

Federated Learning for Vision-based Obstacle Avoidance in the Internet of Robotic Things

Xianjia Yu[†], Jorge Peña Queralta[†], Tomi Westerlund[†]

[†]Turku Intelligent Embedded and Robotic Systems (TIERS) Lab, University of Turku, Finland.
Emails: ¹{xianjia.yu, jopequ, toveve}@utu.fi

Abstract—Deep learning methods have revolutionized mobile robotics, from advanced perception models for an enhanced situational awareness to novel control approaches through reinforcement learning. This paper explores the potential of federated learning for distributed systems of mobile robots enabling collaboration on the Internet of Robotic Things. To demonstrate the effectiveness of such an approach, we deploy wheeled robots in different indoor environments. We analyze the performance of a federated learning approach and compare it to a traditional centralized training process with a priori aggregated data. We show the benefits of collaborative learning across heterogeneous environments and the potential for sim-to-real knowledge transfer. Our results demonstrate significant performance benefits of FL and sim-to-real transfer for vision-based navigation, in addition to the inherent privacy-preserving nature of FL by keeping computation at the edge. This is, to the best of our knowledge, the first work to leverage FL for vision-based navigation that also tests results in real-world settings.

Index Terms—Federated learning; distributed robotic systems; internet of robotic things; vision-based obstacle avoidance; autonomous robots; collaborative learning.

I. INTRODUCTION

As ubiquitous autonomous mobile robots become increasingly interconnected, end-users and applications can benefit from distributed multi-robot systems [1]. Connected robots open a wide variety of opportunities within the Internet of Robotic Things (IoRT) [2], specifically as mobile sensor networks capable of intelligent behavior and multi-modal data acquisition. We are particularly interested in exploring collaborative robot learning within the IoRT context [3], [4], and studying the benefits that federated learning FL can bring to distributed robotic systems.

The relevance of AI and deep learning (DL) in robotic systems has increased significantly as DL methods enable higher degrees of situational awareness [5], [6]. For a variety of tasks in robotics, mobile navigation, human-like walking, teaching through demonstration, and collaborative automation, to mention a few examples, DL or machine learning has attained state-of-the-art performance [5]. One of the areas where DL has had more impact is arguably computer vision, which has led us to put the focus on a relevant use-case for autonomous mobile robots: vision-based obstacle avoidance. Moreover, vision-based obstacle avoidance is instrumental not only for robot navigation but could also aid visually handicapped citizens [7]. Research has shown that DL can

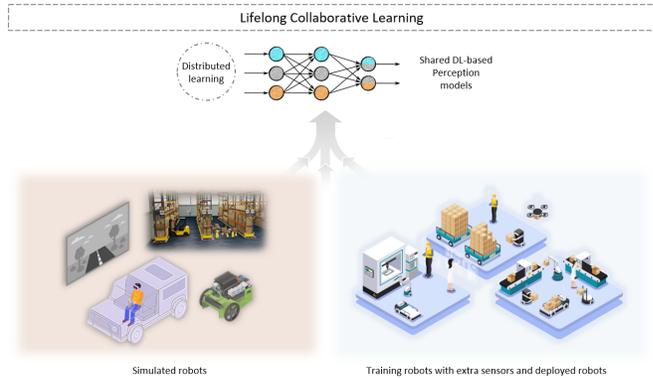


Fig. 1: Conceptual illustration of federated lifelong learning with sim-to-real transfer as a deep learning-based perception model is trained in both a simulation and real robots, with potentially continuous updates.

boost vision-based obstacle avoidance [8] and, especially when combined with specific technologies such as multi-view structure-from-motion, can reach a better degree of precision than other approaches [9].

In this manuscript, we focus on studying a federated learning approach for vision-based obstacle avoidance in distributed IoRT systems. A collective model built from and shared within a team of robots operating in different environments can bring numerous benefits, including more robust performance but more importantly higher readiness levels to operate in new environments or properly react to new situations. Collaborative multi-robot systems can be more efficient and have higher success rates in heterogeneous environments including unknown ones [3].

