

JUNO: Jump-Start Reinforcement Learning-based Node Selection for UWB Indoor Localization

1st Zohreh Hajiakhondi-Meybodi

Electrical & Computer Engineering,
Concordia University,
Montreal, Canada

z_hajiak@encs.concordia.ca

2nd Ming Hou

Defence Research and Development
Canada (DRDC),
Toronto, Canada

ming.Hou@drdc-rddc.gc.ca

3th Arash Mohammadi

Concordia Institute for Inf. Systems Eng.,
Concordia University,
Montreal, Canada

arash.mohammadi@concordia.ca

Abstract—Ultra-Wideband (UWB) is one of the key technologies empowering the Internet of Thing (IoT) concept to perform reliable, energy-efficient, and highly accurate monitoring, screening, and localization in indoor environments. Performance of UWB-based localization systems, however, can significantly degrade because of Non Line of Sight (NLoS) connections between a mobile user and UWB beacons. To mitigate the destructive effects of NLoS connections, we target development of a Reinforcement Learning (RL) anchor selection framework that can efficiently cope with the dynamic nature of indoor environments. Existing RL models in this context, however, lack the ability to generalize well to be used in a new setting. Moreover, it takes a long time for the conventional RL models to reach the optimal policy. To tackle these challenges, we propose the Jump-start RL-based Uwb NLoS selection (JUNO) framework, which performs real-time location predictions without relying on complex NLoS identification/mitigation methods. The effectiveness of the proposed JUNO framework is evaluated in terms of the location error, where the mobile user moves randomly through an ultra-dense indoor environment with a high chance of establishing NLoS connections. Simulation results corroborate the effectiveness of the proposed framework in comparison to its state-of-the-art counterparts.

Index Terms—Indoor Localization, Internet of Things, Anchor Selection, Jump-Start Reinforcement Learning, Ultra-Wideband (UWB).

I. INTRODUCTION

Ultra-WideBand (UWB) technology has been emerged as a solution to meet the phenomenal growth of the need for localizing users in indoor environments [1]. The use of a wide radio spectrum in UWB technologies enables individual multi-path components of the received signal to be efficiently resolved, resulting in high accuracy indoor positioning [2]. To monitor/track users in indoor environments, several localization techniques have been proposed such as Received Signal Strength Indicator (RSSI) [3], [4], Angle of Arrival (AoA) [5]–[7], and Time Difference of Arrival (TDoA) [8], [9], among which time-based solutions [8]–[10] can be considered as more efficient ones. In such scenarios, the time of the received signal from a set of available UWB beacons is required, which can only be estimated accurately if the first arrival path has been properly identified. In this context, one key

challenge is susceptibility to the Non-Line-of-Sight (NLoS) error. In presence of an obstacle between the UWB beacon and the mobile device, the time of the received signal will be delayed, resulting in a positive bias and a significant degradation in the positioning accuracy. Therefore, NLoS mitigation/identification in UWB-based indoor localization is of paramount importance. Conventional indoor localization frameworks reduced such location error via parametric solutions [11], [12], the accuracy of which is dependent on implementation of complex pre-processing techniques adding considerable latency. Furthermore, using a large number of UWB beacons for localizing users is inefficient from the energy consumption perspective. In this regard, Dai *et al.* [13] analytically proved that tracking users' locations through a subset of active beacons offers several benefits, including mitigating the energy consumption of beacons, and improving the location accuracy. Consequently, the main focus of recent researches [14]–[17] has been shifted to use anchor node selection to achieve the best localization performance in terms of location accuracy and resource management. The paper aims to further advance this emerging field.

Related Work: Anchor node selection in the context of indoor localization is utilized to improve the network's performance by setting a set of criteria for selecting a subset of beacons with the highest utilities. One of the most important criteria in indoor localization is to mitigate the location error, caused by NLoS connections. Towards this goal, LoS connections will be selected for monitoring/tracking users' locations instead of extracting location information from all available beacons. Generally speaking, anchor selection frameworks can be classified into two groups, i.e., analytical solutions and Artificial Intelligence (AI) models. Conventionally, the focus of anchor selection frameworks was on the former category, i.e., deriving fixed mathematical/optimization models [15], [17], [20] to meet the acceptable location accuracy. To alleviate the high complexity caused by using complicated mathematical formulations, recent research works [21]–[24] used AI and Machine Learning (ML) models to localize/navigate mobile devices in indoor environments. One of the most commonly used AI-based LoS/NLoS identification approaches is supervised learning models [21], [22]. For instance, Poulouze *et al.* [22]

applied Long Short-Term Memory (LSTM) network to predict the location of mobile devices using the Time of Arrival (ToA)-based UWB method. Despite all the benefits that come from using supervised models, there are several key challenges ahead. On the one hand, supervised models require labelled LoS/NLoS data, which is both costly and time-consuming, limiting general applicability of supervised models within this context. On the other hand, even minor changes in the indoor environment would require updating the training dataset. Unsupervised learning models [24], however, eliminate the necessity for labelling of channel conditions, therefore, allowing generalization while saving time and precious resources. Such models, however, are impractical due to the dynamic nature of indoor environments, such as unknown/varying number of users and adverse environmental conditions. Furthermore, energy consumption efficiency is compromised in scenarios where all the beacons require to transmit/receive signals for LoS/NLoS identification.

