Multi-agent Reinforcement Learning for Distributed Cooperative Vehicular Positioning

Bernardo Camajori Tedeschini, *Student Member, IEEE*, Mattia Brambilla, *Member, IEEE*, Monica Nicoli, *Senior Member, IEEE* and Moe Z. Win, *Fellow, IEEE*

Abstract—With the advent of cooperative intelligent transport systems (C-ITS) and vehicle-to-everything (V2X) communications, cooperative positioning based on V2X sharing of location information has been emerging as a promising augmentation system for conventional satellite navigation. An example is implicit cooperative positioning (ICP) which relies on Bayesian filtering for cooperative sensing of targets that are used as reference points for improving vehicle positioning. ICP methods, however, rely on pre-determined models which makes them sub-optimal in case of non-Gaussian non-linear models or complex cooperation graphs. To address these limitations, the paper proposes a decentralized-partially observable Markov decision process (Dec-POMDP) framework, paired with deep multi-agent reinforcement learning (MARL) algorithms. We introduce a novel ICP-multiagent proximal policy optimization (MAPPO) algorithm where distributed agents (i.e., vehicles) dynamically activate/deactivate the radio links for cooperation with the neighbors to optimize the communication efficiency, still guaranteeing accurate positioning. We reproduce a realistic C-ITS scenario with CARLA simulator, where vehicles move according to real-world dynamics and communicate with each other to cooperatively sense their locations. Results show that the proposed ICP-MAPPO algorithm, with its dynamic-decentralized-execution and centralizedtraining schemes, outperforms state-of-the-art ICP methods by 21% in terms of positioning accuracy, and it can reduce the communication overhead by following the optimal learned policy.

Index Terms—MARL, Dec-POMDP, implicit cooperative positioning, Bayesian-filtering, message passing algorithm.

Manuscript received 21 June 2024; revised 7 September 2024; accepted 19 September 2024. Date of publication XX YYYY 2024; date of current version XX YYYY 2024. The fundamental research described in this paper was supported, in part, by the Roberto Rocca Doctoral Fellowship granted by the Massachusetts Institute of Technology and Politecnico di Milano, by MOST – Sustainable Mobility National Research Center European Union Next-GenerationEU (piano nazionale di ripresa e resilienza (PNRR) – missione 4 componente 2, investimento 1.4 – D.D. 1033 17/06/2022, CN00000023), by the National Science Foundation under Grant CNS-2148251, and by the federal agency and industry partners in the RINGS Program. The material in this paper has been presented in part to the International Conference on Information Fusion (FUSION), Venice, Italy, July 2024. (Corresponding author: Moe Z. Win.)

- B. Camajori Tedeschini is with Dipartimento di Elettronica, Informazione e Bioingegneria (DEIB), Politecnico di Milano, 20133 Milan, Italy and with the Wireless Information and Network Sciences Laboratory, Massachusetts Institute of Technology, Cambridge, MA 02139 USA (e-mail: bernardo.camajori@polimi.it, berni97@mit.edu).
- M. Brambilla is with Dipartimento di Elettronica, Informazione e Bioingegneria (DEIB), Politecnico di Milano, 20133 Milan, Italy (e-mail: mattia.brambilla@polimi.it).
- M. Nicoli is with Dipartimento di Ingegneria Gestionale (DIG), Politecnico di Milano, 20156 Milan, Italy (e-mail: monica.nicoli@polimi.it).
- M. Z. Win is with the Laboratory for Information and Decision Systems (LIDS), Massachusetts Institute of Technology, Cambridge, MA 02139 USA (e-mail: moewin@mit.edu).
- Color versions of one or more figures in this article are available at https://doi.org/10.1109/YYYY.2024.XXXXXXX.

Digital Object Identifier 10.1109/YYYY.2024.XXXXXXX

I. INTRODUCTION

▼OOPERATIVE POSITIONING (CP) represents a key enabling feature for future automated mobility services [1]-[8]. Automated vehicles leverage an on-board sensor suite including global navigation satellite systems (GNSS), light detection and ranging (LIDAR), radio detection and ranging (RADAR), and stereo cameras to perceive the surrounding environment and perform automated maneuvers [9]-[13]. At today, these sensors are not yet able to guarantee highprecision localization in harsh environments such as dense urban areas or canyons and this is a main issue for autonomous driving functions [14]. Recently, methods have been proposed to combine localization sensors with the latest 5th generation (5G) of cellular communications [15]-[20], depicting a new horizon for mobile connectivity and positioning services [21]– [24]. 5G vehicle-to-everything (V2X) communications are envisioned as crucial in the evolution towards cooperative intelligent transport systems (C-ITS) [25]-[28] by enabling simultaneous communication and localization functionalities [29]–[31]. CP among vehicles, by means of sidelink V2X communications, can be used to overcome the GNSS performance degradation and guarantee a seamless high-accuracy positioning (HAP) service [32]-[36]. The complexity lies in the resource-intensive nature of CP [37], which involves vehicles interacting with each other repeatedly to determine positions. In particular, this cooperative process demands significant power and bandwidth [38]–[40], while also facing challenges in scheduling transmissions due to the intricate measurement and information fusion processes [41]-[43]. These factors may cause larger delays and scalability issues in cooperative localization [44], [45].

An emerging approach for cooperative vehicle localization is implicit cooperative positioning (ICP) [32], [46], which integrates GNSS and onboard passive sensor data through Bayesian-filtering, e.g., conventional extended Kalman filter (EKF) or message passing algorithm (MPA), to coherently fuse the measurements at different vehicles. In ICP, passive objects such as poles, road signs, or traffic lights, are cooperatively detected by multiple vehicles and exploited as noisy anchor points to enhance the vehicle location accuracy. In case of a centralized data-processing architecture gathering all vehicles' measurements, convergence can be achieved, but at the expense of high computational complexity. Standard MPA algorithms enable decentralized processing but are optimal only in case of Gaussian-linear models and acyclic factor graphs [47]-[50]. Recent solutions tried to limit the aforementioned problems by either performing fully-distributed

particle-based MPA between vehicles [34] or auto-adjusting the parameters of time-varying models [51]. Still, they rely on particle-based solutions which require high communication and computational loads which limit their scalability.

In recent years, there has been a growing reliance on machine learning (ML) tools to overcome the limits of conventional approaches, especially regarding scalability and non-linear models [52]–[55]. In particular, the reinforcement learning (RL) paradigm [56]–[58] and its deep learning (DL)based version [59]–[61] are notably effective in challenging single-agent Markov decision processes (MDPs) where labeled data are scarce or costly. They also excel in environments where the agent's actions directly impact the state of the environment and long-term rewards are prioritized [62]-[64]. Indeed, RL can be seen as a generalization of Bayesian filtering where the agents do not just predict the state through belief computation but also make decisions to optimize the cooperative process by maximizing long-term rewards, with a policy guiding the decision from state to action. RL is especially well-suited for complex scenarios with extensive state and action spaces, where deep neural networks (DNNs) can efficiently approximate the high-dimensional, nonlinear functions that represent such policies [59], [65]. This approach has been successfully applied in several fields, varying from rate and power control [66]–[69] to dynamic spectrum access in multi-user scenarios and efficient scheduling in vehicular networks [70]-[73].

In case more than one agent acts in the environment and the state is not directly observable, we categorize the framework as multi-agent RL (MARL) [74] and the system as decentralizedpartially observable MDP (Dec-POMDP) [75]-[77]. MARL involves independent agents whose actions influence each other's perception of the environment, and it is often solved with the usage of recurrent neural network (RNN), exploiting histories of observations and actions [78]. MARL algorithms, similarly to RL methods, can be divided into two categories: Q-learning and policy optimization (PO) (which comprises actor-critic methods) [79]-[81]. Q-learning focuses on estimating the long-term reward (i.e., Q-value) of each action, selecting the action with the highest Q-value and indirectly (i.e., not explicitly) formulating the policy [82]-[84]. On the other hand, PO directly optimizes the policy through the gradient of the total reward relative to policy parameters [85]-[88]. Multi-agent PO algorithms, especially when combined with a centralized agent learning and a decentralized execution of the policies (e.g., multi-agent proximal policy optimization (MAPPO) [85]), have shown remarkable performances with respect to Q-learning algorithms. This is mainly due to their being free of learning biases and improved sampling efficiency thanks to training guidelines like parameter sharing [89]–[91].

First attempts to employ MARL for CP, most of the literature works focus on target tracking for intelligent unmanned aerial vehicles (UAVs) [92] or agent scheduling for improving CP [93]. In [92], the RL objective was to maneuver the agents to track passive objects. However, they considered the state (i.e., the location) of the agents as known, while the main challenge is to estimate from the measurements their state jointly with target sensing. In [93], the agent state was

estimated with conventional MPA, while the RL objective was to activate links between agents to optimize cooperative positioning performances (i.e., by minimizing the positioning error bound (PEB)). The drawbacks of this method are that RL is not actively used for positioning but rather as an assistance method to MPA, and that they consider one agent only, i.e., a single link, at the time instead of exploiting the full potential of multi-agent systems (MASs).

Overall, the fundamental unresolved questions related to CP are as follows: i) how to design a decentralized MARL algorithm that simultaneously performs the computation of the agent state beliefs and the scheduling of the agent-to-agent communication resources, optimizing both location accuracy and communication efficiency; ii) what positioning accuracy improvement can be achieved with respect to state-of-the-art Bayesian approaches like ICP that exploit passive object detections between multiple agents; iii) what are the main tradeoffs between positioning improvement and communication resource optimization. Addressing these questions is mandatory for the employment in connected automated vehicles (CAVs), in particular to ensure scalability and handle real-word impairments encountered in vehicular scenarios. In this perspective, the goals of this paper are to develop agent-specific policies for communication scheduling between neighbors and, at the same time, learning a representation of the system dynamics that takes advantage of the selected neighbors' measurements. We propose a MARL-based ICP, a new paradigm in which PO RL algorithms are exploited to extend the conventional Bayesianfiltering approach incorporating the actions of the agents. The main idea is to learn from data the relation between agents' states and passive feature observations (see Fig. 1 for a visualization of the cooperative scenario) by selecting for the cooperation only those links to the neighbors that can provide a significant gain to the positioning accuracy. This approach is shown to not only improve the localization performance but also enhance the communication efficiency.

In this paper, we propose a new MARL algorithm, namely ICP-MAPPO, expressly designed for performing efficient distributed positioning through the MARL framework and extending the conventional Bayesian-filtering ICP to data-driven approaches. The key contributions are as follows:

- We revise the ICP Bayesian-filtering approach analyzing the current limitations and investigating more general frameworks for solution, drawing from the Dec-POMDP system model and MARL methods.
- We reformulate the ICP methodology into a MARL problem and we design the new ICP-MAPPO solution, relying on dynamic-decentralized-execution and training schemes to simultaneously optimize the Bayesian-filtering and MARL objectives.
- We validate the proposed ICP-MAPPO approach in a realistic C-ITS scenario simulated with CARLA [94], where CAVs perform CP by exploiting passive targets, i.e., poles, distributed over the scene.
- We perform a comparison with the state-of-the-art ICP algorithm [32] and single-agent-based algorithms. We prove the superior performances of the proposed algorithm both in terms of positioning error and communication efficiency.

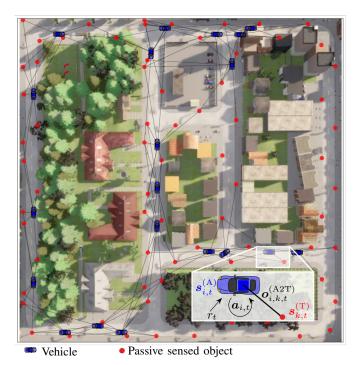


Fig. 1. Cooperative positioning scenario with twenty vehicles (blue vehicle icons), sensed poles acting as ancors (red circles) and detections (black lines).