Federated learning brings multiple benefits to collaborative learning in the IoRT. Compared to cloud-based centralized learning, it allows for optimization of networking resources and the preservation of data privacy by computing model updates directly at the edge. We propose a FL approach that combines data from both simulated and real robotic agents. A conceptual illustration of such system is illustrated in fig. 1. We study the performance of FL over centralized learning in the simulated and real worlds separately, analyze the improvements of merging data from heterogeneous scenarios, and finally the potential for sim-to-real knowledge transfer.

Partly owing to the benefits of sharing knowledge without transferring raw data, FL as a privacy-preserving distributed learning has been utilized in multiple domains in robotics and autonomous system [4]. These domains include, but certainly are not limited to, navigation, cooperative SLAM based on visual-Lidar, trajectory forecasting, human-robot collaborative learning, and robot perception. To the best of our knowledge, however, there is no specific work related to applying FL in robotic vision-based obstacle avoidance in real-world scenarios.

In summary, in this work, we explore the potential for FL within hybrid teams of simulated and real robotic agents. The main contributions of this work are the following. First, the design, implementation, deployment, and evaluation of a vision-based deep-learning approach to obstacle avoidance in mobile robots in heterogeneous simulated and real scenarios. We then evaluate the performance benefits of such an approach over offline learning or learning from more limited data sources. We put an emphasis on the benefits when robots are deployed in heterogeneous environments, showing that collaborative learning improves performance even for robots that do not change their environment. For this work, we deploy robots in highly photorealistic and physically-accurate virtual environments and study the ability of such a setup for sim-to-real transfer.

The remainder of this manuscript is organized as follows. Section II introduces related work in the relevant robot learning and federated learning literature. We then explain the methods and tools used for this work in Section III. Section IV reports the experimental results, while Section V concludes the work and lays down future research directions.

II. RELATED WORK

This section introduces relevant examples in the literature in the areas of FL for robotics, DL for vision-based obstacle avoidance, and sim-to-real transfer. We also discuss different simulation environments that can be used for robot learning.

A. Federated learning in robotics

Robot collaboration has become vital as networked robots have become more ubiquitous [1], and FL has stood as one of the best solutions from the point of view of data privacy, network resource management, and distributed edge computing [10]. Other technologies, such as differential privacy, homomorphic encryption, and distributed ledger technologies (DLTs) have been used in the literature to improve FL from a systematic standpoint, making the collaborative learning process in a multi-robot system safer and privacy-preserving. FL offers potential in a variety of autonomy challenges and robotic subsystems, including cooperative SLAM, human-robot collaborative learning, and navigation, to name a few [4].

B. Deep learning for vision-based obstacle avoidance

Deep learning has been widely applied to robotic applications and has arguably revolutionized the role of AI in robotics [5], [6], [11]. More specifically, it has gained

popularity in vision-based obstacle avoidance because of the wide availability of vision sensors across platforms, and its suitability to different types of environments, including the absence of geometric models and substantial parameter tuning. In a recent relevant study [12], the authors demonstrate the use of deep reinforcement learning (DRL) to perform obstacle avoidance via a monocular camera mounted on a UAV with minimal knowledge of the environment. In addition, DRL could perform more effectively even with irregular input by embedding the image's depth and semantic information. By adding noise to the depth information, the DRL model's resilience was increased to near-state-of-the-art levels in some unseen virtual and real-world circumstances [13].

C. Sim-to-real for robot learning

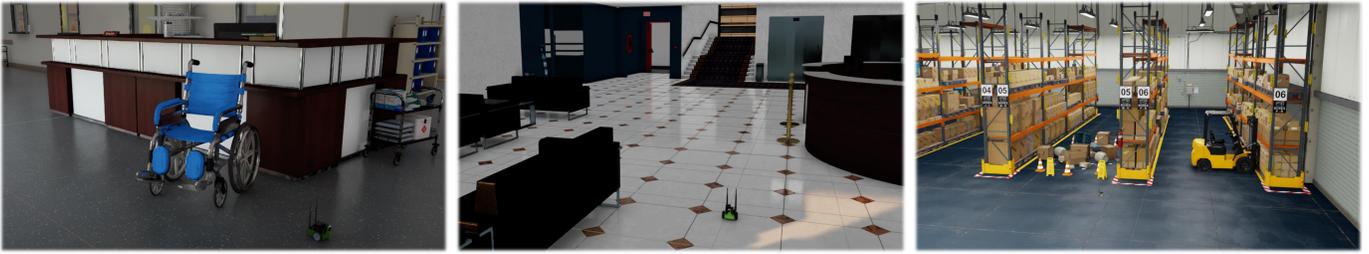
Sim-to-real research has mostly focused on policies obtained through DRL [14], and specifically in areas related to robotic manipulation [15]. Recent studies however have also applied DRL to vision-based obstacle avoidance in robot manipulators, demonstrating that learned models can be adapted efficiently to unseen scenes and unseen objects in the real world [16], [17]. In this work, we take the sim-to-real analysis to mobile robots and explore the potential for collaborative learning of policies that can aid in robust autonomous navigation.