To tackle the above mentioned issues, the main focus of researchers has been shifted to use Reinforcement Learning (RL) models [25]–[28] for indoor localization/navigation application. In our previous work [25], for instance, the mobile user is autonomously trained via an RL model to be localized by a set of UWB beacons with LoS connections. The main objective of the RL model within the indoor localization domain is to learn an optimal or near-optimal policy that maximizes the location accuracy. One of the most important challenges of existing RL models [25], however, is that the optimal policy should be learned by the interaction of the agent (mobile user) with the environment (i.e., via trial and error), without any prior information, especially when the model is just initialized. Consequently, it may take a long time for the RL model to reach the optimal policy. Another challenge is the generalization ability of the pre-trained RL model to be used in a new and different environment, where the density/location of obstacles is changing over the time/environment. To tackle these issues, Uchendu *et al.* [29] proposed the Jump-Start RL (JSRL) model, where the agent use a guide-policy instead of a random one at the beginning of the learning process. Consequently, the learning process is accelerated and the RL generalization ability is highly improved.

Contribution: Motivated by the above discussion, we introduce the Jump-start RL-based Uwb NODe selection (JUNO) framework with the application to indoor localization. The main novelty of this work is the design of an autonomous and real-time anchor node selection, where the key objective is to accelerate the location accuracy improvement. Towards this goal, a combination of the guide and exploration policies is used, where the guide-policy significantly speeds up the early learning phase of the RL model to converge to the optimal location accuracy. Furthermore, since any random guide-policy can be used in the JUNO framework, the generalization ability of the pre-trained RL frameworks improves. Simulation results illustrate that the proposed JUNO framework outperforms its state-of-the-art counterparts in terms of the cumulative rewards

and location error even in an ultra-dense indoor environment.

The rest of this paper is organized as follows: In Section II, the system model is provided. Section III introduces the proposed JUNO framework. Section IV presents experimental results. Finally, Section V concludes the paper.

II. SYSTEM MODEL

In this section, we first introduce the UWB wireless signal model and the TDoA localization formulation. Then, we present the required background on the RL model.

A. UWB Wireless Signal Model

In this study, we consider a multi-user indoor environment (e.g., an office or a hotel building), consisting of N synchronized UWB beacons, denoted by UWB_i , $i = 1, \dots, N$, and several mobile users randomly moving through the environment. To support multiple users, we use the Time Hopping (TH) technique as one of the efficient Code Division Multiple Access (CDMA) schemes, where different codes are assigned to distinct users. Given the Pulse Amplitude Modulation (PAM) modulation, the transmitted signal s_u from user u is obtained as

$$s_u(t) = \sum_{n=-\infty}^{\infty} p_u(n) \sum_{s=0}^{N_c-1} c_u(s)w(t - nT_s - sT_c - \theta_u) \quad (1)$$

where T_s , and N_c represent the symbol time, and the number of chips with duration of T_c in each symbol, respectively. Term $c_u(s) \in \{0, 1\}$ denotes the access code associated with the mobile user u , and $p_u(n) \in \{-1, 1\}$ is the Independent and Identically Distributed (IID) information symbols. Furthermore, term θ_u is the time asynchronism with uniform distribution $[0, T_s]$, and $w(t)$ is the normalized impulse signal. The received signal $r_i(t)$ by UWB_i is expressed as

$$r_i(t) = \sum_{u=1}^{N_u} \sqrt{P_u} \sum_{k=1}^{N(t)} \rho_{u,k}(t, \tau) s_u(t - \tau_{u,k}(t) - \tau_i) + n(t), \quad (2)$$

where $\tau_{u,k}(t)$ and $\rho_{u,k}(t, \tau) = \beta_{u,k}(t, \tau) \exp(j\Phi_{u,k}(t, \tau))$ represent the phase delay and the attenuation associated with k^{th} path, where the total number of detachable paths is denoted by $N(t)$, and N_u is the number of users in the experimental indoor environment. Terms $\beta_{u,k}(t, \tau)$ and $\Phi_{u,k}(t, \tau)$ are the amplitude and phase of k^{th} path, respectively. Term $\beta_{u,k}(t, \tau)$ is a Nakagami- m random variable that reflects the LoS/NLoS link condition, with $m = 1$ representing the Rayleigh fading and $m > 1$ illustrating the Rician channel model. Term $n(t) \sim \mathcal{N}(0, \sigma^2)$ is the Additive White Gaussian Noise (AWGN) channel. Moreover, term $\tau_i = d_i/c$ is the time of arrival of signal $s_u(t)$, which is equivalent to the delay of the first path within the LoS condition. Finally, terms d_i and $c = 3 \times 10^8$ m/s represent the distance between the user u and UWB_i , and the speed of light, respectively. Given the first peak of the estimated Channel Impulse Response (CIR) of at least two UWB beacons UWB_i and UWB_j , denoted by

τ_i and τ_j , the mobile user's location at time slot t , denoted by (x_t, y_t) , is calculated as follows

$$\tau_i - \tau_j = \frac{\sqrt{(x_t - x_i)^2 + (y_t - y_i)^2} - \sqrt{(x_t - x_j)^2 + (y_t - y_j)^2}}{c}, \quad (3)$$

where (x_i, y_i) and (x_j, y_j) denote the locations of UWB_i and UWB_j , respectively. This completes presentation of the UWB wireless signal model, next, we introduce the RL background.

B. RL Background

RL model is an area of learning paradigm, where an agent, interacting with the environment, makes a sequence of decisions to achieve the maximum accumulated rewards. Markov Decision Process (MDP) is used to mathematically describe an RL environment, which includes a set of states \mathcal{S} , a set of actions \mathcal{A} , a reward function \mathcal{R} , and a transition function \mathcal{T} . Based on the current state $s_t \in \mathcal{S}$ at time slot t , the agent takes an action $a_t \in \mathcal{A}$, resulting in a new state $s_{t+1} \in \mathcal{S}$ at time slot $t+1$, which is shown by the transition function $\mathcal{T}(s_t, a_t, s_{t+1})$. The optimum policy π^* , which leads the agent to the maximum accumulated rewards is given by

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\pi} \left\{ \sum_{t=0}^{H-1} \gamma^t r_{t+1} | s_0 = s \right\}, \quad (4)$$

where $\gamma \in [0, 1]$ is the discount factor and H is the total number of steps in one episode. To learn the optimal action, the agent receives a feedback after taking action a_t , known as the reward function $r_t = \mathcal{R}(s_t, a_t)$. Q-learning is a value-based RL framework, where the Q-value associated with the action a_t and the state s_t at time slot t , denoted by $Q(s_t, a_t)$, is expressed as

$$Q(s_t, a_t) = \mathbb{E}_{\pi} \left\{ \sum_{t=0}^{H-1} \gamma^t r_{t+1} | s_0 = s, a_0 = a, a_t = \pi(s_t) \right\}. \quad (5)$$

where π is the policy that leads to taking an action in a given state. The updated Q-value at each time slot is calculated as

$$Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \lambda(r_t + \alpha \max Q(s_{t+1}, a_{t+1})), \quad (6)$$

where $\alpha \in [0, 1]$ represents the learning rate.

III. THE PROPOSED JUNO FRAMEWORK

In this section, we first introduce the JSRL model, and then present details of the proposed JUNO framework.

A. JSRL model

Conventionally, an RL-based agent selects a random action a_t at the beginning of the learning process, where there is no prior information resulting in a long time to reach the optimal policy π . The main difference between the conventional RL and JSRL model [29] is that the agent in the JSRL framework has access to two policies, called guide-policy $\pi^g(a|s)$, and the exploration policy $\pi^e(a|s)$. While the exploration policy is the same policy used in conventional RL models, which

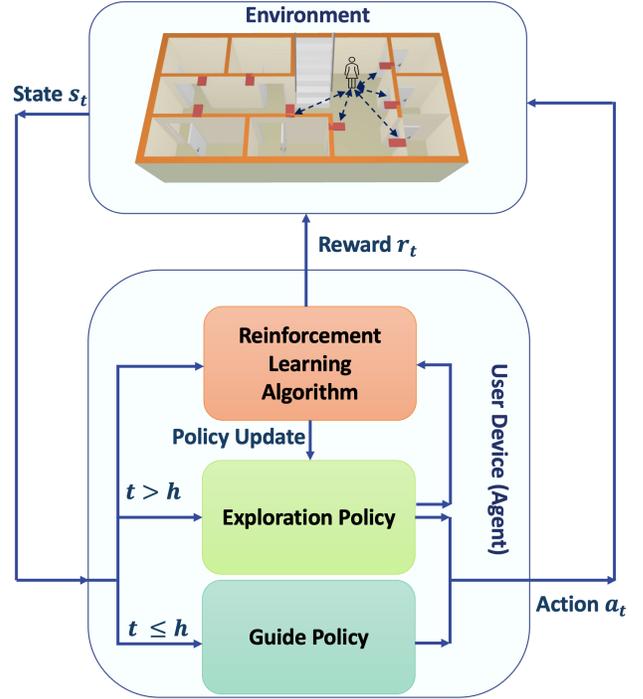


Fig. 1. The block diagram of the proposed JUNO anchor selection framework. will be updated during the training, the guide-policy is a fixed prior policy, learned via an RL model or manually/randomly constructed. Since the Q-table is initialized with zeros, the agent of the JSRL model follows $\pi^g(a|s)$ instead of randomly selecting an action a_t based on the untrained policy $\pi^e(a|s)$. Considering the fact that the distribution of the data under policy $\pi^g(a|s)$ is not exactly the same as the policy $\pi^e(a|s)$, the state space of the policy $\pi^e(a|s)$ would be different. Therefore, it is essential to gradually transit the data collection from the policy $\pi^g(a|s)$ to policy $\pi^e(a|s)$. Consequently, given an RL model with horizon H , we define guide-step $h \leq H$ as the number of steps where the agent uses policy $\pi^g(a|s)$, initialized with H and gradually decreases over the course of training. More precisely, the action is selected based on $\pi^g(a|s)$ for h steps at the initial stage of each training episode, continuing with $\pi^e(a|s)$ for remaining $(H - h)$ steps.

B. JUNO Anchor Selection

Following Reference [29], we use the JSRL model in our proposed JUNO framework to accelerate the learning process. Due to the dynamic nature of indoor venues, caused by varying environmental conditions, the proposed JUNO framework seeks to train the mobile user to autonomously find the optimal LoS connections at any given time and place. The proposed JUNO framework consists of the following main components:

(i) **Agent:** In the proposed JUNO framework, the mobile user operates as the agent, interacting with the environment via a set of actions.

(ii) **State-Space:** The state $s_t \in \mathcal{S}$ is defined as the user's location (x_t, y_t) at time slot t . Following Reference [28], we

discretize the indoor environment into $N_l = N_x \times N_y$ points, where x_t and y_t are obtained as

$$x_t = x_{t-1} + \zeta_x, \quad 0 \leq x_t \leq N_x \quad (7)$$

$$\text{and } y_t = y_{t-1} + \zeta_y, \quad 0 \leq y_t \leq N_y. \quad (8)$$

where $\zeta_x, \zeta_y \in \{-1, 0, 1\}$ are random numbers, indicating the user's movement in x -axis and y -axis [28], respectively. It should be noted that although RL models with a continuous and high-dimensional state-space [30] provide higher resolution and precise localization, they suffer from high complexity, making them inefficient for real-time applications requiring low latencies. The benefits of the discrete RL models in the context of indoor localization [27], [28], therefore, make us to choose discrete RL over its continuous counterpart.

(iii) Action-Space: The action space is defined as a set of nearby UWB beacons, where the user's location can be determined by extracting the time information from the received signal of the corresponding beacons. The cardinality of the action space, denote by N_a , is given by

$$N_a = \frac{N_u!}{(N_u - N_r)! N_r!}, \quad (9)$$

where N_r represents the number of required beacons for localization, depending on some parameters, such as the indoor localization framework (i.e., ToA, TDoA, and Two Way Ranging (TWR)), and the dimension of the experimental environment, i.e., 2D or 3D area. Considering the fact that at least two UWB beacons are required for the TDoA-based localization scheme in a 2D environment, the selected action is a vector, denoted by $\mathbf{a} = [a_i, a_j]$, where a_i, a_j represent UWB_i and UWB_j , respectively.

(iv) Reward: As stated previously, the main objective of the proposed JUNO framework is to minimize the location error caused by UWB beacons with NLoS connections. Therefore, after taking action \mathbf{a}_t , the estimated user's location $(x_{es,t}^{(i,j)}, y_{es,t}^{(i,j)})$ is calculated, where the superscript (i, j) indicates that the estimated location is obtained by the received signals of UWB_i and UWB_j . Considering the fact that even if one of these two connections is NLoS, we will face with a remarkable location error, the combination of UWB_i and UWB_j at any location/time is of paramount importance. For this reason, the reward function $\mathcal{R}(s_t, a_t)$ is defined as

$$\mathcal{R}(s_t, a_t) = \begin{cases} \frac{1}{\mathcal{E}_t}, & \mathcal{E}_t \leq \mathcal{E}_{th}, \\ -\mathcal{E}_t, & \text{o.w.} \end{cases}, \quad (10)$$

where \mathcal{E}_{th} is a pre-defined threshold value for the maximum acceptable location error [27], and \mathcal{E}_t denotes the location error at time slot t , calculated as

$$\mathcal{E}_t = \sqrt{(x_t - x_{es,t}^{(i,j)})^2 + (y_t - y_{es,t}^{(i,j)})^2}. \quad (11)$$

This completes presentation of the proposed JUMP anchor selection framework, next, we will describe our testbed and simulation results.

