

# Removing the need for ground truth UWB data collection: self-supervised ranging error correction using deep reinforcement learning

Dieter Coppens, Ben Van Herbruggen, Adnan Shahid, *Senior member*, IEEE, Eli De Poorter

**Abstract**—Indoor positioning using UWB technology has gained interest due to its centimeter-level accuracy potential. However, multipath effects and non-line-of-sight conditions cause ranging errors between anchors and tags. Existing approaches for mitigating these ranging errors rely on collecting large labeled datasets, making them impractical for real-world deployments. This paper proposes a novel self-supervised deep reinforcement learning approach that does not require labeled ground truth data. A reinforcement learning agent uses the channel impulse response as a state and predicts corrections to minimize the error between corrected and estimated ranges. The agent learns, self-supervised, by iteratively improving corrections that are generated by combining the predictability of trajectories with filtering and smoothening. Experiments on real-world UWB measurements demonstrate comparable performance to state-of-the-art supervised methods, overcoming data dependency and lack of generalizability limitations. This makes self-supervised deep reinforcement learning a promising solution for practical and scalable UWB-ranging error correction.

**Index Terms**—indoor positioning, UWB, reinforcement learning, self-supervised, error-correction

## I. INTRODUCTION

PRECISE indoor positioning technology has attracted significant research interest in recent years due to its role in overcoming the limitations of global positioning system (GPS) in indoor environments for Internet of Things (IoT) applications such as assistive healthcare systems [1], sports tracking [2], smart logistics [3] and various location-based services [4]. Following this trend, Ultra-Wideband (UWB) technology has seen a surge in interest and become one of the more promising technologies for indoor positioning systems (IPS). UWB IPS can achieve centimeter-level positioning accuracy due to the wide bandwidth (>500 MHz) and very short time duration of the pulse (around 2 ns) [5]. While these signal characteristics make UWB more resilient to multipath effects (compared to traditional narrowband techniques such as SigFox, LoRa, Narrowband Internet of Things (NB-IoT), etc. [6], [7]). However, a major remaining challenge is correcting ranging errors caused by this multipath behavior in non-line-of-sight (NLOS) conditions [8], [9]. Current methods to detect and reduce errors caused by NLOS conditions rely mostly on machine/deep learning models trained using large datasets of UWB ranges and raw physical data like the channel impulse

response (CIR) [10]–[12] or calculated features [7] (e.g. amplitude of the signal, energy, power ratio, etc.) labeled with the true positions. While these approaches can lead to high performance, it comes with two major disadvantages. First, collecting such labeled data requires a tedious labeling effort and dataset collection, which requires specialized equipment and personnel with expertise in UWB positioning and ground truth data collection. Second, the usability is limited by the generalization problem, the accuracy of trained solutions drops severely in unseen environments. The unseen environments have different anchor topologies, different sizes, and different UWB hardware, or contain different types of objects that degrade the performance due to (1) variations in the CIR and (2) different UWB physical layer (PHY) properties. The generalization problem worsens the data collection problem as each unique environment requires the labeling effort to be repeated and even so, the environment may have changed by then. These two disadvantages have previously been addressed using (1) semi-supervised learning [11] [13], to reduce the data collection and (2) transfer learning to enable better performance in unseen environments while using only a few labeled samples [14], [15]. However, all these approaches still require some labeled samples, thus a tedious data collection effort. One other research proposes a self-supervised ranging error correction [16] that does not require ground truth or label collection. It uses classical location approaches to jointly estimate the location and range with a deep network in a Time Difference Of Arrival (TDoA) system. However, the learning here is limited as they do not use signal features to aid and improve the learning process, and they cannot correct separate ranges.

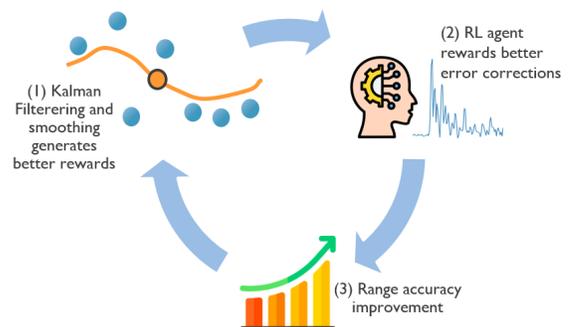


Fig. 1: Conceptual illustration of the idea behind UWB ranging error correction

D. Coppens, B. Van Herbruggen, A. Shahid and E. De Poorter are with ID-Lab, Department of Information Technology, Ghent University—imec, 9052 Ghent, Belgium (e-mail: dieter.coppens@ugent.be; adnan.shahid@ugent.be)

To address these shortcomings, we propose a novel approach based on deep reinforcement learning (RL) which relies on using iteratively improving information automatically derived, removing the need for exact labels. We assume occasional movements of people or vehicles in the environment, which follow sufficiently predictable trajectories. By combining this predictability with filtering, smoothing, and error correction, improvements in error correction are rewarded over time. This iterative process continually enhances the filtered and corrected positions, leading to continuously improving the available information for ranging correction. This concept is illustrated in Figure 1. Finally, the filtering and smoothing can be removed to provide real-time range error correction.

The main contributions of this paper are:

- Introduction of the first self-supervised deep RL framework for CIR-based UWB ranging error correction in a two-way ranging (TWR) system.
- The self-supervised nature of this framework eliminates the requirement for data collection or reliance on ground truth for successful implementation.
- Analyzing the performance of our self-supervised deep RL framework compared to a state-of-the-art supervised convolutional neural network (CNN)

The remainder of the paper is organized as follows. Section II discusses the related work for UWB range error correction. In Section III, the environment in which the dataset is gathered and how the measurements are performed is described. Next, in Section IV the UWB ranging error system model and problem are described. Section V describes the proposed RL methodology and Section VI discusses the performance of the proposed algorithms. The future work follows this in Section VII and finally the conclusion in Section VIII.

## II. RELATED WORK

In this section, an overview of related papers for UWB range error correction in the literature is provided. The related work is split up into four categories, (1) supervised learning, (2) semi-supervised learning, (3) transfer learning, and (4) self-supervised learning.

### A. Supervised learning

The authors of [7] propose a feature-based approach with both support vector machine (SVM) regression and a Gaussian process (GP) to form an estimate of the ranging error. The authors of [12] propose an approach using latent variables that encapsulate information from the CIR about both distance and environmental features to then employ variational inference techniques with neural networks to perform approximate inference in a supervised manner. The authors of [10] propose a supervised deep learning approach for UWB ranging error correction. It leverages a probabilistic deep learning architecture by combining variational inference with probabilistic neural networks. The approach uses a variational autoencoder to learn features from the CIR. [11] uses a similar autoencoder approach for feature extraction from the CIR, but the models are trained in a dual-loss fashion to jointly optimize unsupervised autoencoding and supervised prediction. While

both [11] and [10] leverage unsupervised pre-training of the autoencoder layers, the key ranging error prediction task is formulated as a supervised learning problem. Here, labeled data is used to train a model to directly map inputs to known target outputs. These papers show that supervised machine learning approaches using both raw physical data (CIR) or features can be used to significantly improve the UWB ranging performance. However, none of them address the problem of data collection or the generalization problem, meaning that real-world usability is limited.

### B. Semi-supervised learning

The authors of [13] propose a semi-supervised approach for UWB-ranging error mitigation. Similar to [12] it formulates the problem with a latent variable that encapsulates information about both ranging error and environment. It utilizes a loss function composed of supervised and unsupervised terms, meaning it can use information from both labeled and unlabeled data. This paper addressed the data collection problem and partly succeeded by using semi-supervised learning, meaning that less labeled data is necessary. However, it is not complete unsupervised learning and still requires some supervising (data labeling).