TABLE I MAIN ABBREVIATIONS

| Acronym | Definition |
|-----------|--|
| A2A | Agent-to-agent |
| A2T | Agent-to-target |
| Dec-POMDP | Decentralized-partially observable Markov decision process |
| EKF | Extended Kalman Filter |
| ICP | Implicit cooperative positioning |
| LSTM | Long short-term memory |
| MLP | Multi-layer perceptron |
| MAPPO | Multi-agent proximal policy optimization |
| MARL | Multi-agent reinforcement learning |
| MPA | Message passing algorithm |

For easy reference, Table I lists the main abbreviations used throughout the paper.

The rest of this paper is structured as follows. Sec. II describes the system model of cooperative agents. Sec. III reviews the ICP Bayesian-filtering. Sec. IV presents the MARL framework and the proposed ICP-MAPPO execution and training schemes. Sec. V provides information about the simulated scenario and the results. Finally, Sec. VI draws the conclusions.

Notations

Random variables are displayed in sans serif, upright fonts; their realizations in serif, italic fonts. Vectors and matrices are denoted by bold lowercase and uppercase letters, respectively. For example, a random variable and its realization are denoted by \times and x; a random vector and its realization are denoted by

TABLE II LIST OF NOTATIONS

| Notation | Definition |
|---|---|
| N, K | Number of agents and passive objects |
| $\mathbf{s}_{i,t},\mathbf{a}_{i,t},\mathbf{o}_{i,t}$ | State, action and observation of agent i at time t |
| $oldsymbol{h}_{i,t}^{\mathrm{b}}, oldsymbol{h}_{i,t}^{\mathrm{V}}$ | History in belief and critic NNs of agent i at time t |
| $	au, 	au_t$ | Trajectory and transition at time t |
| r_t,R_t | Reward and reward-to-go at time t |
| $\pi_{\boldsymbol{\theta}}, V_{\boldsymbol{\phi}}, b_{\boldsymbol{\psi}}$ | Actor, critic and beliefs NNs |
| H, L_{τ} | Horizon and trajectory length |
| $A_{i,t}$ | Advantage function of agent i at time t |
| α , β , ϵ | Entropy, reward and clipping coefficients |
| γ , μ | Discount factor and learning rate |

x and x; a random matrix and its realization are denoted by **X** and X, respectively. Random sets and their realizations are denoted by up-right sans serif and calligraphic font, respectively. For example, a random set and its realization are denoted by X and \mathcal{X} , respectively. The function $p_{\mathsf{x}}(x)$, and simply p(x) when there is no ambiguity, denotes the probability density function (PDF) of x. Notations \mathbf{X}^{T} , \mathbf{X}^{*} and \mathbf{X}^{H} indicate the matrix transposition, conjugation and conjugate transposition. With the notation $x \sim \mathcal{N}(\mu, \sigma^2)$ we indicate a Gaussian random variable x with mean μ and standard deviation σ , whose PDF is denoted by $\mathcal{N}(x; \mu, \sigma^2)$. We use $\mathbb{E}\{\cdot\}$ and $\mathbb{V}\{\cdot\}$ to denote the expectation and the variance of a random variable, respectively. \mathbb{R} and \mathbb{C} stand for the set of real and complex numbers, respectively. Finally, we define with $blockdiag(\cdot)$ the block diagonal matrix whose diagonal contains the input blocks matrices.

Notations and definitions of important quantities used in the paper are summarized in Table II.

II. SYSTEM MODEL

We consider a vehicular network where a set of N vehicles engage in cooperative localization as depicted in Fig. 1. The connectivity graph for vehicle cooperation at time t is $\mathcal{G}_t = (\mathcal{V}, \mathcal{E}_t)$, with $\mathcal{V} = \{1, 2, \dots, N\}$ representing the set of agents (vehicles), and \mathcal{E}_t the edges (communication links) among them. Each agent $i \in \mathcal{V}$ in the network at time t has a set of neighbors $\mathcal{N}_{i,t}$, and it is assigned a state $\mathbf{s}_{i,t}^{(A)} = \left[\mathbf{u}_{i,t}^{(A)^{\top}} \mathbf{v}_{i,t}^{(A)^{\top}}\right]^{\top}$, where $\mathbf{u}_{i,t}^{(A)}$ and $\mathbf{v}_{i,t}^{(A)}$ are the 2D position and velocity vectors, respectively, defined in a global coordinate system. We denote with $\mathbf{s}_{t}^{(A)} = \left[\mathbf{s}_{i,t}^{(A)}\right]_{i=1}^{N}$ the aggregate state of all the vehicles at time t. The kinematic state transition of vehicle i at time t is modelled as $\mathbf{s}_{i,t}^{(\mathrm{A})} = f^{(\mathrm{A})}\big(\mathbf{s}_{i,t-1}^{(\mathrm{A})}, \mathbf{w}_{i,t-1}^{(\mathrm{A})}\big)$

$$\mathbf{s}_{i,t}^{(A)} = f^{(A)} \left(\mathbf{s}_{i,t-1}^{(A)}, \mathbf{w}_{i,t-1}^{(A)} \right)$$
 (1)

where $f^{(\mathrm{A})}(\cdot)$ is a nonlinear function that governs the dynamics of the vehicle's state and $\mathbf{w}_{i,t-1}^{(\mathrm{A})}$ represents the driving noise process, incorporating the uncertainty in motion. The model in (1) is associated to a state-transition PDF denoted as $T\left(\boldsymbol{s}_{i,t}^{(\mathrm{A})}|\boldsymbol{s}_{i,t-1}^{(\mathrm{A})}\right)\triangleq p\left(\boldsymbol{s}_{i,t}^{(\mathrm{A})}|\boldsymbol{s}_{i,t-1}^{(\mathrm{A})}\right).$

The scenario includes a set $\mathcal{F} = \{1, 2, \dots, K\}$ of K static and passive objects (or targets, denoted as red circles in Fig. 1) that vehicles can detect and localize by on-board sensors. To facilitate detection by vehicle sensors, specific objects easily identifiable and suitable for the purpose should be used. In this study, poles have been selected due to their ubiquity (especially in urban areas), ease of recognition, and fixed nature. Each pole k is described by a 2D position state $\mathbf{s}_{k,t}^{(T)}$, which is assumed to be constant over time. As before, we denote with $\mathbf{s}_t^{(\mathrm{T})} = \left[\mathbf{s}_{k,t}^{(\mathrm{T})}\right]_{k \in \mathcal{F}}$ the aggregate state of all passive objects at time t.

Vehicles are equipped with three distinct types of sensors. The first is a GNSS receiver, providing an estimate of the vehicle's state $\mathbf{s}_{i,t}^{(\mathrm{A})}$, modelled as

$$\mathbf{o}_{i,t}^{(\mathrm{GNSS})} = \mathbf{H}^{(\mathrm{GNSS})} \, \mathbf{s}_{i,t}^{(\mathrm{A})} + \mathbf{n}_{i,t}^{(\mathrm{GNSS})}$$
 (2)

vehicle's state $\mathbf{s}_{i,t}^{\circ}$, modelled as $\mathbf{o}_{i,t}^{(\mathrm{GNSS})} = \boldsymbol{H}^{(\mathrm{GNSS})} \mathbf{s}_{i,t}^{(\mathrm{A})} + \mathbf{n}_{i,t}^{(\mathrm{GNSS})}$ (2) where $\mathbf{n}_{i,t}^{(\mathrm{GNSS})} \sim \mathcal{N} \left(\mathbf{0}_{2 \times 2}, \boldsymbol{R}_{i,t}^{(\mathrm{GNSS})} \right) \in \mathbb{R}^{2 \times 1}$ is a zeromean Gaussian noise with covariance $\boldsymbol{R}_{i,t}^{(\mathrm{GNSS})} = \sigma^{(\mathrm{GNSS})^2} \boldsymbol{I}_2$, and $\boldsymbol{H}^{(\mathrm{GNSS})} = [\boldsymbol{I}_2 \, \mathbf{0}_{2 \times 2}] \in \mathbb{R}^{2 \times 4}$. From (2), we define the GNSS likelihood function as $p(\boldsymbol{o}_{i,t}^{(\mathrm{GNSS})} | \boldsymbol{s}_{i,t}^{(\mathrm{A})})$, and with $\boldsymbol{o}_{t}^{(\mathrm{GNSS})} = [\boldsymbol{o}_{i,t}^{(\mathrm{GNSS})}]_{i=1}^{N}$ the aggregate GNSS measurements of all the vehicles at time t. of all the vehicles at time t.

The second sensor refers to an active sensing technology for sidelink positioning offering relative agent-to-agent (A2A) location measurements for any pair of vehicles $(i, j) \in \mathcal{E}_t$

$$\mathbf{o}_{i,j,t}^{(\mathrm{A2A})} = \boldsymbol{H}^{(\mathrm{A2A})} \big(\mathbf{s}_{i,t}^{(\mathrm{A})} - \mathbf{s}_{j,t}^{(\mathrm{A})} \big) + \mathbf{n}_{i,j,t}^{(\mathrm{A2A})}$$
(3)

 $\begin{aligned} \mathbf{o}_{i,j,t}^{(\mathrm{A2A})} &= \boldsymbol{H}^{(\mathrm{A2A})} \big(\mathbf{s}_{i,t}^{(\mathrm{A})} - \mathbf{s}_{j,t}^{(\mathrm{A})} \big) + \mathbf{n}_{i,j,t}^{(\mathrm{A2A})} \\ \text{where} \quad \boldsymbol{H}^{(\mathrm{A2A})} &= \quad [\boldsymbol{I}_2 \, \mathbf{0}_{2 \times 2}] \quad \in \quad \mathbb{R}^{2 \times 4} \quad \text{and} \quad \mathbf{n}_{i,j,t}^{(\mathrm{A2A})} \end{aligned}$ $\mathcal{N}(\mathbf{0}_{2\times 2}, \boldsymbol{R}_{i,j,t}^{(\mathrm{A2A})})$ is a zero-mean Gaussian noise with covariance $\boldsymbol{R}_{i,j,t}^{(\mathrm{A2A})} = \sigma^{(\mathrm{A2A})^2} \boldsymbol{I}_2$. Additionally, agents have the capability to communicate with their neighbors to share location-

The third sensor type is a passive technology (e.g., RADAR, LIDAR, camera, or any combination), used by vehicle i to detect a set of passive objects $\mathcal{F}_{i,t} \subseteq \mathcal{F}$ in proximity at time t, and produce agent-to-target (A2T) measurements for each object $k \in \mathcal{F}_{i,t}$ as

$$\mathbf{o}_{i,k,t}^{(\mathrm{A2T})} = \boldsymbol{H}^{(\mathrm{A2T})} \mathbf{s}_{i,t}^{(\mathrm{A})} - \mathbf{s}_{k,t}^{(\mathrm{T})} + \mathbf{n}_{i,k,t}^{(\mathrm{A2T})}$$
(4)

where $m{H}^{(ext{A2T})} = [m{I}_2 \, m{0}_{2 imes 2}] \in \mathbb{R}^{2 imes 4}$ and $m{n}_{i,k,t}^{(ext{A2T})}$ $\mathcal{N}ig(\mathbf{0}_{2 imes2}, oldsymbol{R}_{i,k,t}^{(\mathrm{A2T})}ig)$ is a zero-mean Gaussian noise with covariance $oldsymbol{R}_{i,k,t}^{(\mathrm{A2T})} = \sigma^{(\mathrm{A2T})^2}oldsymbol{I}_2$.