D. Photorealistic simulation platforms

Availability of data is one of the key challenges for the application of DL to robotics. Synthetic data and data from simulated environments have therefore played a key role in numerous studies. However, traditional robotics simulators such as Gazebo, widely used in the development of mobile robots, are incapable of generating realistic visuals. Therefore, while they are efficient for testing autonomy stacks based on other types of sensors such as lidars, which only need accurate geometric data, they are unsuitable for vision-based approaches and limit the sim-to-real knowledge transferability. A number of simulation platforms have emerged in recent years to address this issue. We have listed a relevant set of simulators in table I. Among these, Carla [18] has been widely used in research for self-driving cars, while AirSim [19] has been used in a wide variety of application scenarios. We have chosen for this study, however, the more recent NVIDIA Isaac Sim platform owing to the high-quality visuals but also tools that enable seamless generation of randomized environments for synthetic data acquisition. Randomization and the ability to alter the environments has been shown to be a key parameter to collaborative learning approaches [15], [20]. An additional advantage is a common ecosystem of tools with the embedded NVIDIA Jetson platforms, the state-of-the-art in embedded computing for robots that need discrete GPUs for DL inference.

III. METHODOLOGY

This section covers the different tools, simulation environments and robots utilized in the experiments. We describe the



(a) Hospital

(b) Office

(c) Warehouse

Fig. 2: Customized simulation environments in NVIDIA Isaac Simulator used in the experiments.

TABLE I: Comparison of existing 3D robotics simulators

Simulator platform	Photorealistic visuals	ROS support	Multi-modal sensors
Gazebo [21]		✓	✓
VRKitchen [22]	✓		
MINOS [23]			
Gibson Env [24]	✓	✓	
Habitat [25]	✓		
PreSim [26]	✓		
Carla [18]	✓	✓	✓
AirSim [19]	✓	✓	✓
Nvidia Isaac [27]	✓	✓	✓

approaches to centralized and federated learning, and the deep learning models used for vision-based obstacle avoidance.

A. Simulation settings

As a platform to validate the proposed approach in simulated robots, we have used NVIDIA Isaac Sim, powered by Omniverse. Isaac Sim is a scalable robotics simulation application and synthetic data generation tool that enables the creation of photorealistic, physically accurate virtual environments for developing, testing, and managing AI-based robots [27].

We have set up a series of simulation environments to gather data and validate the vision-based approach to obstacle avoidance. Three main environments are used to analyze the performance of the trained model, with a focus on heterogeneity of objects and backgrounds. The datasets used in this study include data from environments that replicate a hospital (see fig. 2a), a office room (see fig. 2b), and a warehouse (see fig. 2c).

To obtain sufficient data for model training, we used the NVIDIA Isaac Simulator’s domain randomization (DR) and synthetic data recorder (SDR) features. By utilizing DR, we can select a specific object and randomly set its properties such as movement, rotation, light, and texture within a defined range. We can easily record the data generated by DR by using SDR. These operations were carried out in the three customized environments mentioned previously and pictured in fig. 2. The distribution of the datasets is shown in table II, with S_i representing the dataset associated to environment i .