IV. SIMULATION RESULTS

To evaluate the performance of the proposed JUNO framework, we consider an experimental indoor area such as an office building with the size of $(60 \times 50) m^2$, which is compromised of several non-overlapping sub-areas [20], [28]. Following Reference [28], each sub-area is discretized into several square zones, where the dimension of each zone is $(1 \times 1) m^2$. Although the location resolution, i.e., the number of discretized points in the indoor environment N_l , is proportional to the location accuracy, it also results in higher state-space, complexity, and the respond-time. Therefore, there should be a trade-off between the location resolution and the respond-time of the learning model. Despite the recent RL-based localization works [27], [28], where the environment is divided into a grid of $(5 \times 5) m^2$ and $(3 \times 3) m^2$ cells, respectively, we assume higher resolution of $(1 \times 1) m^2$ to improve the location accuracy. Mobile users are randomly moving through the network in 8 directions based on Eqs. (7) and (8), where it is assumed that mobile users are placed at zones' center [28]. At each location, the transmitted signal by the mobile user are received by a set of nearby UWB beacons. Due to the obstacles in the environment, the received signal would be LoS or NLoS connections, where it is assumed that the channel condition of UWB_i , for $(1 \leq i \leq N_u)$, is determined randomly at each zone to initialize the environment.

Fig. 2(a) illustrates the effect of the guide-policy on the proposed JUNO framework. As shown in Fig. 2(a), using a random Q-table or the one obtained by an RL model as the guide-policy outperforms the conventional RL approach, accelerating the learning process of the JUNO framework, improving the location accuracy, and increasing the RL's generalization capabilities. We also investigate the effect of learning rate α on the proposed JUNO framework to obtain the best value of α . As shown in Figs. 2(b)-(c), increasing the number of epochs decreases the location error and increases the cumulative rewards, illustrating that the model is well-trained. Moreover, it is evident that the learning rate has not a great impact on the location accuracy.

Fig. 3(a) shows the effect of ϵ on the JUNO framework, as a parameter to maintain a trade-off between exploration and exploitation, where the random action a_t is chosen with the probability of ϵ . Note that $\epsilon = [1, 0.1]$ means that ϵ is initialized with 1, gradually decreasing with time by $\Delta\epsilon = \frac{\epsilon_{max} - \epsilon_{min}}{N_{epoch}}$ to 0.1, with $N_{epoch} = 100$, which is the best strategy according to the results of Fig. 3(a). Moreover, Fig. 3(b) illustrates the maximum step h_{max} that the guide-policy is initially used, where $h_{max} = 0$ represents the conventional RL model with no guide-policy. As shown in Fig. 3(b), larger h_{max} results in the lower location error.

To illustrate the effectiveness of the proposed framework, we compare it with four baseline models: (i) Weighted Least Square (WLS) anchor selection [15]; (ii) Geometric Dilution of Precision (GDOP) anchor selection [17]; (iii) Nearest Neighbor Node Selection (NN-NS), where the mobile user is localized by N_r number of nearest beacons, and; (iv) Random

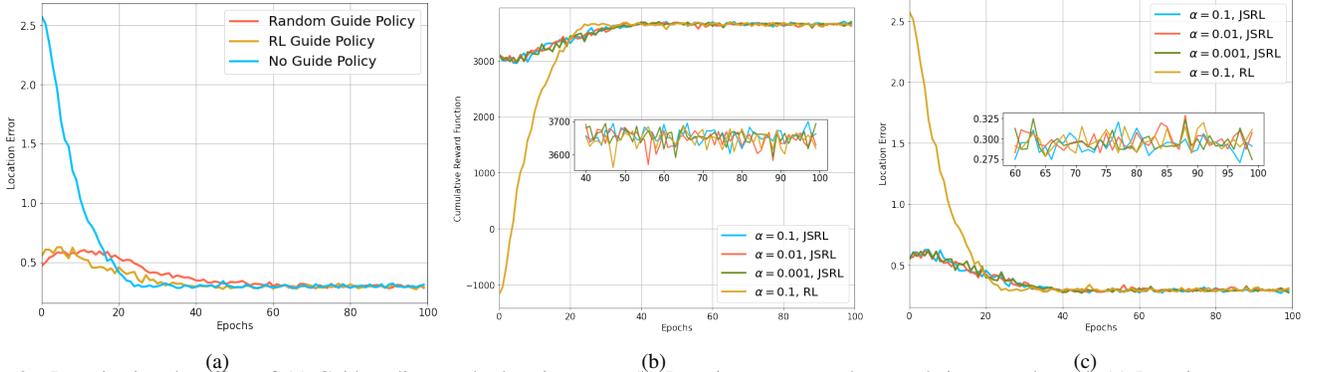


Fig. 2. Investigating the effect of (a) Guide-policy on the location error; (b) Learning rate α on the cumulative rewards, and; (c) Learning rate α on the location error.