 $pig(oldsymbol{o}_{i,j,t}^{(ext{A2A})}|oldsymbol{s}_{i,t}^{(ext{A})},oldsymbol{s}_{j,t}^{(ext{A})}ig)$ denote with $poldsymbol{(o_{i,k,t}^{(ext{A2T})}|oldsymbol{s}_{i,t}^{(ext{A})},oldsymbol{s}_{k,t}^{(ext{T})})}$ the A2A and A2T likeli $p(\mathbf{o}_{i,k,t} \mid \mathbf{s}_{i,t}, \mathbf{s}_{k,t})$ the A2A and A2T likelihoods, respectively. Moreover, we denote with $\mathbf{o}_{i,t} = \begin{bmatrix} \mathbf{o}_{i,t}^{(\mathrm{GNSS})^{\top}} \mathbf{o}_{i,t}^{(\mathrm{A2A})^{\top}} \mathbf{o}_{i,t}^{(\mathrm{A2T})^{\top}} \end{bmatrix}^{\top}$ the vector of all available measurements of vehicle i at time t, where $\mathbf{o}_{i,t}^{(\mathrm{A2A})} = \begin{bmatrix} \mathbf{o}_{i,j,t}^{(\mathrm{A2A})} \end{bmatrix}_{j \in \mathcal{N}_{i,t}}$ and $\mathbf{o}_{i,t}^{(\mathrm{A2T})} = \begin{bmatrix} \mathbf{o}_{i,k,t}^{(\mathrm{A2T})} \end{bmatrix}_{k \in \mathcal{F}_{i,t}}$. The total number of unique A2A and A2T measurements at time t is defined as $N_t^{(\mathrm{A2A})} = \sum_{i=1}^N |\mathcal{N}_{i,t}|$ and $N_t^{(\mathrm{A2T})} = \sum_{i=1}^N |\mathcal{F}_{i,t}|$, respectively. Note that the A2A measurements are not subject to measurement original measurements are not subject to measurement-origin uncertainty, i.e., it is not requested to perform any data association algorithm for pairing them, as the enabling technology is assumed to be active. On the other hand, the A2T observations are unlabelled, as it is unknown which object gives rise to a measurement, being them produced by a passive sensing technology (e.g., RADAR or LIDAR). In

this work, we assume that data association has already been performed at the vehicles (using, e.g., methods [53]) and that each A2T measurement has been correctly labeled with the originating target. We consider perfect data association as we aim to derive the best-case performances on the achievable accuracy of data-driven ICP and compare it with conventional Bayesian ICP in the same conditions. Interested readers can refer to [46] for details on data association and their impact on inference algorithms.

III. BAYESIAN FILTERING

In this section, we describe the Bayesian filtering solution, under the ICP framework, and then we highlight its main drawbacks and improvements.

A. Centralized Implicit Cooperative Positioning

The objective of ICP is to concurrently estimate the state of all vehicles and passive objects in the network. To this aim, we define the set of all available measurements at time t as

where
$$\mathbf{o}_t = H \mathbf{s}_t + \mathbf{n}_t$$
 (5) where $\mathbf{o}_t = \left[\mathbf{o}_{i,t}\right]_{i \in \mathcal{V}} \in \mathbb{R}^{\left(2N+2N_t^{(\mathrm{A2A})}+2N_t^{(\mathrm{A2T})}\right) \times 1}$, H is the matrix modeling the relation to the states, defined as in [32], and $\mathbf{s}_t = \left[\mathbf{s}_t^{(\mathrm{A})^{\top}}\mathbf{s}_t^{(\mathrm{T})^{\top}}\right]^{\top} \in \mathbb{R}^{\left(4N+2K\right) \times 1}$ is the aggregated state of the system. $\mathbf{n}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_t)$ is the overall measurement noise with covariance $\mathbf{R}_t = \mathrm{blockdiag}(\mathbf{R}_t^{(\mathrm{GNSS})}, \mathbf{R}_t^{(\mathrm{A2A})}, \mathbf{R}_t^{(\mathrm{A2T})})$, where $\mathbf{R}_t^{(\mathrm{GNSS})} = \mathrm{blockdiag}(\mathbf{R}_{1,t}^{(\mathrm{GNSS})}, \dots, \mathbf{R}_{N,t}^{(\mathrm{GNSS})})$, $\mathbf{R}_t^{(\mathrm{A2A})} = \mathrm{blockdiag}(\mathbf{R}_{1,t}^{(\mathrm{A2A})}, \dots, \mathbf{R}_{N_t^{(\mathrm{A2A})},t}^{(\mathrm{A2A})})$ with the ℓ -th entry given by $\mathbf{R}_{\ell,t}^{(\mathrm{A2A})} = \mathbf{R}_{i_\ell,j_\ell,t}^{(\mathrm{A2A})}$, and $\mathbf{R}_t^{(\mathrm{A2T})} = \mathrm{blockdiag}(\mathbf{R}_{1,t}^{(\mathrm{A2T})}, \dots, \mathbf{R}_{N_t^{(\mathrm{A2T})},t}^{(\mathrm{A2T})})$ with

The overall state estimate \hat{s}_t is obtained through the minimum mean square error (MMSE) estimator as

$$\widehat{\mathbf{s}}_t = \mathbb{E}\{\mathbf{s}_t|\mathbf{o}_{1:t}\} = \int \mathbf{s}_t \, p(\mathbf{s}_t|\mathbf{o}_{1:t}) \, d\mathbf{s}_t$$
 (6)

where $o_{1:t} = \left[o_{t'}\right]_{t'=1}^{t}$ is the set of all aggregated measurements up to time t and $p(s_t|o_{1:t})$ is the posterior PDF defined

$$p(\mathbf{s}_t|\mathbf{o}_{1:t}) \propto p(\mathbf{o}_t|\mathbf{s}_t) \int p(\mathbf{s}_t|\mathbf{s}_{t-1}) p(\mathbf{s}_{t-1}|\mathbf{o}_{1:t-1}) d\mathbf{s}_{t-1}$$
. (7)

We denote with $b(s_{i,t}|o_{1:t}) \triangleq p(s_{i,t}|o_{1:t})$ the marginal posterior PDF, also called belief of agent i. Given that all the measurements are mutually independent, the likelihood function of s_t is computed as

$$p(\boldsymbol{o}_t|\boldsymbol{s}_t) = p(\boldsymbol{o}_t^{(\text{GNSS})}|\boldsymbol{s}_t^{(\text{A})}) \prod_{i=1}^{N} \prod_{j \in \mathcal{N}_{i,t}} p(\boldsymbol{o}_{i,j,t}^{(\text{A2A})}|\boldsymbol{s}_{i,t}^{(\text{A})}, \boldsymbol{s}_{j,t}^{(\text{A})})$$

$$\times \prod_{i=1}^{N} \prod_{k \in \mathcal{F}_{i,t}} p(\boldsymbol{o}_{i,k,t}^{(\mathrm{A2T})} | \boldsymbol{s}_{i,t}^{(\mathrm{A})}, \boldsymbol{s}_{k,t}^{(\mathrm{T})}). \tag{8}$$

For notation purposes, we will denote the likelihood function also as $O(o_t|s_t) \triangleq p(o_t|s_t)$. In case the dynamic and

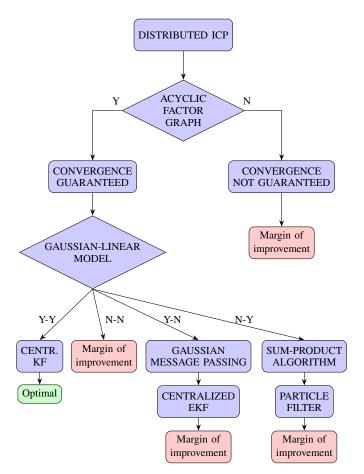


Fig. 2. Convergence conditions in ICP methods.

measurements models in (1) and (5), respectively, are linear and with a Gaussian noise, the state estimate in (6) reduces to a Kalman filter (KF) as described in [32], [46], with efficient resolution in matrix form.

B. Limitations of Bayesian ICP Methods

The centralized ICP approach is impractical for extensive networks due to the following major limitations: the single central computing unit representing a point of failure, and its computational complexity growing cubically with the number of vehicles and passive objects [32]. To overcome such limitations, distributed or consensus-based ICP algorithms have been studied in the past [34]. However, their convergence to the centralized solution is guaranteed only in acyclic (i.e., treestructured) factor graphs. Moreover, even in case of convergence, the result would be optimal only with Gaussian and linear models (i.e., in (1) and (5)). In all the other cases, optimality is not guaranteed. In Fig. 2 we summarized all cases and highlighted those where improvements could be provided by new data-driven designs. We point out that in real-world dynamics, the factor graph is usually not acyclic and the models are typically neither Gaussian nor linear.

The aim of this paper is to address the gap by proposing a new decentralized data-driven solution to the ICP problem suited for non-linear non-Gaussian models, overcoming the limits of parametric Bayesian implementations based on EKF or particle filter (PF) highlighted in Fig. 2. The proposed distributed method also incorporates a data-driven optimization of the cooperation graph by making the agents actively and opportunistically select the cooperating neighbors so as to minimize the communication signaling. In particular, to address the limitations of conventional ICP solutions, we adopt neural networks (NNs)-based models, which are able to learn whatever non-linear function is hidden in the data thanks to the universal approximation theorem. Specifically, a RNN learns the non-linear motion and measurement models, whereas a multi-layer perceptron (MLP) learns the non-linear relation between link activation and state estimate. Moreover, NNs have proven effective even in non-Gaussian settings [53], given their ability to model complex probability distributions without assuming any specific form. The centralized ICP method reviewed in this section will be used as a benchmark to assess the proposed method.

IV. MARL FOR COOPERATIVE POSITIONING

In this section, we first introduce the MARL framework (Sec. IV-A) that will be used later for the design of the ICP-MAPPO solution (Sec. IV-B). The ICP-MAPPO execution and training schemes are reported in Sec. IV-C and IV-D, respectively.

A. MARL Framework

model the cooperative MAS finitedefined horizon Dec-POMDP [75] by the tuple $\langle \mathcal{V}, \mathcal{S}, \mathcal{A}, T_0, T, \mathcal{O}, O, R, \gamma, H \rangle$. We recall that the set \mathcal{V} refers to the cooperative agents, while the sets S and \mathcal{A} denote the state and action spaces, respectively. T_0 is the initial state distribution at time t = 0, while $T(s_t|s_{t-1},a_t) \triangleq p(s_t|s_{t-1},a_t)$ is the state transition PDF that, differently from the Bayesian-filtering system model in Sec. II, now also includes the joint action realization $m{a}_t = [m{a}_{i,t}]_{i \in \mathcal{V}} \in \mathcal{A}$ and the joint state $m{s}_t \in \mathcal{S}$. At each time t, the agents receive the joint observations or measurements $o_t \in \mathcal{O}$ which are sampled from the distribution $O(o_t|a_{t-1},s_t) \triangleq p(o_t|a_{t-1},s_t)$. Note that here, (8) is also function of the previous joint action of the agents a_{t-1} , thus generalizing the concept of Bayesian-filtering. $R(\mathbf{s}_t, \mathbf{a}_t) = \mathbf{r}_t \in \mathbb{R}$ denotes the instantaneous shared reward at time t obtained from the reward function R, while $\gamma \in [0,1)$ and H are the discount factor and time horizon of each episode, respectively.

Since the states and rewards are not directly observable by the agents (partially observable MDP), each agent i keeps track of the so-called *histories* defined as $\boldsymbol{h}_{i,1:t} = \boldsymbol{h}_{i,t} = \left[(\boldsymbol{a}_{i,t'-1}, \boldsymbol{o}_{i,t'}) \right]_{t'=1}^t$. Note that the histories are a generalization of the aggregated measurements up to time t in (6). Given a new observation $\boldsymbol{o}_{i,t}$, the state estimates $\hat{\boldsymbol{s}}_{i,t}$ are produced by MMSE criterion from the belief PDF $b_{\psi}(\boldsymbol{s}_{i,t}|\boldsymbol{o}_{i,t},\boldsymbol{a}_{i,t-1},\boldsymbol{h}_{i,t-1}) = p_{\psi}(\boldsymbol{s}_{i,t}|\boldsymbol{o}_{i,t},\boldsymbol{a}_{i,t-1},\boldsymbol{h}_{i,t-1})$ parameterized by ψ . Moreover, agents adopt a policy $\pi_{\theta}(\boldsymbol{a}_{i,t}|\boldsymbol{h}_{i,t}) = p_{\theta}(\boldsymbol{a}_{i,t}|\boldsymbol{h}_{i,t})$ defined by θ to obtain the action $\boldsymbol{a}_{i,t}$ from histories $\boldsymbol{h}_{i,t}$. A full comparison between Bayesian filtering and RL (i.e., its generalized version) can be found in

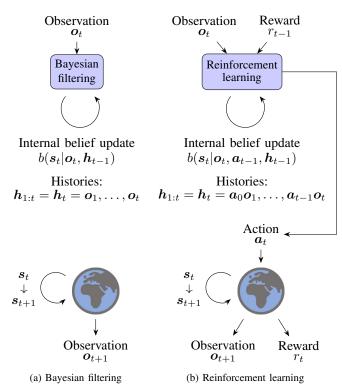


Fig. 3. Comparison between Bayesian filtering and RL.

Fig. 3. By defining the *reward-to-go* $R_t = \sum_{t'=t}^{H-1} \gamma^{t'-t} r_{t'}$ as the cumulative discounted reward from time t to the end of the episode, the objective of the MARL problem is to maximize, over the policy π , the expected cumulative discounted reward from the beginning of the episode

$$\max_{\pi} J(\pi) = \max_{\pi} \mathbb{E}\{\mathsf{R}_0\} \tag{9}$$

 $\max_{\pi} J(\pi) = \max_{\pi} \mathbb{E}\{\mathsf{R}_0\}$ which usually translates into optimizing the parameters of the policy as $\theta^* = \operatorname{argmax}_{\theta} J(\pi_{\theta})$, with π_{θ}^* representing the optimal policy.

B. MARL Solution to the ICP Problem

In standard Dec-POMDP, each agent only knows its local actions and observations, thus resulting in possible nonstationary learning problems from each agent's perspective [96]. By training independent learners to optimize the team reward (i.e., concurrent learning), we induce a change in the dynamics of the environment as teammates continuously adapt their behaviours throughout learning. On the contrary, whenever a fully connected graph with communications is present, the Dec-POMDP collapses to a centralized POMDP, resulting in higher complexity and communication inefficiencies [89], [97], exactly as in centralized ICP. To solve the issues of independent and centralized training-execution, the state-of-the-art works exploit the so called centralized-training and decentralized-execution paradigm. This framework permits to learn the policies in a centralized way and then deploy them in the network graph for decentralized execution [85], [87], [98].

While this approach solves the problem in standard MARL algorithms, in the context of ICP, having access to the neighbors' measurements would allow the positioning accuracy to be significantly improved. Indeed, the objective of ICP is to minimize over the belief b the error on the state estimate as

$$\min_{b} J(b) = \min_{b} \mathbb{E} \left\{ \sum_{t} \left\| \mathbf{s}_{t} - \widehat{\mathbf{s}}_{t} \right\|_{2}^{2} \right\}. \tag{10}$$

Therefore, we here propose to define as actions the agent's selection of the communication links to the neighbors to cooperate with. This allows to optimize the communication efficiency with respect to the centralized solution. Formally, we define the following Dec-POMDP:

- 1) Agents: The agent is identified by vehicle $i \in \mathcal{V}$ that composes the connected network.
- 2) Actions: The action of agent i at time t is $\mathbf{a}_{i,t} = \left[\mathbf{a}_{i,j,t}\right]_{j=1}^{N}$, where $\mathbf{a}_{i,j,t} \in \{0,1\}$ represents the Boolean decision of agent i to communicate with agent j.

 3) States: Only the states of the vehicles $\mathbf{s}_{t}^{(A)}$ are consid-
- ered, while the target states $\mathbf{s}_{t}^{(\mathrm{T})}$ are implicitly learned by the NNs through the hidden features. Indeed, the system does not output or keep track of the states of the targets, since they are not needed as in the ICP Bayesian filtering formulation. In other words, the ICP-MAPPO model just outputs the predicted states of the agents, while the targets' states are contained in the hidden space, i.e., histories. Therefore, from now on, we indicate with \mathbf{s}_t the state of the agents $\mathbf{s}_t^{(A)}$.
- 4) Observations: GNSS, A2A, and A2T measurements described in Sec. II are the observations used in the Dec-POMDP modeling, as they are the only output returned by the world at inference time.

During the centralized training, the agents learn the relation between histories-actions, i.e., policy optimization, and histories-states, i.e., belief optimization, while having access to the full observable state \mathbf{s}_t and measurements \mathbf{o}_t . Conversely, during the decentralized execution, the agents decide how to modify the network graph to achieve the best trade-off between positioning accuracy and communication efficiency. We call this approach centralized-training and dynamic-decentralizedexecution, as during execution, according to the agents' actions, the coordination graph may vary, passing from fullyconnected to fully-decentralized according to the agent's decisions.

C. ICP-MAPPO Execution Scheme

For belief and action prediction, we propose to employ long short-term memory (LSTM) and MLP, respectively. In Fig. 4, we show a compact representation of the execution within each agent. In particular, the NN functions are defined as

$$\hat{s}_{i,t}, h_{i,t}^{b} = b_{\psi}(s_{i,t}|o_{i,t}, \bar{a}_{i,t-1}, \bar{h}_{i,t-1}^{b})$$
 (11)

$$\boldsymbol{a}_{i,t} \sim \pi_{\boldsymbol{\theta}}(\boldsymbol{a}_{i,t}|\boldsymbol{h}_{i,t}^{\mathrm{b}})$$
 (12)

where $o_{i,t}$ is the ordered vector of all measurements of agent i at time t defined as in Sec. II, $\bar{a}_{i,t} = \left[\bar{a}_{i,j,t}\right]_{j=1}^N$ includes the sampled actions from the policy distribution adjusted with the feasibility of the network connectivity as

$$\bar{a}_{i,j,t} = \begin{cases} a_{i,j,t} & \text{if } j \in \mathcal{N}_{i,t} \\ -1 & \text{otherwise} \end{cases}$$
 (13)

and $\bar{h}_{i,t}^{\mathrm{b}}$ are the hidden features of the belief LSTM which contain a compressed representation of the histories of agent i and all selected neighbors at the previous timestep

$$\bar{\boldsymbol{h}}_{i,t}^{b} = \frac{\boldsymbol{h}_{i,t}^{b} + \sum_{j \in \mathcal{V}} \boldsymbol{h}_{j,t}^{b} \, \mathbb{1}(\bar{a}_{i,j,t} == 1)}{1 + \sum_{j \in \mathcal{V}} \mathbb{1}(\bar{a}_{i,j,t} == 1)}$$
(14)

where $\mathbb{1}(\cdot)$ is the indicator function that returns 1 if the condition is true and 0 otherwise. We point out that the hidden features $h_{i,t}^{\rm b}$ include not only past actions and measurements but also the implicit state estimates of the targets $\widehat{s}_t^{({\rm T})}$, which are never explicitly predicted by the system for output space complexity reduction.

The key rationale behind the proposed execution scheme is the following. We employ the average operation in (14) to avoid gradient divergence over the timesteps. Furthermore, the action decision at time t in (12) is mainly based on the previous timestep information $\bar{h}_{i,t-1}^{\text{b}}$, as there is no way for agent i to know a priori the measurements of its neighbors $h_{jt}^{b}, \forall j \in \mathcal{V}$, in order to activate the communications between them. Moreover, the actions $\bar{a}_{i,t}$ are given as input to the belief LSTM for two main reasons. First, the information about which agents were selected for measurements fusion is necessary to coherently predict the state estimate. Second, the negative action values imposed by the lack of possible connectivity permit each agent to implicitly learn its index or identification. In this way, the scalable and efficient parameter sharing approach for training one single NN [89], instead of agent-specific NNs, can be combined with agent differentiation by index learning.

D. ICP-MAPPO Training Scheme

For the reward definition, we propose to use a function that, looking at the future timestep, rewards the actions that gave a predetermined improvement β on the positioning accuracy. In other words, each agent i tries to answer the following question: if I had chosen agent j' instead of agent j, would the performances have improved? This is formalized as

$$r_{t} = \begin{cases} -1 & \text{if } \|\mathbf{s}_{t} - \widehat{\mathbf{s}}_{t}\|_{2}^{2} - \|\mathbf{s}_{t+1} - \widehat{\mathbf{s}}_{t+1}\|_{2}^{2} \leq -\beta \\ +1 & \text{if } \|\mathbf{s}_{t} - \widehat{\mathbf{s}}_{t}\|_{2}^{2} - \|\mathbf{s}_{t+1} - \widehat{\mathbf{s}}_{t+1}\|_{2}^{2} > \beta \\ +2 & \text{if } -\beta < \|\mathbf{s}_{t} - \widehat{\mathbf{s}}_{t}\|_{2}^{2} - \|\mathbf{s}_{t+1} - \widehat{\mathbf{s}}_{t+1}\|_{2}^{2} \leq \beta \end{cases}$$

$$(15)$$

where β is a hyper-parameter which regulates the improvement step. At the beginning of the learning, if the improvement is negative and bigger than β , the reward is negative as the actual agent selection worsen the positioning accuracy. On the other hand, if the improvement is positive and greater than β , the reward is +1. Finally, when the learning starts converging and the improvements become smaller, we introduce a long-term reward of +2. Note that, while in conventional Dec-POMDPs the reward directly depends on the actions, in the proposed system the effect of the actions' choice can be assessed only at the next timestamp and only by measuring the positioning error.

Regarding the type of MARL algorithm, we opted for PO over Q-learning-based methods. This is because Q-learning

algorithms combined with DL have no guarantees of convergence and retain a lot of bias (i.e., inaccurate state-action value or Q-value). On the contrary, PO algorithms retain very low bias since they directly optimize the objective function in (9) and have been proven to outperform Q-learning methods in MARL systems [87]. Moreover, while off-policy RL algorithms use historical data to learn the policy, in the context of CP, where state estimation is crucial, it is essential to utilize the most up-to-date policy available since the action sampling (i.e., radio link activation) directly influences the positioning performances. Despite PO algorithms having an intrinsic high variance, i.e., they require a lot of samples to converge, this can be mitigated by the learning of the value function, either $V^{\pi}(s_t)$ or $Q^{\pi}(s_t, a_t)$, which estimates the long-term reward given a specific state or state-action pair, respectively. Specifically, we employ the state value function defined as

$$V^{\pi}(\mathbf{s}_t) = \mathbb{E}\{\mathsf{R}_t|\mathsf{s}_t = \mathbf{s}_t\}$$

= $\mathbb{E}_{\mathbf{a}_t \sim \pi, \mathbf{s}_{t+1} \sim T} \Big\{ R(\mathbf{s}_t, \mathbf{a}_t) + \gamma V^{\pi}(\mathbf{s}_{t+1}) \Big\}.$ (16)

Usually, $V^{\pi}(s_t)$ cannot be directly computed due to the curse of dimensionality and thus it is estimated by an additional NN $\hat{V}_{\phi}(s_t) = V_{\phi}(s_t)$, with parameters ϕ which are only employed during training.

In standard single-agent RL frameworks, the policy optimization problem is usually defined with the introduction of trajectories $\mathbf{\tau} = (\mathbf{s}_0, \mathbf{a}_0, \dots, \mathbf{s}_H, \mathbf{a}_H)$ by maximizing

$$J(\pi_{\boldsymbol{\theta}}) = \mathbb{E}_{\boldsymbol{\tau} \sim p(\boldsymbol{\tau}|\pi_{\boldsymbol{\theta}})} \left\{ \widetilde{R}(\boldsymbol{\tau}) \right\}$$

$$= \sum_{t=0}^{H} \mathbb{E}_{\mathbf{s}_{t} \sim p(\boldsymbol{s}_{t}|\pi_{\boldsymbol{\theta}}), \mathbf{a}_{t} \sim \pi_{\boldsymbol{\theta}}(\boldsymbol{a}_{t}|\boldsymbol{s}_{t})} \left\{ \gamma^{t} R(\mathbf{s}_{t}, \mathbf{a}_{t}) \right\}$$
(17)

where $\widetilde{R}(\mathbf{\tau}) = \mathsf{R}_0$ is the reward of trajectory $\mathbf{\tau}$, $p(\mathbf{\tau}|\pi) = T_0 \prod_{t=0}^{H-1} T(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t) \, \pi(\mathbf{a}_t|\mathbf{s}_t)$ is the PDF of an H-step trajectory, and $p(\mathbf{s}_t|\pi)$ is the state marginal of the trajectory distribution induced by policy π . Standard REINFORCE PO algorithms [99] update the policy parameters in (17) in the direction of $\nabla_{\theta} J(\pi_{\theta})$, which can be written as (see Appendix A)

$$\nabla_{\boldsymbol{\theta}} J(\pi_{\boldsymbol{\theta}}) = \mathbb{E}_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim p(\mathbf{s}_{t}, \mathbf{a}_{t} | \pi_{\boldsymbol{\theta}})} \left\{ \sum_{t=0}^{H-1} \nabla_{\boldsymbol{\theta}} \log \left(\pi_{\boldsymbol{\theta}}(\boldsymbol{a}_{t} | \boldsymbol{s}_{t}) \right) A_{t} \right\}$$
(18)

where $p(\mathbf{s}_t, \mathbf{a}_t | \pi_{\theta})$ is the state-action marginal of the trajectory distribution induced by policy π and $A_t = A_t(\mathbf{s}_t, \mathbf{a}_t)$ is the generic advantage function at time t [100], which quantifies the convenience of taking a specific action \mathbf{a}_t in a given state \mathbf{s}_t , compared to the average action's expected return for that state.

During successive optimization steps of (18) within the same trajectory, where the objective is to maintain proximity between new and old policy parameters, even minor variations in the NN weights can lead to significant differences in performance. Consequently, a single unfavorable optimization step can drastically deteriorate the policy's effectiveness. Recent state-of-the-art methods, e.g., trust region policy optimization (TRPO) [101] and proximal policy optimization (PPO) [102], tried to solve this problem by taking the largest gradient

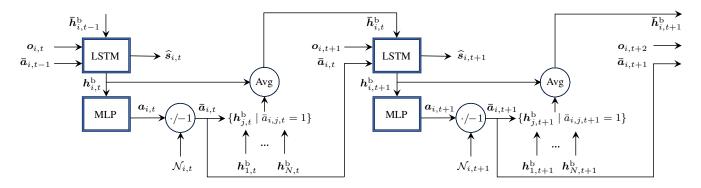


Fig. 4. Dynamic-decentralized-execution scheme of the proposed ICP-MAPPO algorithm.

step size possible to improve performance, while maintaining constraints on how close the new and old policies (i.e., $\pi_{\theta_{\rm old}}$ at previous train step) are allowed to be. The constraint in TRPO is enforced by Kullback–Leibler (KL) divergence and the parameters are obtained by maximizing the *surrogate* objective function as

$$\theta = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \ \mathbb{E}_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim p(\mathbf{s}_{t}, \mathbf{a}_{t} | \pi_{\boldsymbol{\theta}})} \left\{ \frac{\pi_{\boldsymbol{\theta}}(\boldsymbol{a}_{t} | \mathbf{s}_{t})}{\pi_{\boldsymbol{\theta}_{\text{old}}}(\boldsymbol{a}_{t} | \mathbf{s}_{t})} A_{t}(\mathbf{s}_{t}, \mathbf{a}_{t}) \right\}$$
s.t.
$$\mathbb{E}_{\mathbf{s}_{t} \sim p(\mathbf{s}_{t} | \pi_{\boldsymbol{\theta}})} \left\{ D_{\text{KL}} \left(\pi_{\boldsymbol{\theta}}(\cdot | \mathbf{s}_{t}) \| \pi_{\boldsymbol{\theta}_{\text{old}}}(\cdot | \mathbf{s}_{t}) \right) \right\} \leq \epsilon$$
(19)

which resulted in a second-order optimization method. On the contrary, PPO and its recent multi-agent version MAPPO use a much more efficient first-order method that exploits clipping to remove incentives for the new policy to get far from the old policy.

In this paper, we adopt three loss functions: $L(\phi)$ and $L(\theta)$ derived from the MAPPO scheme to train the state-value and policy NNs, respectively, and $L(\psi)$ to train the belief NN. π_{θ} and V_{ϕ} are called *actor* and *critic*, respectively, since the actor is responsible for selecting actions based on the current policy, and the critic evaluates the quality of these actions by estimating the value function. In Dec-POMDP, the critic V_{ϕ} is also dependent on the history of action-observation pairs and thus it is usually modelled with a RNN as

$$\widehat{V}_{\phi}(s_{i,t}, h_{i,t-1}^{V}), h_{i,t}^{V} = V_{\phi}(s_{i,t}, h_{i,t-1}^{V})$$
 (20)

where $h_{i,t}^{\rm V}$ are the hidden features of the critic. Given a trajectory of length L_{τ} (subset of the horizon length H), $L(\phi)$ is defined to perform regression on the rewards-to-go as

$$L(\boldsymbol{\phi}) = \frac{1}{NL_{\tau}} \sum_{i \in \mathcal{V}} \sum_{\ell=1}^{L_{\tau}} \left\{ \max \left(\left[\widehat{V}_{\boldsymbol{\phi}}(\boldsymbol{s}_{i,\ell}, \boldsymbol{h}_{i,\ell}^{V}) - R_{\ell} \right]^{2}, \right. \\ \left. \left[\operatorname{clip}(\widehat{V}_{\boldsymbol{\phi}}(\boldsymbol{s}_{i,\ell}, \boldsymbol{h}_{i,\ell-1}^{V}), \widehat{V}_{\boldsymbol{\phi}_{\text{old}}}(\boldsymbol{s}_{i,\ell}, \boldsymbol{h}_{i,\ell-1}^{V}), \epsilon) - R_{\ell} \right]^{2} \right) \right\}$$

$$(21)$$

where the clip prevents the value function from radically changing between iterations, and it is defined as

$$\operatorname{clip}(A, B, \epsilon) = \min\left(\max(A, B - \epsilon), B + \epsilon\right)$$
 where ϵ is the clip coefficient. (22)

The actor π_{θ} is also trained with clipping to discard the KL

constraint in (19) by minimizing

$$L(\boldsymbol{\theta}) = -\frac{1}{NL_{\tau}} \sum_{i \in \mathcal{V}} \sum_{\ell=1}^{L_{\tau}} \left\{ \min \left(\frac{\pi_{\boldsymbol{\theta}}(\boldsymbol{a}_{i,\ell} | \boldsymbol{h}_{i,\ell}^{b})}{\pi_{\boldsymbol{\theta}_{old}}(\boldsymbol{a}_{i,\ell} | \boldsymbol{h}_{i,\ell}^{b})} \widehat{A}_{i,\ell}, \right. \right.$$

$$\left. \operatorname{clip} \left(\frac{\pi_{\boldsymbol{\theta}}(\boldsymbol{a}_{i,\ell} | \boldsymbol{h}_{i,\ell}^{b})}{\pi_{\boldsymbol{\theta}_{old}}(\boldsymbol{a}_{i,\ell} | \boldsymbol{h}_{i,\ell}^{b})}, 1, \epsilon \right) \widehat{A}_{i,\ell} \right) + \alpha S \left(\pi_{\boldsymbol{\theta}}(\cdot | \boldsymbol{h}_{i,\ell}^{b}) \right) \right\}$$

$$(23)$$

where $\widehat{A}_{i,\ell} = R_\ell - \widehat{V}_{\phi_{\mathrm{old}}}(s_{i,\ell}, \boldsymbol{h}^{\mathrm{V}}_{i,\ell-1})$ is the advantage function estimate, $S(p_{\mathrm{x}}) = \mathbb{E}_{\mathrm{x} \sim p_{\mathrm{x}}} \Big\{ -\log \big(p_{\mathrm{x}}(x) \big) \Big\}$ is the entropy function which encourages the exploration by inducing stochastic policies, and α is the temperature hyper-parameter which balances the trade-off between exploiting the best actions and exploring new actions. Finally, the beliefs b_{ψ} adopt a MSE loss function to minimize J(b) in (10) as

$$L(\psi) = \frac{1}{NL_{\tau}} \sum_{i \in \mathcal{V}} \sum_{\ell=1}^{L_{\tau}} \|\widehat{\mathbf{s}}_{i,t} - \mathbf{s}_{i,t}\|_{2}^{2}.$$
 (24)

All the NNs are trained with maximum likelihood estimation (MLE) criterion. However, while $b_{\psi}(s_{i,t}|o_{i,t},\bar{a}_{i,t-1},\bar{h}_{i,t-1}^{\mathrm{b}})$ directly outputs $\hat{s}_{i,t}$, $\pi_{\theta}(a_{i,t}|h_{i,t}^{\mathrm{b}})$ predicts the probability of communication among agents through sigmoid activation functions, from which actions $a_{i,t}$ are sampled. The full training algorithm can be found in Algorithm 1, where we defined a transition as $\tau_t = (s_t, o_t, h_t^{\mathrm{b}}, \bar{h}_t^{\mathrm{b}}, h_t^{\mathrm{V}}, a_t, \bar{a}_t, r_t, s_{t+1}, o_{t+1}, \hat{s}_{t+1})$. Since our approach combines the usage of passive targets to improve the position estimate and MAPPO MARL to perform an efficient agent selection, we call this algorithm ICP-MAPPO.

The main characteristics of ICP-MAPPO are the following. ICP-MAPPO is a low-bias on-policy algorithm since the data used to train the agents are collected from the policy currently being learned or improved. For value regression, we adopted a centralized value function that takes as input extra global information (i.e., the states) not present in the agent's local observation to accurately estimate the values state. The beliefs are computed as in model-based value estimation (MBVE) RL [103], [104], leveraging the learned dynamics to predict the state estimate. This additionally reduces the variance of the PO method without introducing additional biases by avoiding performing rollouts [105]. Finally, as opposed to conventional MARL algorithms, the rewards are not directly dependent on the action, but only implicitly through the beliefs of the next

Algorithm 1 Implicit Cooperative Positioning Multi-Agent Proximal Policy Optimization (ICP-MAPPO)

```
1: Input: actor, critic and belief parameters \theta = \theta_{\rm old},
        \phi = \phi_{\rm old}, and \psi.
 2: for each training step n = 1, 2, ..., N_{\text{step}} do
                Initialize empty batch \mathcal{B} = \{\} and trajectory \boldsymbol{\tau} = []
 3:
                Initialize histories h_{i,0}^{V} and h_{i,1}^{b} for critic and beliefs
 4:
                Initialize state estimate \hat{s}_0
  5:
               for t = 1, 2, ..., H do
  6:
                        for all agents i \in \mathcal{V} in parallel do
  7:
                                Sample action \boldsymbol{a}_{i,t} \sim \pi_{\boldsymbol{\theta}_{\mathrm{old}}}(\boldsymbol{a}_{i,t}|\boldsymbol{h}_{i,t}^{\mathrm{b}})
  8:
                               Send \boldsymbol{h}_{i,t}^{\mathrm{b}} and receive \boldsymbol{h}_{i,t}^{\mathrm{b}} \ \forall j \in \mathcal{N}_{i,t}
  9:
                               Get value estimate \widehat{V}_{\phi_{\mathrm{old}}}(s_{i,t}, \boldsymbol{h}_{i,t-1}^{\mathrm{V}}) with (20) Compute \bar{\boldsymbol{a}}_{i,t} and \bar{\boldsymbol{h}}_{i,t}^{\mathrm{b}} with (13) and (14)
10:
11:
                                Observe s_{i,t+1}, o_{i,t+1}
12:
                                Get state estimate \hat{s}_{i,t+1} with (11)
13:
                        end for
14:
                        Observe r_t and store \tau_t in \tau
15:
16:
                Compute advantage estimate \widehat{A}_{i,t} \ \forall t and agent i on \boldsymbol{\tau}
17:
                Compute reward-to-go R_t for each \forall t on \tau
18:
                Split trajectory 	au into chunks of length L_{	au}
19:
               \begin{aligned} & \textbf{for each } \ell = 0, 1, \dots, \lfloor H/L_{\tau} \rfloor \, \textbf{do} \\ & \mathcal{B} = \mathcal{B} \cup \left\{ \boldsymbol{\tau}_t, \widehat{A}_t, R_t \right\}_{t=\ell}^{\ell+L_{\tau}} \\ & \text{Adam update of } \boldsymbol{\psi} \text{ on } L(\boldsymbol{\psi}) \text{ with data } \left\{ \boldsymbol{\tau}_t \right\}_{t=\ell}^{\ell+L_{\tau}} \end{aligned}
20:
21:
22:
23:
                end for
                for each mini-batch do
24:
                       Sample \{\tau_\ell\}_{\ell=1}^{L_\tau} \sim \mathcal{B}
Adam update of \theta on L(\theta) with data \{\tau_\ell\}_{\ell=1}^{L_\tau}
25:
26:
                       Adam update of \phi on L(\phi) with data \{\tau_\ell\}_{\ell=1}^{L_\tau}
27:
28:
                \boldsymbol{\theta}_{\mathrm{old}} \leftarrow \boldsymbol{\theta}, \, \boldsymbol{\phi}_{\mathrm{old}} \leftarrow \boldsymbol{\phi}
29:
30: end for
```

timestep. This permits to effectively decouple the evaluation of actions based on the improvement of state predictions rather than immediate outcomes, focusing on long-term strategic benefits rather than short-term gains.

V. SIMULATION EXPERIMENTS

In this section, we first introduce the scenario and the training procedures, and then we describe the baseline methods, and the main simulation results.

A. Simulation Setup

To evaluate the performances of the proposed ICP-MAPPO algorithm, we simulate a C-ITS scenario with the CARLA software [94] in an urban map (i.e., Town02 of CARLA) that spans an area of $200 \times 200 \, \mathrm{m}^2$. Fig. 1 shows a bird-eyeview representation of the map. CARLA takes into account inter-vehicle dynamics, such as acceleration, braking behavior, and collision physics, as well as communication constraints given by the environment. Within the area, 20 CAVs move for 1500 timesteps sampled every 0.2 s, while 72 fixed objects (poles) are detected by the vehicles if in line-of-sight (LoS)

and within a sensing range of 70 m. The same coverage area applies to A2A measurements. For the communications, we only consider the direct LoS path, as if the vehicles were equipped with LIDAR technology that could be blocked by obstacles such as buildings or other vehicles. The absolute driving speed adopted in the testing scenario ranges from 0 to about 60 km/h, with a mean and standard deviation speed of 0.2 km/h and 14 km/h, respectively. We point out that the motion models of the vehicles are not linear and that the factor graph to solve the distributed ICP method contains cycles. For the GNSS, A2A, and A2T observations, measurement errors are simulated as additive independent Gaussian noises with standard deviations of 2 m each.

For the training and testing of the ICP-MAPPO algorithm, we create two different simulations each composed of H=1500 timesteps. Model training is performed over $N_{\rm step}=2000$ episodes (or training steps), each characterized by a different realization of the measurements. For testing, 40 Monte Carlo (MC) evaluations are considered, unless otherwise specified. During training, we adopt a trajectory length $L_{\tau}=H/2$ to use at most 2 mini-batches, as suggested by [87], [106]. The entropy, reward and clipping coefficients have been chosen to be $\alpha=0.01$, $\beta=0.05$ and $\epsilon=0.2$, respectively. Note that $\beta=0.05$ would correspond to an improvement step of the reward function of 5 cm in a non-standardized state scenario. The discount factor is $\gamma=0.99$, while the Adam [107] learning rate is $\mu=10^{-5}$ with standard hyper-parameters.

Regarding the NN architectures, we adopt a critic network with three layers: a fully-connected (FC) linear layer with 256 neurons, a gated recurrent unit (GRU) with hidden size of 256 and a final FC linear layer. The actor is an MLP with two hidden linear layers of [128, 64] neurons and rectified linear unit (ReLU) activation functions, and an output layer with sigmoid activation function. Lastly, the belief network employs two bidirectional LSTM layers of 256 hidden neurons each and ReLU activation functions, followed by a Maxout unit with 128 output features and two linear layers of [64, 32] neurons.

B. Computational Complexity and Latency

To access the real-time processing capabilities of the proposed method in fulfilling the CAVs requirements on latency, we here investigate the computational complexities and communication delays of the proposed ICP-MAPPO solution with respect to the ICP algorithm. We specify that the number of floating point operations (FLOPs) for V_{ϕ} , π_{θ} and b_{ψ} are $0.82 \cdot 10^6$, $0.54 \cdot 10^6$, and $11.3 \cdot 10^6$, respectively. For comparison, the computational complexity of particle-based ICP methods is estimated with $O(N_{\rm mp} \cdot N \cdot K \cdot N_{\rm p})$, where $N_{\rm mp}$ and $N_{\rm p}$ are the number of message passing iterations and particles, respectively. The experiments are performed on a workstation machine with Intel(R) Xeon(R) Silver 4210R CPU @ 2.40 GHz, 96 GB RAM, and a Quadro RTX 6000 24 GB GPU, capable of achieving about $16.3 \cdot 10^{12}$ floating point operations per second (FLOPS) with just CPU performances. This implies a maximum latency for sample-inference of around $1\,\mu s$, which is expected to be truthful and accurate since the computational capabilities of CAVs are planned to far exceed our workstation capabilities with more than $4\cdot 10^{15}$ FLOPS for L5 SAE level [108].

When considering the communication delays with a hidden LSTM size of 256 bytes for ICP-MAPPO and about $N_{\mathrm{mp}} = 1000$ particles (each with 2 bytes for 2D position and 1 byte for the weight) in the ICP method, the data transmission would require approximately 1 and 10 packets, respectively. This estimate is based on 5G vehicle-to-vehicle (V2V) communications with a typical packet size of 300 bytes. Two communication scenarios are possible: direct V2V [109] or vehicle-to-network-to-vehicle (V2N2V) [110] when under cellular coverage. For direct V2V communication, the end-toend (E2E) packet latency is around 1 ms [109], resulting in 10 ms for ICP and 1 ms for ICP-MAPPO. In the V2N2V case, assuming the distances and scenarios described in [110], the E2E packet latency is around 4 ms, resulting in 40 ms for ICP and 4 ms for ICP-MAPPO. We note that the ICP E2E communication delay exceeds the 5 ms latency requirements of fully CAVs [111] in both scenarios, especially if a message passing procedure with multiple belief exchanges is considered. On the contrary, the ICP-MAPPO method meets the stringent latency requirements needed for fully CAVs.

C. Baseline Methods

As benchmark algorithms, we consider the following implementations:

- 1) KF-GNSS: Non-cooperative single-agent GNSS-based KF only using GNSS observations and perfect knowledge of the measurement standard deviation $\sigma^{(\text{GNSS})} = 2 \,\text{m}$. For the motion dynamics (1), we adopt a constant velocity model with standard deviation of the Gaussian-distributed velocity driving process calibrated on the data and equal to $0.5 \,\text{m/s}^2$.
- 2) ICP: Centralized ICP method from [32] with known A2A and A2T standard deviations, i.e., $\sigma^{({\rm A2A})} = \sigma^{({\rm A2T})} = 2\,{\rm m}$, and same motion model as for the KF-GNSS. Note that the use of the exact measurement statistics in generation and tracking allows to obtain the optimal performance (i.e. with no errors due to mis-modeling). Here the network of agents is fully-connected, i.e., all the agents share the same measurements.
- 3) Ego ICP-MAPPO: Proposed ICP-MAPPO method, with no-cooperation, i.e., only comprising the belief LSTM and imposing no connectivity with other agents, i.e., $\bar{a}_{i,j,t} = -1 \ \forall t \in \{0,\ldots,H-1\}, i \in \mathcal{V}, j \in \mathcal{N}_{i,t}$. In this way, each agent has to rely just on its measurements without performing aggregation of the neighbors' hidden features.

D. Results

1) Training performances: In the first assessment, we aim at verifying the convergence of the proposed ICP-MAPPO algorithm during the training episodes. In Fig. 5, 6 and 7, we report the mean belief LSTM loss, reward, and state value function, respectively, along with the 5-95 percentile as error bounds. The metrics are computed among agents and trajectory over the whole episode. From the figures, we

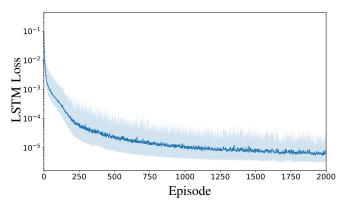


Fig. 5. Belief LSTM loss varying the number of training episodes.

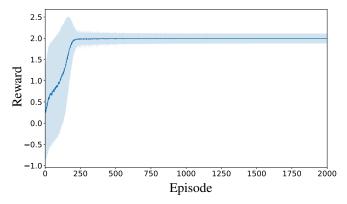


Fig. 6. Achieved reward varying the number of training episodes.

notice two distinct phases of the training: before and after reward convergence. In the first phase, i.e., before episode 250, the exploration is encouraged, leading to a much higher variability of the reward and a very rapid decrease of the LSTM loss function. After passing into the second phase, the positioning improvement becomes smaller, with a consequent convergence of the reward to the value of 2. Notably, also the mean value function converges after about 250 episodes, but with a high variance between agents and trajectories. This may be indicative of a rich and complex environment where the optimal policy may not be static, but rather dynamic and contingent on the interactions between agents and the environment. Indeed, the complexity of the state, e.g., each agent has a different trajectory in the space, can lead to a wide range of value function estimates as different states are visited with varying frequencies.

2) Cooperative positioning testing: This experiment has the objective of comparing the positioning capabilities of ICP-MAPPO with respect to the baselines in an unseen testing trajectory. To this aim, Fig. 8 shows the root mean square error (RMSE) on the vehicle position at each timestep of the trajectory (Fig. 8a) and the corresponding cumulative density function (CDF) of the absolute error (Fig. 8b). The RMSE is computed among the agents at the single timestep, while the mean and error bounds are computed within the MC evaluations. From the results, we observe that the Ego

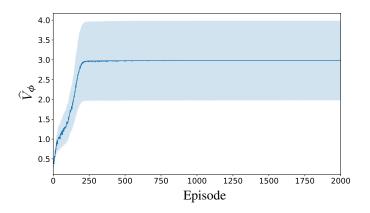


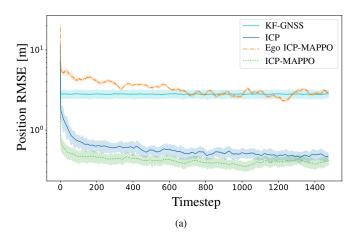
Fig. 7. Mean value function varying the number of training episodes.

ICP-MAPPO, which only relies on GNSS measurements, converges to the KF-GNSS method, indicating a correct usage of the observations to estimate the position. Passing to the cooperative methods, we notice a higher speed of convergence of ICP-MAPPO with respect to the conventional ICP. This is mainly due to the learned vehicles' dynamics and to the effective combination of neighbors' observations. As a consequence, the ICP-MAPPO algorithm outperforms the ICP method in terms of absolute error by 21%, passing from a median of 42 cm to 33 cm.

3) Generalization capabilities: This experiment aims at assessing the generalization capabilities of the proposed method in unseen scenarios. To evaluate the environmental dependence of our model, we tested the pre-trained ICP-MAPPO on a different CARLA map, specifically *Town10*. In Fig. 9, we plotted the position RMSE on testing trajectories in both *Town02* (used for training) and *Town10* (unseen environment), varying the number of passive objects in the respective map. We shall notice that the numbers of poles in *Town10* and *Town02* are 146 and 72, respectively. Since ICP-MAPPO was trained with a maximum input size of 72 measurements, we adjusted the number of targets up to 72 for this experiment.

The results in Fig. 9 confirm that, even in the unseen scenario, a higher number of vehicles increases the positioning accuracy thanks to the cooperation among vehicles. Comparing the results on *Town02* and *Town10*, we note that in the limit-case of no measurements shared among agents, the performances in the two scenarios coincide. On the contrary, when the number of features increases, the performances on the unseen scenario are slightly lower (i.e., about 10 cm) despite the completely new environment.

4) Communication efficiency: In this last assessment, we test the effectiveness of the policy choices in terms of cooperation power and communication efficiency. In Fig. 10 we report the position RMSE at convergence (Fig. 10a) and the mean number of selected agents from the policy (Fig. 10b) varying the maximum degree of connectivity allowed in the network. In Fig. 10a we observe an intuitive inverse relation between the maximum cooperative agents and the RMSE, with a rapid decrease under 1 m of RMSE with just 2 agents. Notably, after 8 cooperative agents, the improvement in RMSE is negligible, with convergence to about 40 cm. To study this behaviour, in



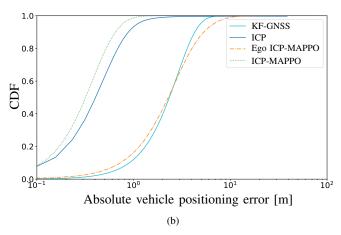


Fig. 8. Testing performances on the cooperative scenario. (a) RMSE of the position over time for the single-agent KF-GNSS, ICP, proposed single agent and cooperative ICP-MAPPO. (b) CDF of the absolute error.

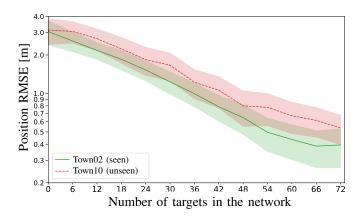
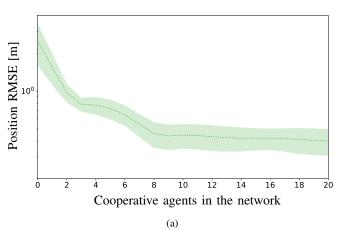


Fig. 9. RMSE on the position achieved by ICP-MAPPO varying the number of targets (i.e., poles) in two distinct environments.

Fig. 10b we notice that the policy tends to select no more than 9 agents for cooperation. This likely occurs because the marginal benefits of additional cooperation diminish beyond this point, leading agents to prefer collaboration with only their closest neighbors. Indeed, incorporating data from distant agents that do not observe common targets results in only slight enhancements in positional accuracy. Lastly, we highlight that the ICP-MAPPO has higher performance than



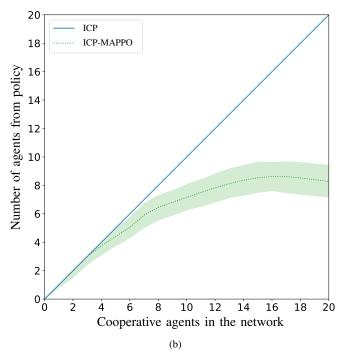


Fig. 10. Communication efficiency comparison between the ICP and the proposed ICP-MAPPO methods. (a) RMSE on the position varying the maximum number of cooperative agents in the network. (b) Mean number of neighbor agents selected by the policy varying the maximum connectivity of the graph.

the ICP method for the same number of cooperative agents in the network.

To evaluate the trade-off between positioning accuracy and communication overhead, in Fig. 11, we plot the mean number of A2A links, considering varying numbers of cooperative vehicles in {2, 6, 10, 15, 20}. We observe that with a smaller number of cooperative agents, such as 2, the ICP-MAPPO tends to employ all available agents, leveraging neighbors' measurements to rapidly reduce GNSS uncertainty. Conversely, with a higher number of agents, particularly beyond 10, the benefits of additional cooperation decrease (as shown in Fig. 10a). This is because only the closest neighbors with a significant number of shared targets substantially enhance positioning accuracy. Notably, with 10 and 20 agents, ICP-

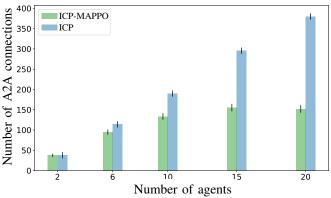


Fig. 11. Mean number of A2A connections in the network graph, for the ICP and the proposed ICP-MAPPO algorithms, and different maximum number of cooperative agents.

MAPPO reduces the number of links by 30% and 60%, respectively, compared to ICP.

VI. CONCLUSION

In this paper, we addressed the problem of CP in a distributed network of agents that exploit passive detected targets to improve the positioning accuracy according to the ICP framework. We provided a generalization of the Bayesian ICP solution by exploiting the MARL approach, which enables the dynamic optimization of the A2A links used for cooperation accounting for partial observability of the state. We presented a novel ICP-MAPPO algorithm where the agents actively select the neighbors to communicate with by following their optimized policy. This allows to minimize the communication overhead for cooperation, while improving the positioning accuracy of ego-agent systems. The proposed solution is proven to outperform single and multi-agent conventional approaches thanks to DL-based states' belief and policy models.

Realistic simulations of a C-ITS scenario created with CARLA simulator demonstrate the superior performances of ICP-MAPPO with state-of-the-art ICP methods, both in terms of positioning accuracy and efficiency of communications. The cooperation is indeed intelligently exploited to enhance the performances and, at the same time, the communication efficiency, by selecting ad-hoc neighbors that are relevant for the task. The benefits of the approach look promising for applications where groups of agents have a common inference objective and predictions/decisions need to be taken based on incomplete or uncertain data.

As future work, we envision the extension of the proposed method to decentralized frameworks [112], incorporating also data association of the targets to the measurements. Additionally, performances could be enhanced by exploiting a higher dimension of latent features within object detectors, instead of filtering specific objects such as poles. This approach would allow vehicles to exchange much more meaningful information in a compressed manner. Furthermore, including motion planning [113] could enable the system to not only estimate but also modify the vehicles' states according to their destinations. Finally, introducing safe RL [114] by adding

safety constraints related to communication resources, such as maximum available bandwidth, would ensure that the policies learned by the agents remain efficient under real-world communication constraints.

APPENDIX A PROOF OF (18)

To prove (18), we start by writing the gradient of the RL objective function in (17) as

$$\nabla_{\boldsymbol{\theta}} J(\pi_{\boldsymbol{\theta}}) = \nabla_{\boldsymbol{\theta}} \mathbb{E}_{\boldsymbol{\tau} \sim p(\boldsymbol{\tau}|\pi_{\boldsymbol{\theta}})} \left\{ \widetilde{R}(\boldsymbol{\tau}) \right\} = \nabla_{\boldsymbol{\theta}} \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau}|\pi_{\boldsymbol{\theta}}) \, \widetilde{R}(\boldsymbol{\tau})$$
$$= \sum_{\boldsymbol{\tau}} \nabla_{\boldsymbol{\theta}} p(\boldsymbol{\tau}|\pi_{\boldsymbol{\theta}}) \, \widetilde{R}(\boldsymbol{\tau}). \tag{A1}$$

Now, we can rewrite the gradient of the trajectory PDF $\nabla_{\theta} p(\tau | \pi_{\theta})$ using the log-derivative trick as

$$\nabla_{\boldsymbol{\theta}} p(\boldsymbol{\tau}|\boldsymbol{\pi}_{\boldsymbol{\theta}}) = p(\boldsymbol{\tau}|\boldsymbol{\pi}_{\boldsymbol{\theta}}) \nabla_{\boldsymbol{\theta}} \log (p(\boldsymbol{\tau}|\boldsymbol{\pi}_{\boldsymbol{\theta}})). \tag{A2}$$

Given that the gradient of the log-trajectory PDF $\nabla_{\theta} \log (p(\tau | \pi_{\theta}))$ is

$$\nabla_{\boldsymbol{\theta}} \log \left(p(\boldsymbol{\tau} | \boldsymbol{\pi}_{\boldsymbol{\theta}}) \right) = \nabla_{\boldsymbol{\theta}} \log \left(T_0 \prod_{t=0}^{H-1} T(\boldsymbol{s}_{t+1} | \boldsymbol{s}_t, \boldsymbol{a}_t) \, \pi_{\boldsymbol{\theta}}(\boldsymbol{a}_t | \boldsymbol{s}_t) \right)$$
$$= \sum_{t=0}^{H-1} \nabla_{\boldsymbol{\theta}} \log \left(\pi_{\boldsymbol{\theta}}(\boldsymbol{a}_t | \boldsymbol{s}_t) \right) \tag{A3}$$

we can rewrite (A1) as

$$\nabla_{\boldsymbol{\theta}} J(\pi_{\boldsymbol{\theta}}) = \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau}|\pi_{\boldsymbol{\theta}}) \nabla_{\boldsymbol{\theta}} \log \left(p(\boldsymbol{\tau}|\pi_{\boldsymbol{\theta}}) \right) \widetilde{R}(\boldsymbol{\tau})$$

$$= \mathbb{E}_{\boldsymbol{\tau} \sim p(\boldsymbol{\tau}|\pi_{\boldsymbol{\theta}})} \left\{ \nabla_{\boldsymbol{\theta}} \log \left(p(\boldsymbol{\tau}|\pi_{\boldsymbol{\theta}}) \right) \widetilde{R}(\boldsymbol{\tau}) \right\}$$

$$= \mathbb{E}_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim p(\mathbf{s}_{t}, \mathbf{a}_{t}|\pi_{\boldsymbol{\theta}})} \left\{ \sum_{t=0}^{H-1} \nabla_{\boldsymbol{\theta}} \log \left(\pi_{\boldsymbol{\theta}}(\boldsymbol{a}_{t}|\mathbf{s}_{t}) \right) \times \sum_{t=0}^{H-1} \gamma^{t} R(\mathbf{s}_{t}, \mathbf{a}_{t}) \right\}.$$

Since the action \mathbf{a}_t at time t only influences the future rewards and not the past ones, (A4) can be equivalently rewritten as

$$\nabla_{\boldsymbol{\theta}} J(\pi_{\boldsymbol{\theta}}) = \mathbb{E}_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim p(\mathbf{s}_{t}, \mathbf{a}_{t} | \pi_{\boldsymbol{\theta}})} \left\{ \sum_{t=0}^{H-1} \nabla_{\boldsymbol{\theta}} \log \left(\pi_{\boldsymbol{\theta}}(\boldsymbol{a}_{t} | \boldsymbol{s}_{t}) \right) \mathsf{R}_{t} \right\}$$
(A5)

where we used the reward-to-go at time $R_t = \sum_{t'=t}^{H-1} \gamma^{t'-t} R(\mathbf{s}_{t'}, \mathbf{a}_{t'})$, as opposed to R_0 .

Since it can be proven that for any function of the state $B(\mathbf{s}_t)$ called baseline, we have that $\mathbb{E}_{\mathbf{a}_t \sim \pi_{\theta}(\mathbf{a}_t|\mathbf{s}_t)} \Big\{ \nabla_{\theta} \log \big(\pi_{\theta}(\mathbf{a}_t|\mathbf{s}_t) \, B(\mathbf{s}_t) \big) \Big\} = 0$, then we can reduce the variance of the PO algorithm, while remaining unbiased, by subtracting the baseline from the reward-to-go as

$$\nabla_{\boldsymbol{\theta}} J(\pi_{\boldsymbol{\theta}}) = \mathbb{E}_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim p(\mathbf{s}_{t}, \mathbf{a}_{t} | \pi_{\boldsymbol{\theta}})} \left\{ \sum_{t=0}^{H-1} \left[\nabla_{\boldsymbol{\theta}} \log \left(\pi_{\boldsymbol{\theta}}(\boldsymbol{a}_{t} | \boldsymbol{s}_{t}) \right) \times \left(\mathsf{R}_{t} - B(\mathbf{s}_{t}) \right) \right] \right\}.$$
(A6)

Finally, R_t and $B(\mathbf{s}_t)$ are usually substituted with their estimates $Q^\pi(\mathbf{s}_t, \mathbf{a}_t)$ and $V^\pi(\mathbf{s}_t)$, respectively, leading to the definition of the advantage function $A_t = Q^\pi(\mathbf{s}_t, \mathbf{a}_t) - V^\pi(\mathbf{s}_t)$. Recently, more advanced versions of the advantage function, as the generalized advantage estimator (GAE) function A_t^{GAE} have been proposed in the literature [100] to regulate the bias-variance trade-off, increase stability, efficiency, and obtain faster convergence. We want to point out that usage of the baseline and/or the estimate of R_t are not necessary, and thus any function $F_t \in \left\{\mathsf{R}_t, Q^\pi(\mathbf{s}_t, \mathbf{a}_t), \mathsf{R}_t - V^\pi(\mathbf{s}_t), A_t, A_t^{\mathrm{GAE}}\right\}$ is a valid choice.

REFERENCES

- [1] M. Z. Win *et al.*, "Network localization and navigation via cooperation," *IEEE Commun. Mag.*, vol. 49, no. 5, pp. 56–62, May 2011.
- [2] L. Gao et al., "Cooperative localization in transportation 5.0," IEEE Trans. Intell. Veh., pp. 1–6, Mar. 2024.
- [3] M. Z. Win, Y. Shen, and W. Dai, "A theoretical foundation of network localization and navigation," *Proc. IEEE*, vol. 106, no. 7, pp. 1136– 1165, Jul. 2018.
- [4] Y. Gao, H. Jing, M. Dianati, C. M. Hancock, and X. Meng, "Performance analysis of robust cooperative positioning based on GPS/UWB integration for connected autonomous vehicles," *IEEE Trans. Intell. Veh.*, vol. 8, no. 1, pp. 790–802, Jan. 2023.
- [5] A. Mahmoud, A. Noureldin, and H. S. Hassanein, "Integrated positioning for connected vehicles," *IEEE Trans. Intell. Transp. Syst.*, vol. 21, no. 1, pp. 397–409, Jan. 2020.
- [6] A. Conti, S. Mazuelas, S. Bartoletti, W. C. Lindsey, and M. Z. Win, "Soft information for localization-of-things," *Proc. IEEE*, vol. 107, no. 11, pp. 2240–2264, Nov. 2019.
- [7] F. Gustafsson and F. Gunnarsson, "Mobile positioning using wireless networks: possibilities and fundamental limitations based on available wireless network measurements," *IEEE Signal Process. Mag.*, vol. 22, no. 4, pp. 41–53, Jul. 2005.
- [8] C. A. Gómez-Vega, Z. Liu, C. A. Gutiérrez, M. Z. Win, and A. Conti, "Efficient deployment strategies for network localization with assisting nodes," *IEEE Trans. Mobile Comput.*, vol. 23, no. 5, pp. 6272–6287, May 2024.
- [9] L. Chen et al., "Milestones in autonomous driving and intelligent vehicles: Survey of surveys," *IEEE Trans. Intell. Veh.*, vol. 8, no. 2, pp. 1046–1056, Feb. 2023.
- [10] D. Cao et al., "Future directions of intelligent vehicles: Potentials, possibilities, and perspectives," *IEEE Trans. Intell. Veh.*, vol. 7, no. 1, pp. 7–10, Mar. 2022.
- [11] P. Yang, D. Duan, C. Chen, X. Cheng, and L. Yang, "Multi-sensor multi-vehicle (MSMV) localization and mobility tracking for autonomous driving," *IEEE Trans. Veh. Technol.*, vol. 69, no. 12, pp. 14 355–14 364, Oct. 2020.
- [12] J. Ji, A. Khajepour, W. W. Melek, and Y. Huang, "Path planning and tracking for vehicle collision avoidance based on model predictive control with multiconstraints," *IEEE Trans. Veh. Technol.*, vol. 66, no. 2, pp. 952–964, Feb. 2017.
- [13] T. G. Reid et al., "Localization requirements for autonomous vehicles," SAE Int. J. Connected Automated Veh., vol. 2, no. 3, pp. 12–02–03– 0012, Oct. 2019.
- [14] S. Kuutti, S. Fallah, K. Katsaros, M. Dianati, F. Mccullough, and A. Mouzakitis, "A survey of the state-of-the-art localization techniques and their potentials for autonomous vehicle applications," *IEEE Internet Things J.*, vol. 5, no. 2, pp. 829–846, Mar. 2018.
- [15] M. H. C. Garcia et al., "A tutorial on 5G NR V2X communications," IEEE Commun. Surveys & Tuts., vol. 23, no. 3, pp. 1972–2026, Feb. 2021.
- [16] F. Morselli, S. Modarres Razavi, M. Z. Win, and A. Conti, "Soft information-based localization for 5G networks and beyond," *IEEE Trans. Wireless Commun.*, vol. 22, no. 12, pp. 9923–9938, Dec. 2023.
- [17] H. Zhou, W. Xu, J. Chen, and W. Wang, "Evolutionary V2X technologies toward the internet of vehicles: Challenges and opportunities," *Proc. IEEE*, vol. 108, no. 2, pp. 308–323, Feb. 2020.
- [18] G. Torsoli, M. Z. Win, and A. Conti, "Blockage intelligence in complex environments for beyond 5G localization," *IEEE J. Sel. Areas Commun.*, vol. 41, no. 6, pp. 1688–1701, Jun. 2023.

- [19] A. Conti, G. Torsoli, C. A. Gómez-Vega, A. Vaccari, G. Mazzini, and M. Z. Win, "3GPP-compliant datasets for xG location-aware networks," *IEEE Open J. Veh. Technol.*, vol. 5, pp. 473–484, Dec. 2024.
- [20] A. Conti, G. Torsoli, C. A. Gómez-Vega, A. Vaccari, and M. Z. Win, "xG-Loc: 3GPP-compliant datasets for xG location-aware networks," *IEEE Dataport*, Dec. 2023.
- [21] L. Italiano, B. Camajori Tedeschini, M. Brambilla, H. Huang, M. Nicoli, and H. Wymeersch, "A tutorial on 5G positioning," *IEEE Commun. Surveys & Tuts.*, pp. 1–48, Aug. 2024.
- [22] B. Camajori Tedeschini, G. Kwon, M. Nicoli, and M. Z. Win, "Real-time Bayesian neural networks for 6G cooperative positioning and tracking," *IEEE J. Sel. Areas Commun.*, vol. 42, no. 9, pp. 2322–2338, Aug. 2024.
- [23] A. Conti et al., "Location awareness in beyond 5G networks," IEEE Commun. Mag., vol. 59, no. 11, pp. 22–27, Nov. 2021.
- [24] B. Camajori Tedeschini and M. Nicoli, "Cooperative deep-learning positioning in mmWave 5G-advanced networks," *IEEE J. Sel. Areas Commun.*, vol. 41, no. 12, pp. 3799–3815, Dec. 2023.
- [25] M. A. Javed, S. Zeadally, and E. B. Hamida, "Data analytics for cooperative intelligent transport systems," *Veh. Commun.*, vol. 15, pp. 63–72, Jan. 2019.
- [26] A. Alalewi, I. Dayoub, and S. Cherkaoui, "On 5G-V2X use cases and enabling technologies: A comprehensive survey," *IEEE Access*, vol. 9, pp. 107710–107737, Jul. 2021.
- [27] K. Sehla, T. M. T. Nguyen, G. Pujolle, and P. B. Velloso, "Resource allocation modes in C-V2X: From LTE-V2X to 5G-V2X," *IEEE Internet Things J.*, vol. 9, no. 11, pp. 8291–8314, Mar. 2022.
- [28] M. Driusso, C. Marshall, M. Sabathy, F. Knutti, H. Mathis, and F. Babich, "Vehicular position tracking using LTE signals," *IEEE Trans. Veh. Technol.*, vol. 66, no. 4, pp. 3376–3391, Apr. 2017.
- [29] G. Kwon, Z. Liu, A. Conti, H. Park, and M. Z. Win, "Integrated localization and communication for efficient millimeter wave networks," *IEEE J. Sel. Areas Commun.*, vol. 41, no. 12, pp. 3925–3941, Dec. 2023.
- [30] F. Liu et al., "Integrated sensing and communications: Toward dualfunctional wireless networks for 6G and beyond," *IEEE J. Sel. Areas Commun.*, vol. 40, no. 6, pp. 1728–1767, Jun. 2022.
- [31] G. Kwon, A. Conti, H. Park, and M. Z. Win, "Joint communication and localization in millimeter wave networks," *IEEE J. Sel. Topics Signal Process.*, vol. 15, no. 6, pp. 1439–1454, Nov. 2021.
- [32] G. Soatti, M. Nicoli, N. Garcia, B. Denis, R. Raulefs, and H. Wymeer-sch, "Implicit cooperative positioning in vehicular networks," *IEEE Trans. Intell. Transp. Syst.*, vol. 19, no. 12, pp. 3964–3980, Dec. 2018.
- [33] M. Brambilla, M. Nicoli, G. Soatti, and F. Deflorio, "Augmenting vehicle localization by cooperative sensing of the driving environment: Insight on data association in urban traffic scenarios," *IEEE Trans. Intell. Transp. Syst.*, vol. 21, no. 4, pp. 1646–1663, Apr. 2020.
- [34] M. Brambilla et al., "Cooperative localization and multitarget tracking in agent networks with the sum-product algorithm," *IEEE Open J. Signal Process.*, vol. 3, pp. 169–195, Mar. 2022.
- [35] F. Jin, K. Liu, C. Liu, T. Cheng, H. Zhang, and V. C. S. Lee, "A cooperative vehicle localization and trajectory prediction framework based on belief propagation and transformer model," *IEEE Trans. Consum. Electron.*, pp. 1–1, Feb. 2024.
- [36] M. Z. Win, W. Dai, Y. Shen, G. Chrisikos, and H. Vincent Poor, "Network operation strategies for efficient localization and navigation," *Proc. IEEE*, vol. 106, no. 7, pp. 1224–1254, Jul. 2018.
- [37] A. Conti, M. Guerra, D. Dardari, N. Decarli, and M. Z. Win, "Network experimentation for cooperative localization," *IEEE J. Sel. Areas Commun.*, vol. 30, no. 2, pp. 467–475, Feb. 2012.
- [38] W. Dai, Y. Shen, and M. Z. Win, "Distributed power allocation for cooperative wireless network localization," *IEEE J. Sel. Areas Commun.*, vol. 33, no. 1, pp. 28–40, Jan. 2015.
- [39] P. Yang, C. Xiang, and S. Zhang, "Distributed joint power and bandwidth allocation for multiagent cooperative localization," *IEEE Commun. Lett.*, vol. 26, no. 11, pp. 2601–2605, Aug. 2022.
- [40] T. Zhang, A. F. Molisch, Y. Shen, Q. Zhang, H. Feng, and M. Z. Win, "Joint power and bandwidth allocation in wireless cooperative localization networks," *IEEE Trans. Wireless Commun.*, vol. 15, no. 10, pp. 6527–6540, Oct. 2016.
- [41] J. Nie, J. Yan, H. Yin, L. Ren, and Q. Meng, "A multimodality fusion deep neural network and safety test strategy for intelligent vehicles," *IEEE Trans. Intell. Veh.*, vol. 6, no. 2, pp. 310–322, Sep. 2020.
- [42] Z. Liu et al., "Robust target recognition and tracking of self-driving cars with radar and camera information fusion under severe weather conditions," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 7, pp. 6640– 6653, Feb. 2021.

- [43] T. Wang, Y. Shen, A. Conti, and M. Z. Win, "Network navigation with scheduling: Error evolution," *IEEE Trans. Inf. Theory*, vol. 63, no. 11, pp. 7509–7534, Nov. 2017.
- [44] B. Teague, Z. Liu, F. Meyer, A. Conti, and M. Z. Win, "Network localization and navigation with scalable inference and efficient operation," *IEEE Trans. Mobile Comput.*, vol. 21, no. 6, pp. 2072–2087, Jun. 2022.
- [45] S. Dwivedi, D. Zachariah, A. De Angelis, and P. Handel, "Cooperative decentralized localization using scheduled wireless transmissions," *IEEE Commun. Lett.*, vol. 17, no. 6, pp. 1240–1243, Jun. 2013.
- [46] L. Barbieri, B. Camajori Tedeschini, M. Brambilla, and M. Nicoli, "Deep learning-based cooperative LiDAR sensing for improved vehicle positioning," *IEEE Trans. Signal Process.*, vol. 72, pp. 1666–1682, Mar. 2024.
- [47] Y. Weiss and W. T. Freeman, "Correctness of belief propagation in Gaussian graphical models of arbitrary topology," *Neural Comput.*, vol. 13, no. 10, pp. 2173–2200, Oct. 2001.
- [48] J. Yedidia, W. Freeman, and Y. Weiss, "Constructing free-energy approximations and generalized belief propagation algorithms," *IEEE Trans. Inf. Theory*, vol. 51, no. 7, pp. 2282–2312, