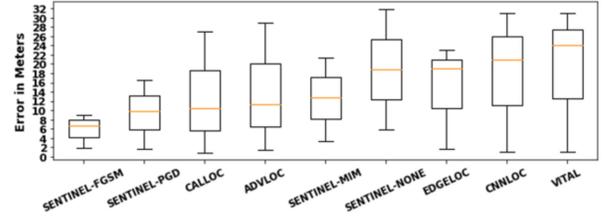


Fig. 12. Performance comparisons of all SENTINEL variants against state-of-the-art indoor localization frameworks.

857 varying levels of ϵ , while the y -axis denotes the prediction
 858 error in meters. Each bar in the plot signifies the average
 859 prediction error across all testing devices, building floorplans,
 860 and φ values. Additionally, error bars are included to depict
 861 the range between the best (lower whisker) and worst-case
 862 (upper whisker) prediction errors. Our analysis reveals that
 863 as ϵ increases, there is a slight rise in prediction errors.
 864 However, we observe that all SENTINEL variants stabilize
 865 at approximately $\epsilon = 0.2$ (except SENTINEL-NONE, lack-
 866 ing adversarial training). This suggests that regardless of
 867 the increase in perturbation strength, all SENTINEL models
 868 demonstrate consistent performance. Furthermore, we observe
 869 that the SENTINEL-FGSM variant consistently outperforms
 870 all other SENTINEL variants. On average, SENTINEL-
 871 FGSM demonstrates 1.48 \times , 2.81 \times , and 1.90 \times lower average
 872 prediction errors compared to SENTINEL-PGD, SENTINEL-
 873 NONE, and SENTINEL-MIM, respectively. The superior
 874 performance of the SENTINEL-FGSM variant, even as ϵ
 875 increases during testing, can be attributed to the robustness
 876 gained through FGSM-based adversarial training. Although
 877 the model was trained with a fixed ϵ value, the adversarial
 878 training process encourages the model to capture underlying
 879 patterns in feature positions that are susceptible to adversarial
 880 attacks. This enables the model to generalize and adapt to
 881 perturbations even on varying ϵ . In contrast, other methods
 882 like PGD and MIM often induce significant perturbations
 883 in underlying features, leading to overfitting and reduced
 884 resilience during testing. The chosen epsilon range of 0–0.5
 885 represents a practical attack range for indoor localization [39].

886 E. Comparison Against State-of-the-Art Frameworks

887 In this section, we compare the performance of all
 888 SENTINEL variants against state-of-the-art indoor localization
 889 frameworks across various parameters, including different

TABLE II
MODEL PARAMETERS, SIZE OF ALL FRAMEWORKS

Framework	Total Parameters	Model Size
CALOC	652,390	2.48 MB
CNNLOC	858,720	3.27 MB
ADVLOC	1,746,752	6.99 MB
SENTINEL	2,117,687	8.07 MB
EDGELOC	2,317,687	8.84 MB
VITAL	2,347,006	8.95 MB

890 devices, building floorplans, ϵ (ranging from 0 to 0.5), and φ
 891 (ranging from 0 to 100). Fig. 12 presents a box and whiskers
 892 plot, showcasing the comparison of the best case (lower
 893 whisker), worst case (upper whisker), and average (orange
 894 line) errors across all frameworks. This enhanced resilience
 895 can be attributed to the adversarial training and capsule
 896 neural network employed by the SENTINEL framework.
 897 The FGSM-based adversarial training introduces optimal
 898 adversarial features and feature dispositions (magnitude and
 899 positions), contrasting with other adversarial training meth-
 900 ods that may lead to overfitting. The proposed capsule
 901 neural network treats each feature as a vector, effectively
 902 recognizing and capturing underlying patterns between the
 903 original (clean) and adversarial samples during training.
 904 This enables the SENTINEL-FGSM model to demonstrate
 905 lower prediction errors across various scenarios and metrics
 906 compared to the other frameworks. The SENTINEL-FGSM
 907 model demonstrates $1.47\times$, $1.55\times$, $1.68\times$, $1.91\times$, $2.82\times$,
 908 $2.83\times$, $3.13\times$, and $3.5\times$ lower average errors compared to
 909 SENTINEL-PGD, CALLOC, ADVLOC, SENTINEL-MIM,
 910 SENTINEL-NONE, EDGELOC, CNNLOC, and VITAL,
 911 respectively. Additionally, recognizing the need for lightweight
 912 frameworks adaptable for resource-constrained edge devices,
 913 we analyze the parameter count and memory footprint of the
 914 various frameworks as shown in Table II. SENTINEL yields
 915 a compact model size of 8.07 MB.

916 F. Evaluation on the New Real-World Rogue AP 917 Attack Dataset

918 In this section, we introduce a novel Wi-Fi RSS finger-
 919 print dataset named *RSSRogueLoc* [35], designed to capture
 920 the detrimental effects of rogue APs for indoor localization
 921 systems. Unlike prior works which primarily rely on sim-
 922 ulated adversarial attacks introduced by methods, such as
 923 FGSM, PGD, and MIM, *RSSRogueLoc* delves into real-world
 924 adversarial scenarios, particularly those involving rogue APs.
 925 Building on the dataset outlined in Section V-A, *RSSRogueLoc*
 926 introduces a secondary testing dataset comprising up to five
 927 new devices configured as rogue APs (devices detailed in
 928 Table III), designed to execute evil twin attacks as discussed
 929 in Section III-A, where each rogue is configured to impact one
 930 legitimate AP. The *RSSRogueLoc* fingerprints were collected
 931 by incrementally introducing rogue APs across all RPs within
 932 each building floorplan. This sequential escalation started from
 933 Rogue 0, signifying the absence of all rogues, followed by
 934 Rogue 1 with one rogue per RP per floorplan, Rogue 2 with
 935 two rogues per RP per floorplan, Rogue 3 with three rogues

TABLE III
ROGUE AP DEVICES USED IN *RSSRogueLoc*

Device Name	Wi-Fi Chipset	Device Type
Samsung G991U	Samsung Exynos 2100	Smartphone
Apple A2789	Apple U2	Laptop
HP 840 G6	Intel Wi-Fi AX201	Laptop
Vivo V2025	Qualcomm Snapdragon 720G	Smartphone
HP 840 G10	Intel Wi-Fi AX211	Laptop

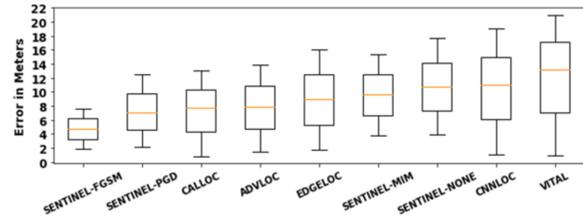


Fig. 13. Performance comparisons of all SENTINEL models against state-of-the-art on the *RSSRogueLoc* dataset.

per RP per floorplan, Rogue 4 with four rogues per RP per 936
 floorplan, and finally Rogue 5 with five rogues per RP per 937
 floorplan. The testing fingerprints were collected using all 938
 eight devices mentioned in Table I. This process unfolded 939
 over several weeks, to thoroughly capture the complexities of 940
 rogue AP configurations across numerous RPs and building 941
 floorplans. 942

To provide additional insights into the performance 943
 of all SENTINEL variants and state-of-the-art baseline 944
 frameworks on the *RSSRogueLoc* dataset, we present Fig. 13. 945
 The SENTINEL-FGSM model demonstrates $1.51\times$, $1.65\times$, 946
 $1.68\times$, $1.91\times$, $2.04\times$, $2.27\times$, $2.34\times$, and $2.80\times$ 947
 lower average error compared to SENTINEL-PGD, CALLOC, 948
 ADVLOC, EDGELOC, SENTINEL-MIM, SENTINEL- 949
 NONE, CNNLOC, and VITAL, respectively. 950

951 VI. CONCLUSION

The SENTINEL framework proposed in this work exhibits 952
 resilience against RSS fluctuations arising from environmental 953
 noise, edge device heterogeneity, and challenging adversarial 954
 attacks, due to its novel combination of adversarial training 955
 and modified capsule neural networks, while being relatively 956
 lightweight for edge device deployment. Through rigorous 957
 evaluation, we found that the SENTINEL-FGSM variant 958
 consistently achieves the lowest indoor localization errors, 959
 outperforming all baseline frameworks by $1.47\times$ – $3.5\times$ in 960
 average errors and $1.83\times$ – $3.4\times$ in worst-case errors on sim- 961
 ulated adversarial attacks. Moreover, our introduction of the 962
RSSRogueLoc dataset, designed to capture real-world effects of 963
 rogue APs (performing evil twin attacks in real-time), further 964
 highlights the superiority of the SENTINEL-FGSM variant. 965
 With $1.51\times$ – $2.8\times$ lower average errors and $1.63\times$ – $2.74\times$ 966
 lower worst-case errors compared to other state-of-the-art 967
 frameworks. 968

969 REFERENCES

- [1] S. Tiku and S. Pasricha, "An overview of indoor localization techniques," 970
IEEE Consum. Electron. Mag., to be published. 971