B. Real-world experimental settings

We utilize three mobile robots for real-world experiments and a local computing server for training the deep learning

TABLE II: Distribution of simulation datasets

	Hospital (S_1)		Office (S_2)		Warehouse (S_3)	
	blocked	free	blocked	free	blocked	free
Prop.	44%	56%	64%	46%	60%	40%
Total	27%		54%		19%	

TABLE III: Distribution of real-world datasets

	Room 1 (\mathcal{R}_1)		Room 2 (\mathcal{R}_2)		Room 3 (\mathcal{R}_3)	
	blocked	free	blocked	free	blocked	free
Prop.	40%	60%	50%	50%	50%	50%
Total	11%		44%		45%	

models. The platforms used in the experiments are three Jetbot robots from Waveshare (depicted in fig. 3), equipped by default with a wide-angle lens RGB camera and an embedded NVIDIA Jetson Nano development kit. In addition, a Rplidar A1 2D rangefinder has been installed on a 3D printed frame for automated data labeling when real-world data is collected. The local edge server used for training the obstacle avoidance models is equipped with an 8-core Intel i7-9700K processor, 64 GB of memory, and an NVIDIA GeForce RTX 2080 Ti discrete graphics card.

We deployed the mobile robots in three different indoor environments (office spaces, hallways and laboratory environments) to validate the obstacle avoidance policies trained with the different approaches. The three rooms vary in terms of objects present as obstacles, material texture, layout, and style. The distribution of the images acquired by the three robots is shown in table III, where $\mathcal{R}_i, i \in \{0,1,2\}$ represent each of the acquired datasets for the respective real-world environments. We intentionally imbalance the data in order to imitate a situation in which robots would not be able to collect data equally across different environments and operational conditions. We also do this to evaluate whether there is a performance impact in environments based on the amount of collected data. As such, *Room 1* only accounts for 11% of the total amount of collected images.

C. Vision-based obstacle avoidance models

In order to achieve vision-based obstacle avoidance for the operation of mobile robots, we have chosen to train a generic



Fig. 3: Customized Jetbot platform

DL model to assist robots in discriminating between different types of obstacles across heterogeneous environments. This approach comes in contrast to other options, including the detection of individual objects or semantic segmentation (e.g., for segmenting free floor from objects and walls). The selected approach enables us to focus on analyzing the performance of a federated learning approach and the ability for sim-to-real transfer rather than on the design of a specific obstacle avoidance strategy, which is the main objective of this study

More precisely, we utilized a deep convolutional neural network (CNN) to carry out a vision-based obstacle classifier for two only two classes, that define whether the environment ahead is *blocked* or *free* for the robot to navigate. Owing to the relative low level of complexity of the classifier and the limited size of the collected datasets, we have selected the AlexNet [28] architecture as appropriate for such binary classification task. AlexNet has been established as the precedent for deep CNN as one of the most widely used backbones for executing various tasks across multiple domains.

We train different models for each separate dataset with both simulated and real data, as well as combinations of these. The models are trained using two approaches: a centralized learning approach that aggregates data from all robots in the local edge server and trains them at once; and a federated learning approach that only fuses the individual models trained in each of the different scenarios.

IV. EXPERIMENTAL RESULTS

Through this section, we report the experimental results obtained with data from both simulated and real robots. We show first the performance of the different approaches, with the latter part shifting towards the potential for sim-to-real knowledge transferability.

A. Centralized training vs. federated learning

The first objective of our experiments is to analyze the performance improvements that a federated learning approach brings over a centralized training with traditional data aggregation. To do this, we used the data we collected in the simulated hospital (\mathcal{S}_0), office (\mathcal{S}_1), and warehouse (\mathcal{S}_2) to train our model on each dataset and all possible combinations of two or three of the datasets. Equivalently for the federated

learning approach, we run different training rounds in which we simulate that a different subset of robots is collaboratively learning without sharing any actual raw data. In this approach, only the models are fused and a common model updated iteratively. fig. 4a and fig. 4b report the accuracy of the different models for the centralized and federated approaches, respectively. For the FL results, the training happens only with combination of datasets from different environments.

The accuracy of models trained with real-world data is then shown in fig. 4c and fig. 4d. The data has been obtained with Jetbots navigating in three different office and laboratory indoor environments ($\mathcal{R}_i, i \in \{0,1,2\}$).

In addition to the accuracy, we also calculate the area under the ROC curve (AUC) for each of the scenarios where training is carried out through either the centralized or federated approaches. The results are reported in table IV This metric gives a better understanding of the reliability of the models. In this particular application scenario of robotic navigation, there is indeed a disparity in the cost of false negatives over false positives in terms of the robot’s integrity. However, from the point of view of performance, false positives can degrade significantly the navigation speed and time, while low-frequency collisions can be mitigated with, e.g., bumper sensors. Therefore, the classification-threshold invariance of the AUC metric is relevant to this use case.

Through both the accuracy and AUC results, we see that there is a clear improvement when the models and not the data are aggregated. In addition to the better navigation results, this also brings other advantages. First, we optimize the networking resources, allowing for intermittent connectivity and potentially lower bandwidth requirements when the size of the data in training batches is significantly smaller than the models. Moreover, this allows for privacy-preserving collaboration between different end-users or robot operators, as the raw data does not need to be exposed to a central authority or service.

B. Sim-to-real performance evaluation

In the last part of our experiments, we evaluate the ability of both the centralized and federated learning approaches to transfer knowledge from simulation environments to the real world. To do this, we rely on the same simulation environments but introduce an independent data validation set (\mathcal{R}^*) from a navigation mission across different types of indoor spaces. The accuracy for each of the trained models is shown in fig. 5. We can observe that relatively low performance is achieved with either approach when only one of the simulation environments is used for training the models. This may result from overfitting the model to non-realistic features in the simulated worlds. However, when heterogeneous data is introduced in training, the federated learning approach significantly improves. Our results also show that only when aggregating data from all three simulation environments the centralized learning approach is able to improve the performance to the level of federated learning. The specific reason behind this