### C. Transfer learning

To address the generalization problem, the authors of [14] propose a transfer learning (TL) framework for UWB error correction using feature- and raw CIR-based approaches. The framework allows for automatic optimizations for TL deep learning models towards new environments while keeping the number of labeled training samples small. The authors demonstrated high accuracy improvements (643 mm to 245 mm) with minimal data collection in challenging environments. The authors of [15] propose an unsupervised TL method based on domain adversarial training and adaptive encoder-decoders. Domain adversarial training is applied to reduce the distribution mismatch between source and target environments. The method still requires labeled data for training the source model. Transfer learning addresses the generalization problem, but still requires data collection for the pre-trained model and/or (minimal) data collection for transferring the knowledge to a new environment.

### D. Self-supervised learning

To the best of our knowledge, [16] is the only self-supervised approach for UWB error mitigation. The authors propose a deep location and ranging correction (DLRC) network to jointly estimate the tag position corrections and distance corrections in a UWB localization system using TDoA. The approach tries to extract the high-level spatial features from ten consecutive ranging measurements received at all anchors in the system, by using the topological information of the UWB system, location loss, and ranging loss are minimized together, allowing training without the ground truth. This approach varies significantly from ours as it was developed for a TDoA system, utilizing measurements from various anchors,

TABLE I: Comparison of our proposed UWB ranging error correction approach with related work. The table mentions the learning method and inputs for learning that are used.

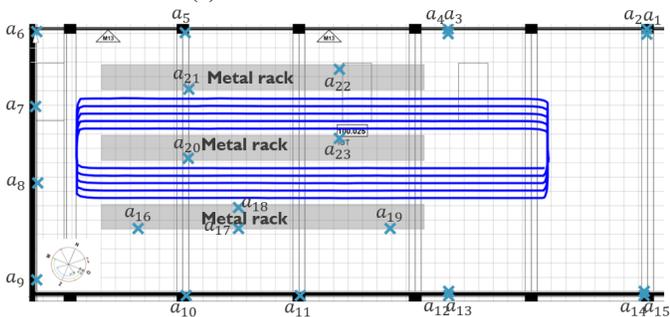
Paper	Self-supervised	ML approach	Localization technique	Environment type	Input for learning	Final outcome of model
[7]		SVM	TWR	LOS/NLOS	Features, range	Ranging error
[12]		Gaussian process Inter-Instance Variational Auto-Encoder	TWR	LOS/NLOS	CIR	Ranging Error Environment label
[10]		Variational inference Probabilistic learning	TWR	LOS/NLOS	CIR	Ranging error
[11]		Variational Auto-Encoder	TWR	LOS/NLOS	CIR	Ranging error
[13]		Variational Bayesian process	TWR	LOS/NLOS	CIR	Ranging Error Environment label
[14]		Transfer learning	TWR	LOS/NLOS	Features, CIR	Ranging Error LOS/NLOS label
[15]		Transfer learning Domain Adversarial Training	TWR	LOS/NLOS	CIR	Ranging Error Environment label
[16]	✓	CNN	TD0A	LOS	Range	Ranging error Positioning error
<b>Our work</b>	✓	<b>RL</b>	<b>TWR</b>	<b>LOS/NLOS</b>	<b>CIR</b>	<b>Ranging Error</b>

without leveraging the valuable information in the CIR data for error correction. This highlights a research gap where our proposed self-supervised reinforcement learning method specifically targets single-range measurements within a TWR ranging system utilizing CIR data and removes the need for ground truth data collection.

### III. DATASET DESCRIPTION



(a) The IIoT lab environment



(b) Floorplan of the IIoT lab with the position of each anchor indicated as light blue X and the ground truth trajectory of the dataset as a dark blue line.

The dataset is collected in an industrial lab environment, which is part of the Industrial Internet of Things (IIoT) testbed [17] of the IDLab research group at Ghent University. The lab is a 240 m<sup>2</sup> warehouse environment, representative of many Industry 4.0 use cases. The IIoT testbed consists of an open space area and an area with metal racks, leading to line-of-sight (LOS) and NLOS situations, pictured in Figure 2a. The environment is equipped with 18 Qualisys Miquis M3 Motion Capture (MOCAP) cameras, capable of tracking hundreds of passive infrared reflective MOCAP markers with a quantified uncertainty in the millimeter range at speeds up to 340 Hz, enabling accurate ground truth determination for evaluation purposes (not used for training in this research). In addition, the MOCAP system is used in combination with a mobile robotic platform to drive repeatable trajectories through the lab. A total of 23 anchors are distributed over the environment, the placement is illustrated in Figure 2b with the light blue crosses. The dataset was collected using Wi-PoS devices [18] that carry the Qorvo DW1000 UWB transceivers. During measurement, the CIR information used for learning was captured at the anchor nodes. To capture the data a mobile robot drives around the lab at 0.1 m/s, the trajectory of the robot is shown in Figure 2b. This trajectory leads to 3463 UWB ranging samples with the different anchors. The ranging method used in the system is called Asymmetric double-sided TWR (ADS-TWR) [19]

#### A. Data pre-processing

Before we use the CIR as state information in the RL algorithm, proposed in section V, we process the raw CIR data, in a pre-processing phase. The pre-processing of the raw CIR involves three distinct steps. First, the complex-valued IQ-sampled array is converted to an RSSI-sampled array. The RSSI is the absolute value of the complex IQ sample, by representing the real (I) and imaginary (Q) components on a Cartesian coordinate system, the RSSI value can thus be determined using the Euclidean distance from the origin:

$$RSSI = \sqrt{I^2 + Q^2} \quad (1)$$

Second, the RSSI-array is trimmed to 150 samples, 50 samples before and 100 after the estimated first path by the DW1000 using the leading edge algorithm. Lastly, the remaining array is normalized using min-max normalization. This means that the highest value in the array becomes 1 and the lowest 0:

$$CIR_{norm} = \frac{CIR - \min(CIR)}{\max(CIR) - \min(CIR)} \quad (2)$$

The normalization step results in smaller numerical values, which is better for training the RL algorithm because it improves the generalization capabilities. This approach tries to make the algorithm to learn to focus on learning signal-to-noise ratio (SNR) and peak features of the CIR, instead of absolute signal strength features. This is important as these can vary significantly across different settings and environments (for example, the average distance between tag and anchor in the environment or higher transmit power configurations) and may not necessarily indicate larger errors or (N)LOS signal propagation. The pre-processing steps, significantly reduce the complexity of the required models, making it computationally more efficient and faster to train. Additionally, a more focused input can help the model generalize better to new, unseen data, as it emphasizes learning essential features and patterns. While reducing the input size, we focused the data around the first path where most errors occur [14].

#### IV. PROBLEM AND SYSTEM DESCRIPTION

In this paper, the purpose is to correct the ranging measurements between the tag and anchor. For a better understanding of the problem, we first provide an overview of the UWB system.