- 972 [2] D. Lymberopoulos and J. Liu, "The microsoft indoor localization
973 competition: experiences and lessons learned," *IEEE Signal Process.*
974 *Mag.*, vol. 34, no. 5, pp. 125–140, Sep. 2017.
- 975 [3] A. Petrenko et al., "Exploring mobility indoors: An application of
976 sensor-based and GIS systems," *Trans. GIS*, vol. 18, no. 3, pp. 351–369,
977 2014.
- 978 [4] W. K. Zegeye, S. B. Amsalu, Y. Astatke, and F. Moazzami, "WiFi RSS
979 fingerprinting indoor localization for mobile devices," in *Proc. IEEE*
980 *UEMCON*, 2016, pp. 1–6.
- 981 [5] M. Shu, G. Chen, and Z. Zhang, "Efficient image-based indoor local-
982 ization with MEMS aid on the mobile device," *ISPRS J. Photogrammet.*
983 *Remote Sens.*, vol. 185, pp. 85–110, Mar. 2022.
- 984 [6] X. Tian, R. Shen, D. Liu, Y. Wen, and X. Wang, "Performance analysis
985 of RSS fingerprinting based indoor localization," *IEEE Trans. Mobile*
986 *Comput.*, vol. 16, no. 10, pp. 2847–2861, Oct. 2017.
- 987 [7] B. Yang, L. Guo, R. Guo, M. Zhao, and T. Zhao, "A novel trilateration
988 algorithm for RSSI-based indoor localization," *IEEE Sensors J.*, vol. 20,
989 no. 14, pp. 8164–8172, Jul. 2020.
- 990 [8] X. Hou, T. Arslan, A. Juri, and F. Wang, "Indoor localization for
991 Bluetooth low energy devices using weighted off-set triangulation
992 algorithm," in *Proc. ION GNSS*, 2016, pp. 2286–2292.
- 993 [9] F. Alhomayani and M. H. Mahoor, "Deep learning methods for
994 fingerprint-based indoor positioning: A review," *J. Location Based*
995 *Services*, vol. 14, no. 3, pp. 129–200, 2020.
- 996 [10] N. Singh, S. Choe, and R. Punmiya, "Machine learning based
997 indoor localization using Wi-Fi RSSI fingerprints: An overview," *IEEE*
998 *Access*, vol. 9, pp. 127150–127174, 2021.
- 999 [11] D. Duong, Y. Xu, and K. David, "The influence of fast fading and
1000 device heterogeneity on Wi-Fi fingerprinting," in *Proc. IEEE VTC*, 2018,
1001 pp. 1–5.
- 1002 [12] I. Alshami, N. Ahmad, and S. Sahibuddin, "RSS certainty: An efficient
1003 solution for RSS variation due to device heterogeneity in WLAN
1004 fingerprinting-based indoor positioning system," in *Proc. PICICT*, 2021,
1005 pp. 71–76.
- 1006 [13] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing
1007 adversarial examples," 2014, *arXiv:1412.6572*.
- 1008 [14] C. Szegedy et al., "Going deeper with convolutions," in *Proc. IEEE*
1009 *CVPR*, 2015, pp. 1–9.
- 1010 [15] W. Xue, X. Hua, Q. Li, W. Qiu, and X. Peng, "Improved clustering
1011 algorithm of neighboring reference points based on KNN for indoor
1012 localization," in *Proc. IEEE UPINLBS*, 2018, pp. 1–4.
- 1013 [16] J. Jadidi, M. Patel, and J. V. Miro, "Gaussian processes online observa-
1014 tion classification for RSSI-based low-cost indoor positioning systems,"
1015 in *Proc. IEEE ICRA*, 2017, pp. 6269–6275.
- 1016 [17] A. B. Adege, H. P. Lin, G. B. Tarekegn, and S. S. Jeng, "Applying deep
1017 neural network (DNN) for robust indoor localization in multi-building
1018 environment," *Appl. Sci.*, vol. 8, no. 7, p. 1062, 2018
- 1019 [18] M. Dakkak, B. Daachi, A. Nakib, and P. Siarry, "Multi-layer perceptron
1020 neural network and nearest neighbor approaches for indoor localization,"
1021 in *Proc. IEEE SMC*, 2014, pp. 1366–1373.
- 1022 [19] W. Liu, H. Chen, Z. Deng, X. Zheng, X. Fu, and Q. Cheng, "LC-DNN:
1023 Local connection based deep neural network for indoor localization with
1024 CSI," *IEEE Access*, vol. 8, pp. 108720–108730, 2020.
- 1025 [20] Y. Kim, H. Shin, Y. Chon, and H. Cha, "Smartphone-based Wi-Fi
1026 tracking system exploiting the RSS peak to overcome the RSS variance
1027 problem," *Pervas. Mobile Comput.*, vol. 9, no. 3 pp. 406–420, Jun. 2013.
- 1028 [21] X. Song et al., "CNNLoc: Deep-learning based indoor localization with
1029 WiFi fingerprinting," in *Proc. IEEE SUI*, 2019, pp. 589–595.
- 1030 [22] D. Gufran, S. Tiku, and S. Pasricha, "SANGRIA: Stacked autoencoder
1031 neural networks with gradient boosting for indoor localization," *IEEE*
1032 *Embedded Syst. Lett.*, vol. 16, no. 2, pp. 142–145, Jun. 2024.
- [23] S. Tiku, D. Gufran, and S. Pasricha, "Multi-head attention neural
1033 network for smartphone invariant indoor localization," in *Proc. IEEE*
1034 *IPIN*, 2022, pp. 1–8.
- [24] Z. Zhang, H. Du, S. Choi, and S. H. Cho., "Tips: Transformer based
1035 indoor positioning system using both CSI and doa of WiFi signal," *IEEE*
1036 *Access*, vol. 10, pp. 111363–111376, 2022.
- [25] G. Elsayed, P. Ramachandran, J. Shlens, and S. Kornblith, "Revisiting
1037 spatial invariance with low-rank local connectivity," in *Proc. PMLR*,
1038 2020, pp. 1–19.
- [26] O. Parkhi, A. Ved, and A. Zisserman, "Deep face recognition," in *Proc.*
1039 *BMVC*, 2015, pp. 1–12.
- [27] D. Gufran, S. Tiku, and S. Pasricha, "VITAL: Vision transformer
1040 neural networks for accurate smartphone heterogeneity resilient indoor
1041 localization," in *Proc. IEEE DAC*, 2023, pp. 1–6.
- [28] Q. Ye et al., "EdgeLoc: A robust and real-time localization system
1042 toward heterogeneous IoT devices," *IEEE Internet Things J.*, vol. 9,
1043 no. 5, pp. 3865–3876, Mar. 2022, doi: [10.1109/JIOT.2021.3101368](https://doi.org/10.1109/JIOT.2021.3101368)
- [29] A. Shafahi et al., "Adversarial training for free!," in *Proc. NIPS*, 2019,
1044 pp. 1–11.
- [30] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial machine learning
1045 at scale," 2016, *arXiv:1611.012366*.
- [31] X. Wang, X. Wang, S. Mao, J. Zhang, S. C. G. Periaswamy, and
1046 J. Patton, "Adversarial deep learning for indoor localization with channel
1047 state information tensors," *IEEE Internet Things J.*, vol. 9, no. 19,
1048 pp. 18182–18194, Oct. 2022.
- [32] D. Gufran and S. Pasricha, "CALLOC: Curriculum adversarial learning
1049 for secure and robust indoor localization," 2023, *arXiv:2311.06361*.
- [33] S. Finlayson, H. W. Chung, I. Kohane, and A. Beam, "Adversarial
1050 attacks against medical deep learning systems," 2018, *arXiv:1804.05296*.
- [34] Y. Dong et al., "Boosting adversarial attacks with momentum," in *IEEE*
1051 *CVPR*, 2018, pp. 1–12.
- [35] "EPIC-CSU: Heterogeneous RSSI indoor navigation," GitHub. 2023.
1052 [Online]. Available: [https://github.com/EPIC-CSU/heterogeneous-rssi-](https://github.com/EPIC-CSU/heterogeneous-rssi-indoor-nav)
1053 [indoor-nav](https://github.com/EPIC-CSU/heterogeneous-rssi-indoor-nav)
- [36] Q. Lu, H. Qu, Y. Zhuang, X. Lin, Y. Zhu, and Y. Liu, "A passive client-
1054 based approach to detect evil twin attacks," in *Proc. IEEE ICESSE*, 2017,
1055 pp. 233–239.
- [37] A. Mallik, "Man-in-the-middle-attack: Understanding in simple
1056 words," in *Cyberspace J. Pendidikan Teknologi Informatika*, vol. 2, no. 2,
1057 p. 109, 2019
- [38] S. Sabour, N. Frosst, and G. E. Hinton, "Dynamic routing between
1058 capsules," in *Proc. NeurIPS*, 2017, pp. 1–11.
- [39] M. Patil, X. Wang, and S. Mao, "Adversarial attacks on deep learning-
1059 based floor classification and indoor localization," in *Proc. ACM WSML*,
1060 2021, pp. 7–12.
- [40] D. Gufran and S. Pasricha, "FedHIL: Heterogeneity resilient fed-
1061 erated learning for robust indoor localization with mobile devices,"
1062 *ACM Trans. Embedded Comput. Syst.*, vol. 22, no. 5S, p. 125,
1063 2023.
- [41] D. Gufran, S. Tiku, and S. Pasricha, "STELLAR: Siamese multiheaded
1064 attention neural networks for overcoming temporal variations and device
1065 heterogeneity with indoor localization," in *Proc. IEEE ISPIN*, 2023,
1066 pp. 1–15.
- [42] S. Tiku, D. Gufran, and S. Pasricha, "Smartphone invariant indoor
1067 localization using multi-head attention neural network," in *Machine*
1068 *Learning for Indoor Localization and Navigation*. Cham, Switzerland:
1069 Springer Int., 2023.
- [43] D. Gufran, S. Tiku, and S. Pasricha, "Heterogeneous device resilient
1070 indoor localization using vision transformer neural networks," in
1071 *Machine Learning for Indoor Localization and Navigation*. Cham,
1072 Switzerland: Springer Int., 2023.