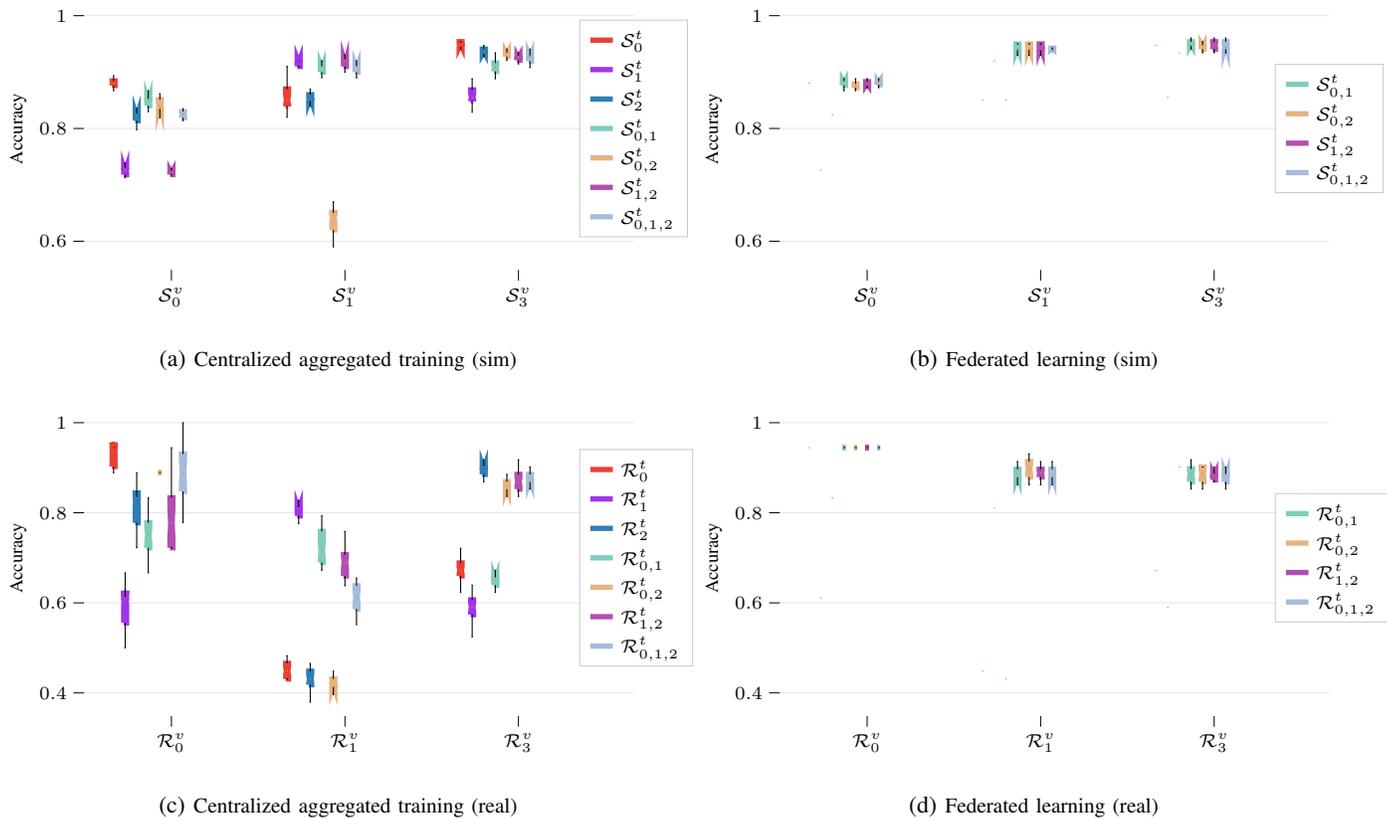


Fig. 4: Accuracy of the different models obtained through centralized learning with aggregated data or federated learning with fused local models. These results are trained (t) and validated (v) with respective simulation datasets ($\mathcal{S}_i^t, \mathcal{S}_i^v$) and real datasets ($\mathcal{R}_i^t, \mathcal{R}_i^v$) independently.

TABLE IV: Area under ROC curve (AUC) values for the aggregated centralized learning and federated learning approaches.

		Training datasets											
		Centralized learning with aggregated data							Federated learning				
		\mathcal{S}_0^t	\mathcal{S}_1^t	\mathcal{S}_2^t	$\mathcal{S}_{0,1}^t$	$\mathcal{S}_{0,2}^t$	$\mathcal{S}_{1,2}^t$	$\mathcal{S}_{1,2,3}^t$	$\mathcal{S}_{0,1}^t$	$\mathcal{S}_{0,2}^t$	$\mathcal{S}_{1,2}^t$	$\mathcal{S}_{1,2,3}^t$	
Validation datasets	Sim	\mathcal{S}_0^v	0.28	0.56	0.56	0.33	0.50	0.71	0.52	0.85	0.85	0.85	0.85
		\mathcal{S}_1^v	0.33	0.50	0.75	0.33	0.50	0.75	0.60	0.94	0.93	0.93	0.95
		\mathcal{S}_2^v	0.43	0.94	0.46	0.42	0.50	0.23	0.62	0.96	0.95	0.95	0.95
Real	\mathcal{R}_0^v	0.63	0.58	0.42	0.75	0.63	0.08	0.58	0.88	0.88	0.88	0.88	
	\mathcal{R}_1^v	0.31	0.66	0.60	0.35	0.25	0.50	0.70	0.83	0.85	0.85	0.85	
	\mathcal{R}_2^v	0.69	0.57	0.87	0.46	0.51	0.48	0.56	0.90	0.92	0.90	0.92	

behaviour requires further study and will be the object of future research.

V. DISCUSSION AND ANALYSIS

We have presented a federated learning approach for vision-based obstacle avoidance in mobile robots that leverages data from both simulated agents and real robots with additional sensors. We have shown that interconnected robots relying on deep learning for vision-based navigation can aid each other without sharing raw data. Specifically, we show how training the same model with data from heterogeneous environments

improves performance across the simulated and real worlds. More importantly, the performance improvements are better when the models are trained through a federated learning approach compared to centralized learning. In addition to the application-specific improvements, the federated learning approach brings inherent benefits in terms of communication optimization and preservation of data privacy, enabling collaboration across organizations or users. Finally, we have shown that the presented approach is able to transfer knowledge from simulation to reality effectively.