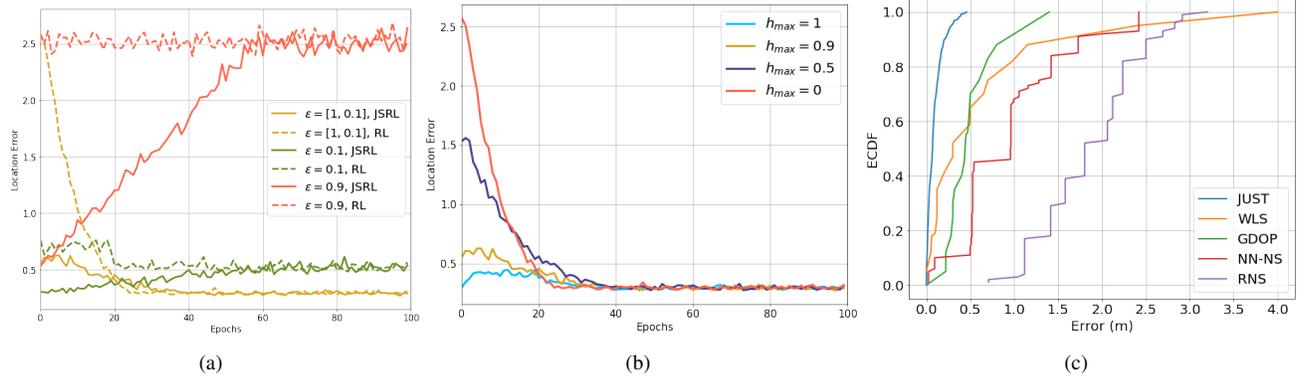


Fig. 3. Investigating the effect of (a) ϵ on the location error; (b) h_{max} on the location error, and; (c) ECDF on the location error.

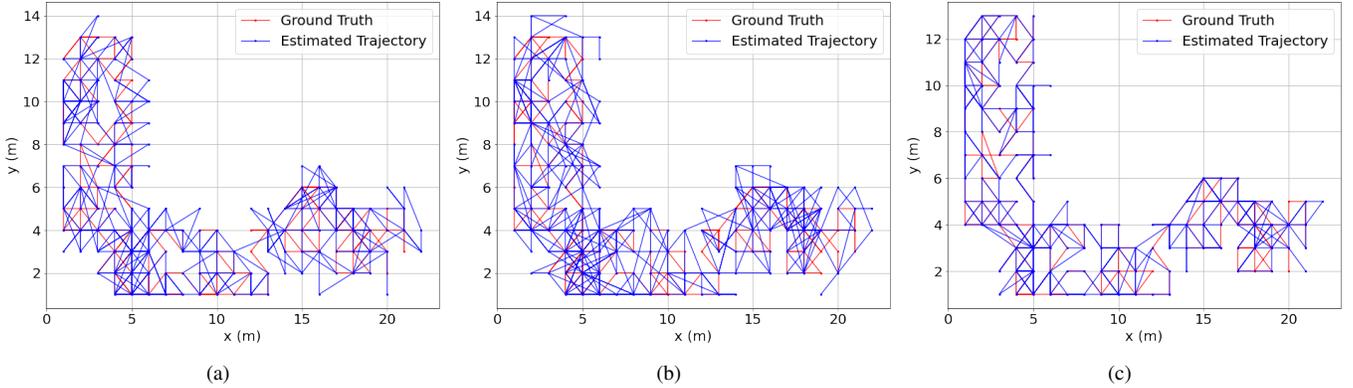


Fig. 4. Ground truth and estimated random trajectories by using: (a) Nearest neighbor; (b) Random, and; (c) JUNO frameworks.

Node Selection (RNS), where a set of UWB beacons are randomly selected for localization. Fig. 3(c) compares the Empirical Cumulative Distribution Function (ECDF) of different frameworks. According to the results shown in Fig. 3(c), the location error caused by the proposed JUNO framework is considerably lower than that of its counterparts. Finally, in Fig. 4 we consider a random trajectory in a $(24 \times 15) m^2$ rectangular indoor environment, and compare the ground truth with the estimated trajectories by our proposed JUNO framework and two other baselines. As shown in Fig. 4, the estimated path by our proposed framework in the most points

closely follows that of the ground truth.

Finally, we use the Root Mean Squared Error (RMSE), which is a generally used performance metric within the localization domain, calculated as follows

$$RMSE = \sqrt{\frac{\sum_{t=1}^N (L_t - L_{est,t})^2}{N}}, \quad (12)$$

where $L_t = (x_t, y_t)$ and $L_{est,t} = (x_{es,t}, y_{es,t})$ represent the exact and the estimated location of the user at time t , and N denotes the total number of steps once the proposed JUNO

framework reaches the steady state. The proposed JUNO framework achieves the localization error of 0.32 m, while the CNN [21] and LSTM-based [22] localization frameworks track a mobile user with 0.58, 0.48 location errors, respectively.

V. CONCLUSION

In this paper, we presented the Jump-start RL-based Uwb NNode selection (JUNO) framework with application to indoor localization. The key objective was to overcome several challenges in existing UWB-based anchor selection frameworks, such as not adapting with time-varying environmental conditions, lack of generalization, and taking long times to reach the optimal policy. Using the proposed JUNO framework, the learning process is accelerated, where the mobile user is autonomously trained to identify a set of UWB beacons with LoS connections to be localized based on the 2-D TDoA technique. The effectiveness of the proposed framework is evaluated in terms of the location error and cumulative rewards. Simulation results illustrated that the proposed JUNO framework reduced the location error. According to the simulation results, the proposed framework outperformed its counterparts in tracking a mobile user moving in a random trajectory. With the emphasis on non-uniform distribution of UWB beacons, our future research direction is to propose an adaptive version of the JUNO framework, where the cardinality of the action space can be adjusted during the learning process.

REFERENCES

- [1] R. Huang, J. Tao, L. Yang, Y. Xue and Q. Wu, "Robust TDoA Indoor Tracking Using Constrained Measurement Filtering and Grid-Based Filtering," *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 4895-4899.
- [2] F. Zafari, A. Gkelias and K. K. Leung, "A Survey of Indoor Localization Systems and Technologies," *IEEE Communications Surveys & Tutorials*, vol. 21, no. 3, pp. 2568-2599, thirdquarter 2019.
- [3] M. Atashi, M. S. Beni, P. Malekzadeh, Z. HajiAkhondi-Meybodi, K. N. Plataniotis, and A. Mohammadi, "Orientation-Matched Multiple Modeling for RSSI-based Indoor Localization via BLE Sensors," *European Signal Processing Conference (EUSIPCO)*, 2020.
- [4] M. Salimibeni, Z. HajiAkhondi-Meybodi, M. Atashi, P. Malekzadeh, K. N. Plataniotis, and A. Mohammadi, "IoT-TD: IoT Dataset for Multiple Model BLE-based Indoor Localization/Tracking," *European Signal Processing Conference*, 2020.
- [5] Z. HajiAkhondi-Meybodi, M. S. Beni, K. N. Plataniotis, and A. Mohammadi "Bluetooth Low Energy-based Angle of Arrival Estimation via Switch Antenna Array for Indoor Localization," *International Conference on Information Fusion*, July 2020.
- [6] Z. HajiAkhondi-Meybodi, M. S. Beni, A. Mohammadi, and K. N. Plataniotis, "Bluetooth Low Energy-based Angle of Arrival Estimation in Presence of Rayleigh Fading," *IEEE International Conference on Systems, Man, and Cybernetics*, 2020.
- [7] Z. HajiAkhondi-Meybodi, M. Salimibeni, A. Mohammadi and K. N. Plataniotis, "Bluetooth Low Energy and CNN-Based Angle of Arrival Localization in Presence of Rayleigh Fading," *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 7913-7917.
- [8] M. Khalaf-Allah, "Particle Filtering for Three-Dimensional TDoA-Based Positioning Using Four Anchor Nodes," *Sensors*, vol. 20, pp. 4516-4542, Jan. 2020.
- [9] R. Huang, L. Yang, J. Tao and Y. Xue, "KLD Minimization-Based Constrained Measurement Filtering For Two-Step TDoA Indoor Tracking," *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 4635-4639.
- [10] S. Zhao, X. P. Zhang, X. Cui and M. Lu, "Optimal TOA Localization for Moving Sensor in Asymmetric Network," *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 7928-7932.
- [11] L. Zhang, C. Wang, M. Ma and D. Zhang, "WiDIGR: Direction-Independent Gait Recognition System Using Commercial Wi-Fi Devices," *IEEE Internet of Things Journal*, vol. 7, no. 2, pp. 1178-1191, Feb. 2020.
- [12] J. Zhang, J. Salmi, E. S. Lohan, "Analysis of Kurtosis-based LOS/NLOS Identification using Indoor MIMO Channel Measurement," *IEEE Trans. on Vehicular Technology*, vol. 62, no. 6, pp. 2871-2874, July 2013.
- [13] W. Dai, Y. Shen, and M. Z. Win, "Sparsity-Inspired Power Allocation for Network Localization," *IEEE International Conference on Communications (ICC)*, pp. 2785-2790, June 2013.
- [14] C. Wang, Y. Ning, J. Wang, L. Zhang, J. Wan, Q. He, "Optimized Deployment of Anchors based on GDOP Minimization for Ultra-Wideband Positioning," *Journal of Spatial Science*, pp. 1-18, Nov. 2020.
- [15] A. Albaidhani, A. Morell, and J. L. Vicario, "Anchor Selection for UWB Indoor Positioning," *Transactions on Emerging Telecommunications Technologies* vol. 30, no. 6, June 2019.
- [16] A. Courty, M. L. Gentil, O. Berder, P. Scalart, S. Fontaine and A. Carer, "Anchor Selection Algorithm for Mobile Indoor Positioning using WSN with UWB Radio," *IEEE Sensors Applications Symposium (SAS)*, 2019, pp. 1-5.
- [17] A. Albaidhani, and A. Alsudani, "Anchor Selection by Geometric Dilution of Precision for an Indoor Positioning System using Ultra-Wide Band Technology," *IET Wireless Sensor Systems*, 2020.
- [18] Y. Zhu, Y. Zhang, F. Yan, L. Shen, and Y. Wu, "Node Selection for Cooperative Localization with NLOS Mitigation in Wireless Sensor Networks," Proc. in *IEEE Globecom Workshops*, Dec. 2016, pp. 1-6.
- [19] H. Zhang, X. Qi, Q. Wei, Q., and L. Liu, "TOA NLOS Mitigation Cooperative Localisation Algorithm based on Topological Unit," *IET Signal Processing*, vol. 14, no. 10, pp. 765-773, Jan. 2021.
- [20] S. Wang, and Y. Zhang, "Convex Hull based Node Selection NLoS Mitigation for Indoor Localization," Proc. in *IEEE Wireless Communications and Networking Conference*, Apr. 2016, pp. 1-5.
- [21] D. T. A. Nguyen, et al. "Deep learning-based localization for UWB systems," *Electronics*, val. 9, no. 10, pp. 1712, 2020.
- [22] A. Poulou, D. S. Han, "UWB Indoor Localization Using Deep Learning LSTM Networks," *Appl. Sci.* 2020, pp. 6290.
- [23] M. Salimibeni, Z. Hajiakhondi-Meybodi, A. Mohammadi, and Y. Wang, "TB-ICT: A Trustworthy Blockchain-Enabled System for Indoor COVID-19 Contact Tracing," *arXiv preprint arXiv:2108.08275*, 2021.
- [24] J. Fan and A. S. Awan, "Non-Line-of-Sight Identification Based on Unsupervised Machine Learning in Ultra-Wideband Systems," *IEEE Access*, vol. 7, pp. 32464-32471, 2019.
- [25] Z. Hajiakhondi-Meybodi, A. Mohammadi, M. Hou, and K. N. Plataniotis, "DQLEL: Deep Q-Learning for Energy-Optimized LoS/NLoS UWB Node Selection," accepted in *IEEE Transactions on Signal Processing*, 2022.
- [26] D. Miliotis, "Efficient Indoor Localization via Reinforcement Learning," *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 8350-8354.
- [27] M. Mohammadi, A. Al-Fuqaha, M. Guizani and J. Oh, "Semisupervised Deep Reinforcement Learning in Support of IoT and Smart City Services," *IEEE Internet of Things Journal*, vol. 5, no. 2, pp. 624-635, April 2018.
- [28] Y. Li, X. Hu, Y. Zhuang, Z. Gao, P. Zhang and N. El-Sheimy, "Deep Reinforcement Learning (DRL): Another Perspective for Unsupervised Wireless Localization," *IEEE Internet of Things Journal*, vol. 7, no. 7, pp. 6279-6287, July 2020.
- [29] I. Uchendu, T. Xiao, Y. Lu, B. Zhu, M. Yan, J. Simon, M. Bennice, C. Fu, C. Ma, J. Jiao, and S. Levine, "Jump-Start Reinforcement Learning," *arXiv preprint arXiv:2204.02372*, Apr. 2022.
- [30] E. Marchesini et al., "Discrete Deep Reinforcement Learning for Mapless Navigation," *IEEE Inter. Conf. on Robotics and Automation (ICRA)*, 2020, pp. 10688-10694.
- [31] L. Gen, et al. "Is Q-learning minimax optimal? A tight sample complexity analysis," *arXiv preprint arXiv:2102.06548*, 2021.