##### A. UWB Localization System

An UWB IPS provides 3D positions  $(x, y, z)$ , relative to a reference point  $ref = (0, 0, 0)$ , for a tag  $t_l \in \{t_1, t_2, \dots, t_L\}$ , with L the total number of tags. For this, it needs to know the coordinates of the fixed anchors  $a_k \in \{a_1, a_2, \dots, a_K\}$ , with K the total number of anchors. To determine its position  $t_{lp} = (t_{lx}, t_{ly}, t_{lz})$ , the tag  $t_l$  will measure the range (distance) to available anchors  $a_k$  in the system. The ground truth range  $\Delta a_k t_l$  can be expressed as follows:

$$\Delta a_k t_l = \sqrt{(a_{k_x} - t_{l_x})^2 + (a_{k_y} - t_{l_y})^2 + (a_{k_z} - t_{l_z})^2} \quad (3)$$

To find the position of the tag,  $\Delta a_k t_l$  is estimated using time of flight (ToF). In this paper, ADS-TWR is used to estimate the ToF accurately. The ToF can be converted to  $\Delta a_k t_l$  as follows:

$$\Delta a_k t_l = ToF_{a_k t_l} \cdot c \quad (4)$$

Where  $c$  is the speed of light ( $3 \times 10^8$  m/s). The ToF is typically estimated using the CIR which quantifies how the communication channel alters the UWB pulse, encapsulating its delay, amplitude, and phase changes. Using a leading-edge algorithm, as used in the popular DW1000 UWB chip [20], the time when the arriving signal, from the accumulated UWB pulses, first rises above the noise floor is the detected first

path ( $fp'$ ). When there is no obstacle between anchor and tag, so-called LOS conditions, this detected first path is close to the actual first path ( $fp$ ) and the ToF estimation is accurate. However, in real-world conditions, there are multipath effects and NLOS conditions. These two effects degrade first path detection and thus higher inaccuracies in ToF estimation. This effect can be demonstrated as follows using the CIR, logged at the UWB transceiver, for signal propagation between  $a_k$  and  $t_l$ :

$$CIR_{a_k t_l}(t) = \sum_{s=1}^S \alpha_s \delta(t - \tau_s) + n(t) \quad (5)$$

Where  $t$  represents the timestamp for each value within the CIR (one CIR has 1016 complex values, corresponding to  $\pm 10^{-9}s$  each);  $S$  is the number of multipath components;  $\alpha_s$  is the amplitude of the  $s$ -th multipath component;  $\tau_s$  the time delay of the  $s$ -th multipath component;  $\delta$  the Dirac delta function and  $n$  represents the additive white Gaussian noise (AWGN) present in the channel. In NLOS conditions,  $fp$  can be severely attenuated and the calculated ToF becomes inaccurate:

$$\widehat{ToF}_{a_k t_l} = ToF_{a_k t_l} + \tau_{(fp' - fp)} \quad (6)$$

With  $fp'$  the first detected path above the noise floor  $n(fp-1)$  and the real first path  $fp$  not detected or not yet received.  $\tau_{fp' - fp}$  is the time difference between the detected first path and the real first path. The ranging error this causes can be calculated as:

$$e_{a_k t_l} = \tau_{(fp' - fp)} \cdot c \quad (7)$$

The calculated range becomes:

$$\widehat{\Delta a_k t_l} = \Delta a_k t_l + e_{a_k t_l} \quad (8)$$

The goal of the UWB error correction model is to predict  $e_{a_k t_l}$  as accurately as possible, using the CIR as input, without collecting a labeled dataset for the training process. Because at each time step, there is only one range received and, for simplicity, in the remainder of the paper  $a_k$  and  $t_l$  will be omitted.

#### V. PROPOSED METHODOLOGY

In our methodology, we assume occasional movements of people or vehicles in the environment, which follow sufficiently predictable trajectories. Combining this predictability with filtering, smoothing, and error correction, improvements in error correction are rewarded over time. This is achieved using a RL process that continually enhances the filtered and corrected positions, leading to continuously improving data available for ranging correction.

##### A. Reinforcement learning

A RL framework consists of an agent and an environment interacting with each other. Anything in the area around the anchor and tag UWB devices that could affect range estimation is regarded as the environment. At each time  $t$ , the agent observes a state  $S_t$  that represents all relevant available information of the environment and takes action  $A_t$ . Here,

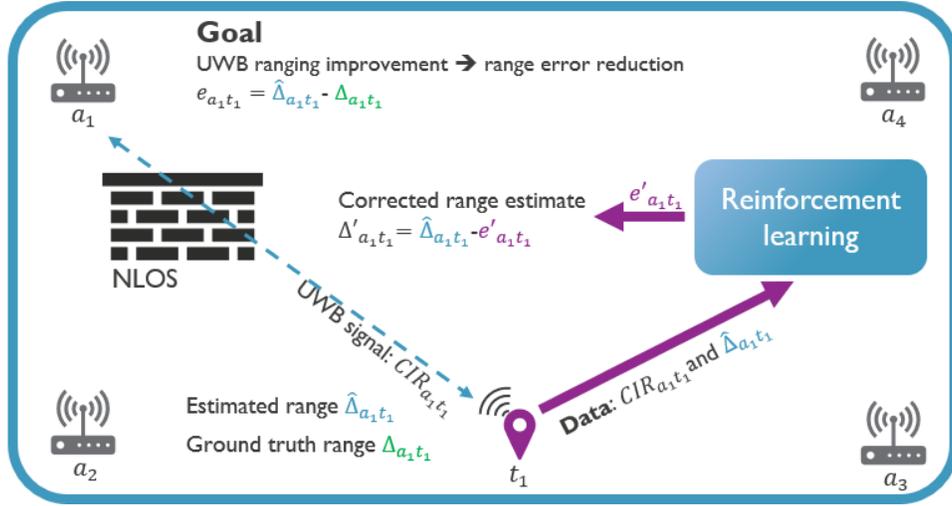


Fig. 3: Illustration of the mathematical UWB localization system description

TABLE II: Mathematical symbols used throughout this article

Symbol	Description
$ref$	Localization system reference point
$t_l$	UWB tag
$a_k$	UWB anchor
$L$	Total number of tags
$K$	Total number of anchors
$\Delta_{a_k t_l}$	Euclidean distance between $a_k$ and $t_l$ , shortened to $\Delta$
$ToF_{a_k t_l}$	Time of flight between $a_k$ and $t_l$ , shortened to $ToF$
$CIR_{a_k t_l}$	Channel impulse response, shortened to $CIR$
$\widehat{ToF}_{a_k t_l}$	Estimated ToF between $a_k$ and $t_l$ , shortened to $\widehat{ToF}$
$e_{a_k t_l}$	Ranging error between $a_k$ and $t_l$ , shortened to $e$
$\hat{\Delta}_{a_k t_l}$	Estimated range between $a_k$ and $t_l$ , shortened to $\hat{\Delta}$
$\hat{e}$	Estimated range error by RL agent
$\Delta'$	Corrected estimated range by RL agent
$S_t$	The state of the environment at time $t$ ( $CIR$ )
$A_t$	The action of the agent at time $t$ ( $\hat{e}$ )
$\pi$	The agent's policy, $\pi : CIR \rightarrow \hat{e}$
$\mu$	The actor network of the RL agent
$\theta$	Weights of the actor network $\mu$
$Q$	The critic network of the RL agent
$\phi$	Weights of the critic network $Q$
$y$	Target Q-value
$R_t$	Reward received at time $t$
$\gamma$	Discount factor, determining weight of target critic
$\hat{Q}$	The target critic network
$\hat{\phi}$	Weights of the target critic network
$\tau_{critic}$	Soft copy factor of the critic
$J$	The sampled policy gradient
$B$	Number of samples in a batch
$\epsilon$	Exploration rate of the RL agent
$\lambda$	Decay factor of the exploration rate
$\hat{\mu}$	The target actor
$\hat{\theta}$	Weights of the target actor
$\tau_{actor}$	Soft copy factor of the actor
$p_{KF}$	Kalman Filter position
$m$	Middle position in smoothing buffer
$N$	Length of circular smoothing buffer
$p_{avg,m}$	Averaged position related to middle data $m$ in buffer
$\Delta_{avg,m}$	Resulting range from filtering and smoothing

at each time step, the UWB localization system estimates the range between a tag and an anchor,  $\hat{\Delta}_t$ . The RL agents corrects the estimate to  $\Delta'_t = \hat{\Delta}_t - \hat{e}_t$ , meaning that  $A_t$  is the error correction  $\hat{e}_t$ .  $S_t$  is the  $CIR_t$  associated with the current range estimation  $\hat{\Delta}_t$ . We assume the UWB ranging error to be between  $\pm 1000$  mm, meaning that the action space  $\mathcal{A}$  can be expressed as:

$$\mathcal{A} = [-1000, 1000] \quad (9)$$

The CIR value received from the DW1000 is pre-processed (described in detail in Section III) to an array of 150 samples with a value between 0 and 1, meaning the state space  $\mathcal{S}$  can be described as:

$$\mathcal{S} = [0, 1]^{150} \quad (10)$$

Due to the continuous action space, we base our custom RL algorithm on the deep deterministic policy gradient (DDPG) algorithm. The behavior of the agent is determined by the policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ . The goal of reinforcement learning is to learn a policy that maximizes the expected rewards. DDPG uses an actor-critic framework, where the policy is determined by the actor network  $\mu(S_t | \theta)$ , with  $\theta$  the weights of the network. The actor network approximates the optimal policy by learning to output the action that maximizes the expected cumulative reward. This expected cumulative reward is determined by an action-value function  $Q(S_t, A_t)$  that is approximated by the critic network  $Q(S_t, A_t | \phi)$ , with  $\phi$  the weights of the network, which takes in state-action pairs and estimates the Q-value. This critic is trained by minimizing the temporal difference between the predicted Q-value and the observed Q-value based on the received rewards  $R_t$ . This is done by minimizing the following loss function, which is an adapted version of the standard DDPG algorithm, as in this problem  $A_t$  does not influence  $S_{t+1}$ :

$$Loss(\phi) = (y_t - Q(S_t, A_t | \phi))^2 \quad (11)$$

With  $y_t$  the target Q-value

$$y_t = R_t + \gamma \hat{Q}(S_t, \mu(S_t | \theta) | \hat{\phi}) \quad (12)$$

This  $y_t$  is dependent on target critic network  $\hat{Q}(S_t, A_t | \hat{\phi})$ , which is a slowly updated version of the main critic by softly copying the weights:  $\hat{\phi} \leftarrow \tau_{critic}\phi + (1 - \tau_{critic})\hat{\phi}$  with  $\tau_{critic} \ll 1$ . This helps stabilize training in DDPG by providing a more consistent target for Q-value predictions. Not using a target critic can lead to increased sensitivity to non-stationary rewards and difficulties in achieving convergence. This would be catastrophic in this research, as the iterative update process causes non-stationary rewards. The actor and critic networks learn collaboratively: the actor network learns to maximize the predicted Q-values by the critic, simultaneously the critic network guides this learning by providing feedback on the quality of the chosen actions in corresponding states. The actor network is updated using a sampled policy gradient to maximize the received expected cumulative rewards:

$$J = -\frac{1}{B} \sum Q(S_t, A_t | \phi) \quad (13)$$

With  $B$  the total number of samples in a batch.

### B. Action selection

At each time step, the agent uses the actor network  $\mu(S_t | \theta)$ , to determine the current best estimate of the correction  $\hat{e}_t$  that will result in the highest reward. However, at the start, the actor network is not well-trained and does not yet know which actions will lead to the best rewards. This leads to two adaptations. First, the actor uses an exploitation/exploration step with an epsilon-greedy policy. In this policy, the correction from  $\mu(S_t | \theta)$  is selected with a probability of  $1-\epsilon$  (exploitation of the actor). With a probability of  $\epsilon$  a random action is chosen uniformly (exploration). The  $\epsilon$  follows an exponential decay during training, at each step:

$$\epsilon = \epsilon_{min} + (\epsilon_{max} - \epsilon_{min}) \cdot e^{-\lambda \cdot step} \quad (14)$$

With  $\epsilon_{min}$  and  $\epsilon_{max}$  the respective minimum and maximum exploration and  $\lambda$  the decay. Exploration is a crucial aspect in RL because it allows the agent to explore which actions lead to good rewards without being constrained by what already has been learned. Second, to avoid bad training data for iterative improvement, we introduce a target actor,  $\hat{\mu}(S_t | \hat{\theta})$  with weights initialized to zero, meaning that the first  $\hat{e}_t$  will also be zero and will not influence the training data. Initializing the weights of a neural network to zero is generally avoided because it leads to a lack of symmetry breaking during training. When all weights are initialized to the same value, neurons in the network will have the same gradients during backpropagation, and they will continue to update in the same way. As a result, the network will fail to learn meaningful representations. However, the  $\hat{\mu}(S_t, \hat{\theta})$  is not intended to be trained on, the weights from  $\mu(S_t, \theta)$  will be "softly" copied to the target network once  $\mu(S_t, \theta)$  has been trained sufficiently to improve the labels rather than deteriorate them. The soft target updates are given by:

$$\hat{\theta} \leftarrow \tau_{actor}\theta + (1 - \tau_{actor})\hat{\theta} \quad (15)$$

With  $\tau_{actor} \ll 1$ . This poses the question of how to define "sufficiently trained". To address this, we employ the "ReduceLRonPlateau" scheduler from PyTorch [21], a dynamic

learning rate adjustment mechanism. The scheduler monitors the loss of the actor, reflecting the quality of the actor's policy, and adjusts the learning rate when a plateau in learning is detected. Once a plateau is identified, indicating that the actor network  $\mu(S_t, \theta)$  has reached a state of sufficient training, the learning rate is reduced. This reduction triggers the soft updating of the target actor network  $\hat{\mu}(S_t, \hat{\theta})$ .

The soft target updates ensure a gradual and controlled transfer of knowledge from the actor network to the target actor network. It is crucial to note that during this soft updating process, the target actor  $\hat{\mu}(S_t, \hat{\theta})$  does not participate in training the actions taken by the actor. Its role is confined to contributing to the data processing pipeline that leads to the calculation of rewards, and maintaining stability in the training process as illustrated in Figure 4. In this way, the dynamic adjustment of the learning rate via "ReduceLRonPlateau" serves as a reliable criterion for defining "sufficiently trained" and triggers the appropriate updates to the target actor network.

### C. Data processing for self-supervised reward

As discussed before, the actor determines ranging correction  $\hat{e}_t$  and the target actor generates correction  $\hat{e}_t$  that leads to the corrected target range  $\hat{\Delta}_t = \hat{\Delta}_t - \hat{e}_t$  and this corrected range estimate is used to iteratively improve the range correction and self-generate better labels.  $\hat{\Delta}_t$  is converted to a position using a Kalman Filter [22].

$$p_{KF,t} = \text{Kalman\_Filter}(\hat{\Delta}_t) \quad (16)$$

Then added to a circular buffer  $C$  of length  $N$  (assumed odd) used for smoothing. If the buffer is full when a new  $p_{KF,t}$  is added to the buffer together with its associated  $\Delta'_t$ ,  $CIR_t$  and  $\hat{e}_t$ , the oldest value in the buffer is removed and the average position of all positions in the buffer is determined and linked to the value at the middle position of the buffer:

$$m = t + \frac{N - 1}{2} \quad (17)$$

And average position:

$$p_{avg,m} = \left( \frac{1}{N} \sum_{i=t}^{t+N-1} x_i, \frac{1}{N} \sum_{i=t}^{t+N-1} y_i \right) \quad (18)$$

This  $p_{avg,m}$  is related to the remaining data at position  $m$  in the circular buffer:  $\Delta'_m$ ,  $CIR_m$  and  $\hat{e}_m$ . Finally, to get an improved range estimate,  $p_{avg,m}$  is converted back to a range:  $\Delta_{avg,m}$ , by calculating the Euclidean distance with the anchor  $a_n$ . This value is our current best estimate of the range and is used in the reward function:

$$R_m = \frac{1}{|\Delta'_m - \Delta_{avg,m}|} \quad (19)$$

The goal of the reward function is to provide a quantitative measure of the success of an agent's actions in the environment. By shaping the reward function appropriately, we can guide the agent to exhibit the desired behavior, namely improved range accuracy. This reward function gives higher rewards the closer the corrected range of the RL agent is to the current best estimate of the range. Updating the neural network at every time step with one sample would be very

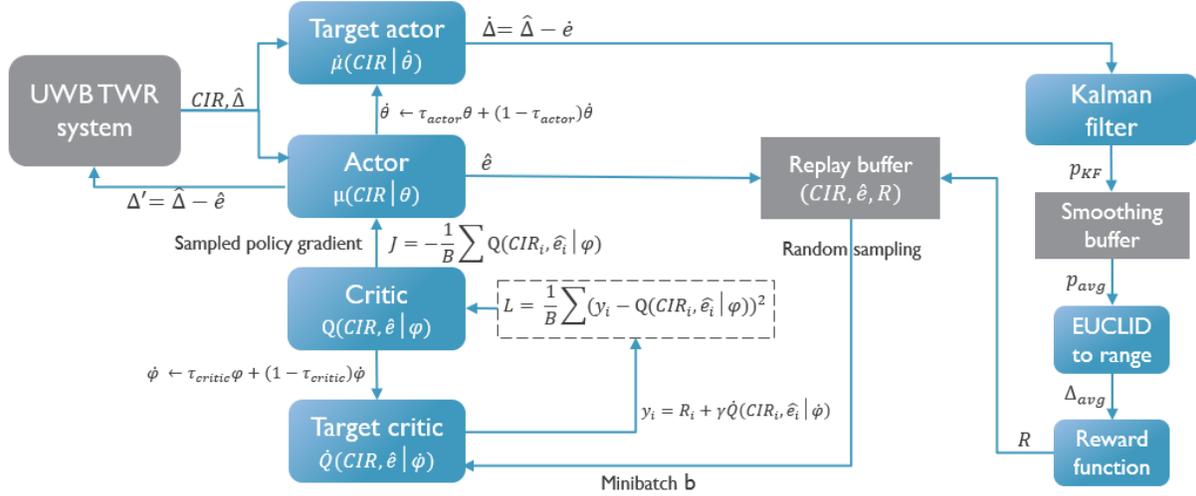


Fig. 4: Complete overview of the proposed (adapted) DDPG algorithm for UWB error correction

inefficient. Therefore, the network is updated on batches of data that are sampled from a replay memory containing experiences  $(CIR_m, \hat{e}_m, R_m)$  generated during the execution of the algorithm. There are several methods to sample from this memory, for this problem we opted for random sampling instead of prioritized sampling as we do not want to overfit on certain experiences or have a lack of diversity in the sampled experiences. An illustration of the complete proposed methodology is shown in Figure 4 and the pseudocode is given in Algorithm 1.

The network architecture of the actor and critic is given in Table III. The actor network ends with scaling the output of a dense layer with a Tanh activation function. The Tanh output is between -1 and 1, which leads to the final output being scaled to -1000 and 1000 which is equal to the action space. The critic network starts similar to the actor network, but the information of the CIR is encapsulated in 4 latent features. These 4 latent features are concatenated with the action selected by the actor. This is then further processed to a final layer with a Tanh activation, meaning that the Q-value is between -1 and 1.

## VI. RESULTS AND ANALYSIS

### A. Baselines and metrics

For performance evaluation, we will use two evaluation metrics: (1) the mean absolute error (MAE) as it encapsulates the performance in a single value and (2) box plots to provide a clear and concise way to see the spread (variability) of ranging errors and thus a more general overview of the performance while also highlighting central tendencies. To evaluate our proposed method, we compare our results against two baselines: the first baseline is the uncorrected UWB performance, and the second is a state-of-the-art supervised CNN-based method [11] trained on the fully labeled dataset. The results of the supervised CNN are not directly adopted from the paper itself, but the developed model has been retrained on the dataset of this research. The range error results of NLOS samples will be shown and discussed separately because of the reduced signal clarity in NLOS situations. For NLOS, the signal propagation

---

### Algorithm 1: Self-supervised RL for error correction

---

**Data:** Initialize replay memory  $D$ ;  
 Initialize circular smoothing buffer  $C$  with length  $N$ ;  
 Initialize actor network  $\mu$  and critic network  $Q$  function with random weights  $\theta$  and  $\phi$ ;  
 Initialize target actor  $\hat{\mu}$  with weights  $\hat{\theta} = 0$ ;  
 Initialize target critic  $\hat{Q}$  with random weights  $\hat{\phi}$

**while** *episode* < *training episodes* **do**

**while** *episode not done* **do**

Get current data  $\hat{\Delta}_t$ ,  $CIR_t$ , and  $\hat{e}_t$ ;  
 With probability  $\epsilon$ , select random correction  $a_t$ ;  
 Otherwise,  $a_t = \hat{e}_t = \hat{\mu}(\hat{\Delta}_t, \hat{\theta})$ ;  
 Correct range estimate  $\hat{\Delta}'_t = \hat{\Delta}_t - \hat{e}_t$ ;  
 Determine  $p_{KF,t} = \text{Kalman\_Filter}(\hat{\Delta}'_t)$ ;  
 Add  $p_{KF,t}$  to circular buffer  $C$ ;  
**if**  $C$  is full **then**

$p_{avg,m} = (\frac{1}{N} \sum_{i=t}^{t+N-1} x_i, \frac{1}{N} \sum_{i=t}^{t+N-1} y_i)$ ;  
 Convert  $p_{avg,m}$  to  $\Delta_{avg,m}$ ;  
 Calculate  $R_m = \frac{1}{|\Delta'_m - \Delta_{avg,m}|}$ ;  
 Store experience  $d_m = (CIR_m, \hat{e}_t, R_m)$  in  $D$ ;

**if** Every  $K$  steps **then**

Sample a random minibatch  $b$  from  $D$ ;  
**foreach**  $d_j$  in  $b$  **do**

$y_j = R_j + \gamma \hat{Q}(CIR_j, \hat{e}_j | \hat{\phi})$ ;  
 Update critic by minimizing the loss:  
 $L = \frac{1}{B} \sum_j (y_j - Q(CIR_j, \hat{e}_j | \phi))^2$ ;  
 Update actor using sampled policy gradient:  
 $J = -\frac{1}{B} \sum_j Q(CIR_j, \hat{e}_j | \phi)$ ;

**if** Every  $T$  steps **then**

Update target critic:  $\hat{\phi} \leftarrow \tau \phi + (1 - \tau) \hat{\phi}$ ;  
**if**  $\mu$  sufficiently trained **then**

Update target actor:  $\hat{\theta} \leftarrow \tau \theta + (1 - \tau) \hat{\theta}$ ;

---

TABLE III: Actor and Critic Network Architectures

Actor network			Critic Network		
Layer	Activation	Output size	Layer	Activation	Output size
Input (State)		(1,150,1)	Input (State)		(1,150,1)
Conv2D(128,16x1)	ReLU	(128,150,1)	Conv2D(128,16x1)	ReLU	(128,150,1)
Maxpool(2x1)		(128,75,1)	Maxpool(2x1)		(128,75,1)
Conv2D(64,8x1)	ReLU	(64, 75, 1)	Conv2D(64,8x1)	ReLU	(64, 75, 1)
Conv2D(32,2x1)	ReLU	(32,75,1)	Conv2D(32,2x1)	ReLU	(32,75,1)
BatchNorm		(32,75,1)	BatchNorm		(32,75,1)
Dropout 25%		(32,75,1)	Dropout 25%		(32,75,1)
Flatten		2400	Flatten		2400
Dense	ReLU	150	Dense	ReLU	150
BatchNorm		150	BatchNorm		150
Dropout 20%		150	Dropout 20%		150
Dense	ReLU	100	Dense	ReLU	100
Dropout 20%		100	Dropout 20%		100
Dense	ReLU	50	Dense	ReLU	50
Dropout 10%		50	Dropout 10%		50
Dense	Sigmoid	25	Dense	Sigmoid	25
Dense	Tanh	1	Dense	ReLU	4
Output Scaling (x1000)		1	Concat (add action)		5
			Dense	ReLU	8
			Dense	ReLU	16
			Dense	ReLU	8
			Dense	Tanh	1

between transmitter and receiver is more complex due to the attenuated first path signal power, leading to a more complicated relationship between CIR and error correction. The NLOS situations are the most vital for error correction, as they are prone to the largest ranging errors. Showing the performance in NLOS situations separately provides insight into how the baselines and our proposed RL algorithm perform in the most challenging conditions.

### B. Training

The RL algorithm was trained for 1000 episodes with  $\gamma = 0.5$ ,  $\tau_{critic} = \tau_{actor} = 0.01$ ,  $\alpha_{actor} = 5e^{-5}$  and  $\alpha_{critic} = 5e^{-4}$ . The patience of the learning rate schedulers was set to 150 episodes.

Figure 5 illustrates the learning curve of the algorithm. During the first 100 episodes, there is a steep decrease in MAE and thus a quickly improving performance. Between episodes 100 and 350, the decrease starts slowing down, which leads the scheduler to reduce the learning rates. The reduced learning rate is visible in Figure 5 from episode 350 onwards. This early learning phase is primarily shaped by the exploration-exploitation trade-off, the exploration is decaying exponentially. Lower and faster decaying exploration would cause an even more steep decrease in MAE, but could come at the cost of worse final performance as the algorithm explores fewer possibilities. Higher and slower declining exploration would come at the cost of slower convergence and more training episodes needed. In the figure, the reduced fluctuations in performance, from episode 350 onward, show the reduced learning rate. At the end of the training, the MAE of the RL algorithm is distinctly lower than the uncorrected UWB ranging and the supervised CNN approach. Figure 6 illustrates the iterative improvement by the RL algorithm during training. The green curve represents the smoothed KF trajectory without any RL correction, this is the data used to calculate the reward before the target actor is updated. Once the target actor starts

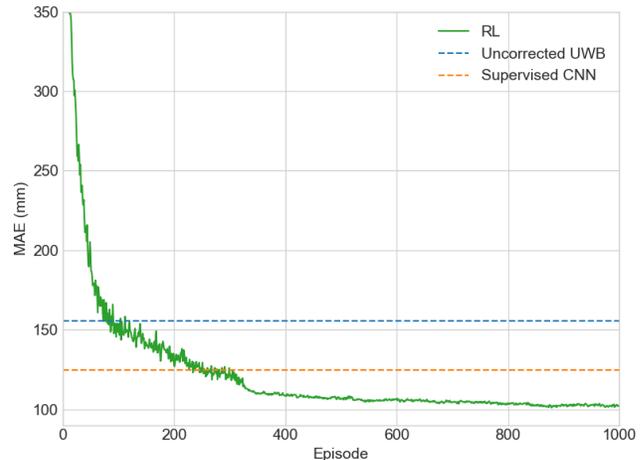


Fig. 5: Performance comparison of our proposed RL algorithm during training with uncorrected UWB ranging and a supervised CNN approach in terms of MAE. The figure shows that our proposed algorithm quickly improves the ranging performance compared to uncorrected UWB ranging, and later surpasses the supervised CNN performance.

getting updates from episode 350 onwards, the data used to calculate the rewards starts improving. The orange trajectory is used for reward calculation at episode 400 and the green trajectory is the improved trajectory at episode 1000, the end of training. This visually illustrates the iterative improvement of the algorithm.

### C. Evaluation

In Figure 7, the box plots show that the proposed RL algorithm performs better than the supervised CNN approach and significantly better than the uncorrected UWB. The median range error of the RL algorithm is lower than the other two methods for both box plots.

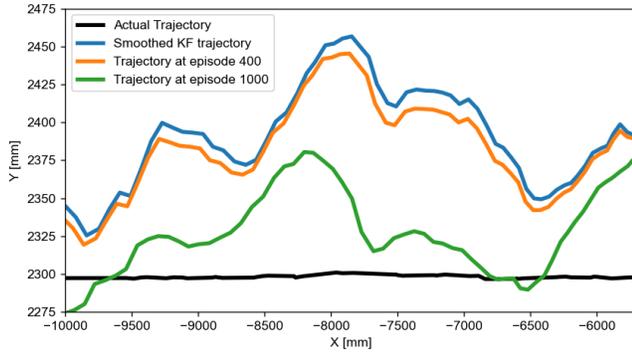


Fig. 6: Trajectory comparison of the original smoothed KF (with no RL correction) with the improved trajectories during training after 400 and 1000 episodes

A first explanation is that **our approach is rewarded when the overall trajectory improves. As such, our approach focuses on those CIR types that most negatively impact the overall trajectory. In contrast, traditional CNNs optimize all CIRs equally, regardless of their importance to the overall localization.** A second explanation relates to an improved design of the (actor) neural network compared to prior work. Our last layer has Tanh activation with output scaling while the prior CNN uses no activation in the last layer so that the output could be any real number which might lead to numerical instability in some cases.

Moreover, the ability to learn online from new experiences without supervision means the RL agent can seamlessly adapt to changes in environments by continuing the training process, shown in Section VI-D, the RL approach can learn the specifics of an environment even if the environment changes. Additionally, the range errors of the RL algorithm are more tightly clustered around the median. This indicates that the RL algorithm is more robust in its performance than the other two methods. Figure 7b highlights the performance of the proposed RL algorithm in NLOS situations, with a lower median error rate and smaller interquartile range, this again indicates a more consistent performance in the more difficult situations. However, it has visibly more outliers than the supervised CNN approach. Table IV tells a similar story. The proposed RL approach reduces MAE by 31.6% compared to uncorrected UWB and by 14.5% compared to the supervised CNN. The separate NLOS results emphasize the increased performance even more, as the proposed algorithm decreases the MAE by 34.8% compared to uncorrected UWB and by 22.8% compared to the supervised CNN.

TABLE IV: Quantitative results of the baselines and proposed algorithm

Method	MAE (mm) all samples	MAE (mm) NLOS
Uncorrected UWB	155	181
Supervised CNN [11]	124	153
Proposed RL algorithm	106	118

#### D. Adaptivity of algorithm

The same environment can change over time. To evaluate the performance of our algorithm when there are sudden changes in the environment, a new dataset was collected in the same warehouse 6 months later. At that time, there were more goods in the racks, additional clutter in the warehouse (obstacles, boxes, ...) and the anchor nodes experienced many small disturbances over time. These combined effects lead to a more challenging environment, which is reflected in the higher MAE for uncorrected UWB. The MAE during training is shown in Figure 8. First, the proposed algorithm is trained on "environment 1", which is the original dataset used in the previous evaluation, between episodes 0-500, the training is similar to Figure 5, except it is halted after 500 instead of 1000 episodes. After 500 episodes, the environment is switched to "environment 2", the same environment but 6 months later and more difficult as discussed before. The learning rates of the RL algorithm are reset to the starting values and the exploration is increased. At first, the RL algorithm leads to worse performance, due to the exploration, but quickly adapts to the environment and surpasses the supervised model trained on the first dataset and later also the supervised CNN trained on the new dataset. This result displays the adaptivity of the RL algorithm compared to the supervised CNN approach, without needing to label a dataset it can adapt to changing environments and continuously lead to better ranging performance. The supervised CNN approach requires a new labeled dataset in a changed environment, while our approach does not.

#### E. Complexity analysis

1) *Algorithmic complexity*: The proposed RL algorithm leverages deep neural networks to approximate the actor and critic functions. Therefore, it is important to analyze their complexities. The networks can be broken down in components different components. First, convolutional layers that can be calculated as follows:  $O(H * W * C_{in} * C_{out} * K_w * K_h)$  with  $H * W$  the input size,  $C_{in}$  the input channels,  $C_{out}$  the output channels and  $K_w * K_h$  the kernel size. Following Table III, this results in:

- Conv1:  $O(1 * 150 * 128 * 16) = O(3.07 * 10^5)$
- Conv2:  $O(75 * 128 * 64 * 8) = O(4.92 * 10^6)$
- Conv3:  $O(75 * 64 * 32 * 2) = O(3.07 * 10^5)$

Showing that the complexity is dominated by the second convolutional layer. The complexity of a linear layer is given by  $O(nm)$  with  $n$  the number of input features and  $m$  the number of output features. The first dense layer will be the largest and given by  $O(3.60 * 10^5)$ . Following linear layers become less and less complex as input and output sizes decrease.

2) *Time complexity*: The algorithm's convergence time duration is limited by the incoming data rate of the UWB localization system. The UWB localization has a sampling rate of 50 Hz, because the dataset has about 3000 samples the time duration of 1 episode is about one minute.

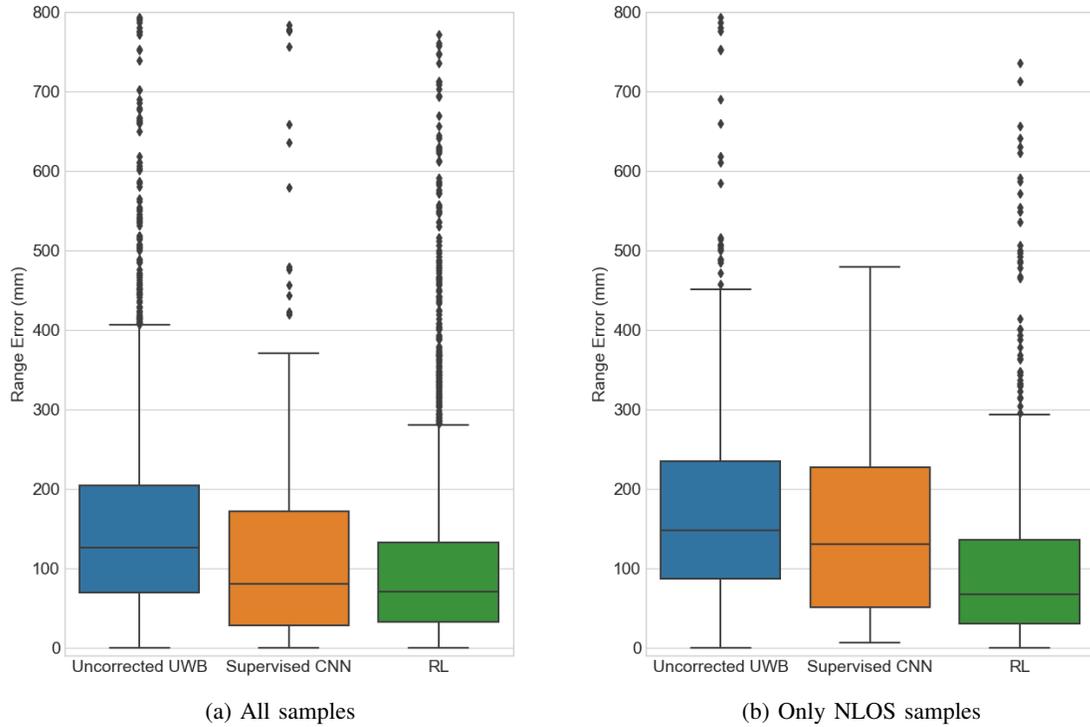


Fig. 7: The ranging errors of uncorrected UWB, the supervised CNN and our proposed RL algorithm during evaluation for (a) all samples and (b) only in NLOS samples. The figures show that our proposed self-supervised RL algorithm performs comparable or better than a supervised CNN approach.

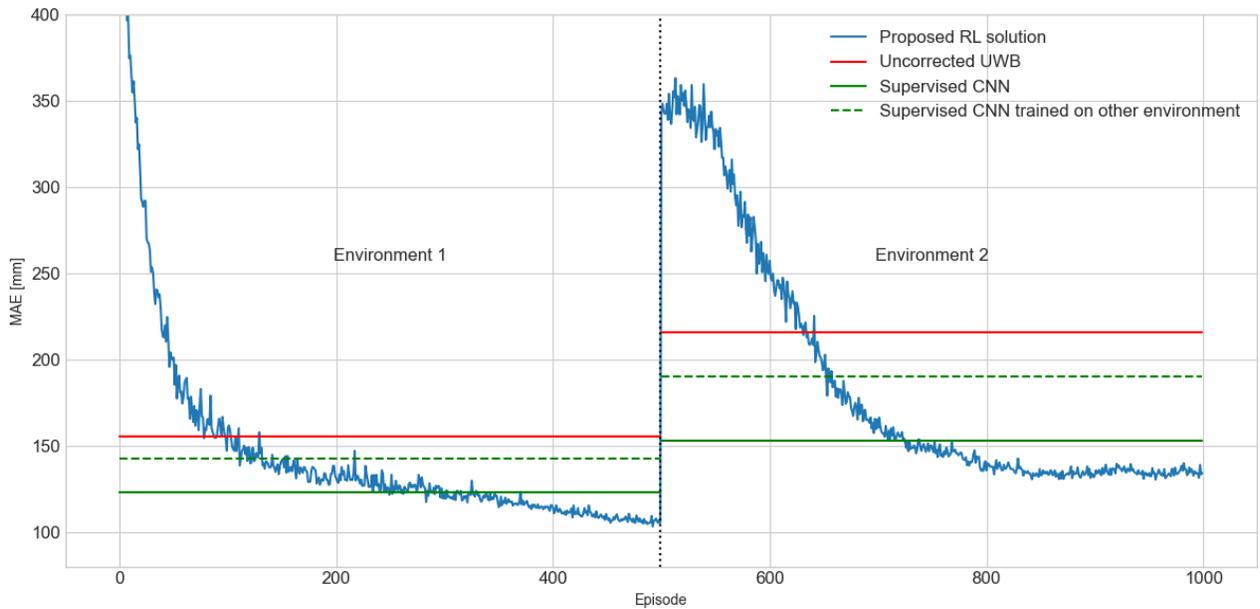


Fig. 8: MAE comparison of our approach with a state-of-the-art CNN for error correction [1] for a changing environment. The proposed RL algorithm is compared with uncorrected UWB and the supervised CNN trained on the current environment and the other environment. The figure shows the generalization problem of the supervised CNN and that the proposed RL algorithm can adapt itself to a changing environment.