Owing to the potential for sim-to-real transfer and account-

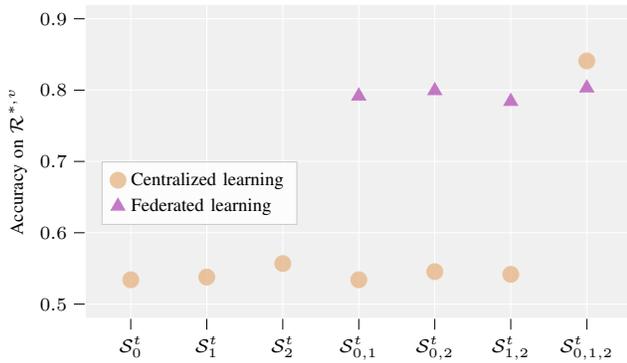


Fig. 5: Sim-to-real accuracy of simulation-trained models validated on an independent real-world navigation dataset.

ing for the better performance of federated over centralized learning, future work will be directed towards lifelong FL through a combination of simulated and real agents. Additionally, we will further explore the reasons behind the differences in performance across the approaches presented in this manuscript.

ACKNOWLEDGMENT

This research work is supported by the Academy of Finland's AutoSOS project (Grant No. 328755) and RoboMesh project (Grant No. 336061).