## VII. FUTURE WORK

There are several avenues to further expand on this research. A first potential enhancement could be to make the Kalman Filter adaptive to the trajectory. Straight trajectories can have more smoothening, while corners need reduced smoothening. By modifying the smoothening and filtering process based on trajectory characteristics, the tracking system's accuracy and robustness can be further investigated. A second area of improvement lies in the selection of points from the smoothed and filtered trajectory to the discrete data points for learning. Currently, the middle point in the buffer is associated with the average position. Future research could investigate the feasibility of defining a continuous trajectory and selecting the closest point. This adjustment could potentially lead to more responsive and accurate error correction processes. Additionally, the system's capabilities could be expanded by integrating various sources of additional information. This includes exploring adding map data, reflections, CIR, and range data between anchor nodes (with known fixed positions). Furthermore, Inertial Measurement Unit (IMU) data could be added as input to the Kalman Filter to make it more robust and allow for better labels. Finally, applying this research to positioning systems using TDoA methods instead of TWR systems. TDoA, which relies on measuring the time delays of signals arriving at different nodes, is a widely used technique in wireless localization, and applying this research there would further improve the real-world practicality of more positioning systems.

## VIII. CONCLUSION

In this work, we propose a novel self-supervised deep reinforcement learning approach for Ultra-Wideband ranging error correction that does not require ground truth data. This is significant because collecting large labeled datasets for model training is impractical for real-world indoor positioning system deployment. The methodology is based on the assumption that there are occasional movements of people or vehicles in the environment, following sufficiently predictable trajectories. Experiments on real-world measurements demonstrate our approach achieves comparable or improved ranging accuracy compared to a state-of-the-art CNN approach for error correction. Specifically, our method reduces errors by up to 31.6% compared to uncorrected UWB in challenging situations without any data labeling. Additionally, the reinforcement learning agent can quickly adapt to changing environments. This makes our self-supervised framework highly practical for use in real indoor scenarios, as it removes the dependency on time-consuming and costly ground truth collection efforts. In summary, by not relying on labeled data, our approach paves the way for more scalable and generalized Ultra-Wideband error mitigation solutions using deep reinforcement learning that can be easily deployed in various indoor spaces.

## REFERENCES

- [1] R. Bazo, C. A. da Costa, L. A. Seewald, L. G. da Silveira, R. S. Antunes, R. d. R. Righi, and V. F. Rodrigues, "A survey about real-time location systems in healthcare environments," *Journal of Medical Systems*, vol. 45, pp. 1–13, 2021.
- [2] K. Minne, N. Macoir, J. Rossey, Q. Van den Brande, S. Lemey, J. Hoebeke, and E. De Poorter, "Experimental evaluation of uwb indoor positioning for indoor track cycling," *Sensors*, vol. 19, no. 9, p. 2041, 2019.
- [3] M. Elsanhoury, P. Mäkelä, J. Koljonen, P. Välisuo, A. Shamsuzzoha, T. Mantere, M. Elmusrati, and H. Kuusniemi, "Precision positioning for smart logistics using ultra-wideband technology-based indoor navigation: A review," *IEEE Access*, vol. 10, pp. 44 413–44 445, 2022.
- [4] H. Huang, G. Gartner, J. M. Krisp, M. Raubal, and N. Van de Weghe, "Location based services: ongoing evolution and research agenda," *Journal of Location Based Services*, vol. 12, no. 2, pp. 63–93, 2018.
- [5] D. Coppens, A. Shahid, S. Lemey, B. Van Herbruggen, C. Marshall, and E. De Poorter, "An overview of uwb standards and organizations (ieee 802.15. 4, fira, apple): Interoperability aspects and future research directions," *IEEE Access*, vol. 10, pp. 70 219–70 241, 2022.
- [6] A. Alarifi, A. Al-Salman, M. Alsaleh, A. Alnafessah, S. Al-Hadhrani, M. A. Al-Ammar, and H. S. Al-Khalifa, "Ultra wideband indoor positioning technologies: Analysis and recent advances," *Sensors*, vol. 16, no. 5, p. 707, 2016.
- [7] H. Wymeersch, S. Marano, W. M. Gifford, and M. Z. Win, "A machine learning approach to ranging error mitigation for uwb localization," *IEEE transactions on communications*, vol. 60, no. 6, pp. 1719–1728, 2012.
- [8] W. M. Gifford, D. Dardari, and M. Z. Win, "The impact of multipath information on time-of-arrival estimation," *IEEE Transactions on Signal Processing*, vol. 70, pp. 31–46, 2020.
- [9] B. Denis, J. Keignart, and N. Daniele, "Impact of nlos propagation upon ranging precision in uwb systems," in *IEEE conference on Ultra Wideband Systems and Technologies, 2003*. IEEE, 2003, pp. 379–383.
- [10] C. Mao, K. Lin, T. Yu, and Y. Shen, "A probabilistic learning approach to uwb ranging error mitigation," in *2018 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2018, pp. 1–6.
- [11] J. Fontaine, M. Ridolfi, B. Van Herbruggen, A. Shahid, and E. De Poorter, "Edge inference for uwb ranging error correction using autoencoders," *IEEE Access*, vol. 8, pp. 139 143–139 155, 2020.
- [12] Y. Li, S. Mazuelas, and Y. Shen, "A variational learning approach for concurrent distance estimation and environmental identification," *IEEE Transactions on Wireless Communications*, 2023.
- [13] —, "A semi-supervised learning approach for ranging error mitigation based on uwb waveform," in *MILCOM 2021-2021 IEEE Military Communications Conference (MILCOM)*. IEEE, 2021, pp. 533–537.
- [14] J. Fontaine, F. Che, A. Shahid, B. Van Herbruggen, Q. Z. Ahmed, W. B. Abbas, and E. De Poorter, "Transfer learning for uwb error correction and (n) los classification in multiple environments," *IEEE Internet of Things Journal*, 2023.
- [15] Z. Li, K. Hu, T. Wang, S. Cui, and Y. Shen, "An unsupervised transfer learning method for uwb ranging error mitigation," *IEEE Communications Letters*, 2023.
- [16] B. Yang, J. Li, Z. Shao, and H. Zhang, "Self-supervised deep location and ranging error correction for uwb localization," *IEEE Sensors Journal*, vol. 23, no. 9, pp. 9549–9559, 2023.
- [17] [Online]. Available: <https://www.ugent.be/ea/idlab/en/research/research-infrastructure/industrial-iot-lab.htm/>
- [18] B. Van Herbruggen, B. Jooris, J. Rossey, M. Ridolfi, N. Macoir, Q. Van den Brande, S. Lemey, and E. De Poorter, "Wi-pos: A low-cost, open source ultra-wideband (uwb) hardware platform with long range sub-ghz backbone," *Sensors*, vol. 19, no. 7, p. 1548, 2019.
- [19] Y. Jiang and V. C. Leung, "An asymmetric double sided two-way ranging for crystal offset," in *2007 International Symposium on Signals, Systems and Electronics*, 2007, pp. 525–528.
- [20] "Decawave - dw1000 ic," Accessed March, 2024. [Online]. Available: <https://www.decawave.com/product/dw1000-radio-ic/>.
- [21] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "PyTorch: An Imperative Style, High-Performance Deep Learning Library," in *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc., 2019, pp. 8024–8035. [Online]. Available: <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
- [22] G. Mao, S. Drake, and B. D. Anderson, "Design of an extended kalman filter for uav localization," in *2007 Information, Decision and Control*. IEEE, 2007, pp. 224–229.