REFERENCES

- [1] Jorge Peña Queralta, Jussi Taipalmaa, Bilge Can Pullinen, Victor Kathan Sarker, Tuan Nguyen Gia, Hannu Tenhunen, Moncef Gabbouj, Jenni Raitoharju, and Tomi Westerlund. Collaborative multi-robot search and rescue: Planning, coordination, perception, and active vision. *IEEE Access*, 8:191617–191643, 2020.
- [2] Pieter Simoons, Mauro Dragone, and Alessandro Saffiotti. The internet of robotic things: A review of concept, added value and applications. *International Journal of Advanced Robotic Systems*, 15(1):1729881418759424, 2018.
- [3] Ertug Olcay, Fabian Schuhmann, and Boris Lohmann. Collective navigation of a multi-robot system in an unknown environment. *Robotics and Autonomous Systems*, 132:103604, 2020.
- [4] Yu Xianjia, Jorge Peña Queralta, Jukka Heikkinen, and Tomi Westerlund. Federated learning in robotic and autonomous systems. *Procedia Computer Science*, 191:135–142, 2021.
- [5] Harry A Pierson and Michael S Gashler. Deep learning in robotics: a review of recent research. *Advanced Robotics*, 31(16):821–835, 2017.
- [6] Artúr István Károly, Péter Galambos, József Kuti, and Imre J Rudas. Deep learning in robotics: Survey on model structures and training strategies. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 51(1):266–279, 2020.
- [7] Mateus Mendes, AP Coimbra, and MM Crisostomoy. Assis-cicerone robot with visual obstacle avoidance using a stack of odometric data. *IAENG Int. J. Comput. Sci.*, 45:219–227, 2018.
- [8] Joel O Gaya, Lucas T Gonçalves, Amanda C Duarte, Breno Zanchetta, Paulo Drews, and Silvia SC Botelho. Vision-based obstacle avoidance using deep learning. In *2016 XIII Latin American Robotics Symposium and IV Brazilian Robotics Symposium (LARS/SBR)*, pages 7–12. IEEE, 2016.
- [9] Shichao Yang, Sandeep Konam, Chen Ma, Stephanie Rosenthal, Manuela Veloso, and Sebastian Scherer. Obstacle avoidance through deep networks based intermediate perception. *arXiv preprint arXiv:1704.08759*, 2017.
- [10] Ahmed Imteaj, Urmish Thakker, Shiqiang Wang, Jian Li, and M Hadi Amini. A survey on federated learning for resource-constrained iot devices. *IEEE Internet of Things Journal*, 9(1):1–24, 2021.
- [11] Yu Xianjia, Sahar Salimpour, Jorge Peña Queralta, and Tomi Westerlund. Analyzing general-purpose deep-learning detection and segmentation models with images from a lidar as a camera sensor. In *International Conference on Intelligent Systems Design and Engineering Applications, Lecture Notes in Electrical Engineering (to appear)*. Springer, 2022.
- [12] Abhik Singla, Sindhu Padakandla, and Shalabh Bhatnagar. Memory-based deep reinforcement learning for obstacle avoidance in uav with limited environment knowledge. *IEEE Transactions on Intelligent Transportation Systems*, 22(1):107–118, 2021.
- [13] Lingping Gao, Jianchuan Ding, Wenxi Liu, Haiyin Piao, Yuxin Wang, Xin Yang, and Baocai Yin. A vision-based irregular obstacle avoidance framework via deep reinforcement learning. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 9262–9269, 2021.
- [14] Wenshuai Zhao, Jorge Peña Queralta, and Tomi Westerlund. Sim-to-real transfer in deep reinforcement learning for robotics: a survey. In *2020 IEEE Symposium Series on Computational Intelligence (SSCI)*, pages 737–744. IEEE, 2020.
- [15] Wenshuai Zhao, Jorge Peña Queralta, Li Qingqing, and Tomi Westerlund. Towards closing the sim-to-real gap in collaborative multi-robot deep reinforcement learning. In *2020 5th International Conference on Robotics and Automation Engineering*, pages 7–12. IEEE, 2020.
- [16] Kefang Zhang, Jiatao Lin, Lv Bi, and Tan Zhang. Sim2real learning of vision-based obstacle avoidance for robotic manipulators. In *RSS Workshop*, 2020.
- [17] Tan Zhang, Kefang Zhang, Jiatao Lin, Wing-Yue Geoffrey Louie, and Hui Huang. Sim2real learning of obstacle avoidance for robotic manipulators in uncertain environments. *IEEE Robotics and Automation Letters*, 7(1):65–72, 2021.
- [18] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. Carla: An open urban driving simulator. In *Conference on robot learning*, pages 1–16. PMLR, 2017.
- [19] Shital Shah, Debadepta Dey, Chris Lovett, and Ashish Kapoor. Airsim: High-fidelity visual and physical simulation for autonomous vehicles. In *Field and service robotics*, pages 621–635. Springer, 2018.
- [20] Wenshuai Zhao, Jorge Peña Queralta, Li Qingqing, and Tomi Westerlund. Ubiquitous distributed deep reinforcement learning at the edge: Analyzing byzantine agents in discrete action spaces. *Procedia Computer Science*, 177:324–329, 2020.
- [21] Nathan Koenig and Andrew Howard. Design and use paradigms for gazebo, an open-source multi-robot simulator. In *2004 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)(IEEE Cat. No. 04CH37566)*, volume 3, pages 2149–2154. IEEE, 2004.
- [22] Xiaofeng Gao, Ran Gong, Tianmin Shu, Xu Xie, Shu Wang, and Song-Chun Zhu. Vrkitchen: an interactive 3d virtual environment for task-oriented learning. *arXiv preprint arXiv:1903.05757*, 2019.
- [23] Manolis Savva, Angel X. Chang, Alexey Dosovitskiy, Thomas Funkhouser, and Vladlen Koltun. MINOS: Multimodal indoor simulator for navigation in complex environments. *arXiv:1712.03931*, 2017.
- [24] Fei Xia, Amir R. Zamir, Zhi-Yang He, Alexander Sax, Jitendra Malik, and Silvio Savarese. Gibson env: real-world perception for embodied agents. In *Computer Vision and Pattern Recognition (CVPR), 2018 IEEE Conference on*. IEEE, 2018.
- [25] Andrew Szot, Alex Clegg, Eric Undersander, Erik Wijmans, Yili Zhao, John Turner, Noah Maestre, Mustafa Mukadam, Devendra Chaplot, Oleksandr Maksymets, Aaron Gokaslan, Vladimir Vondrus, Sameer Dharur, Franziska Meier, Wojciech Galuba, Angel Chang, Zsolt Kira, Vladlen Koltun, Jitendra Malik, Manolis Savva, and Dhruv Batra. Habitat 2.0: Training home assistants to rearrange their habitat. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [26] Honglin Yuan and Remco C Veltkamp. Presim: A 3d photo-realistic environment simulator for visual ai. *IEEE Robotics and Automation Letters*, 6(2):2501–2508, 2021.
- [27] Nvidia. Nvidia isaac sim. <https://developer.nvidia.com/isaac-sim>. [Online] - Last access: 2022-01-26.
- [28] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25:1097–1105, 2012.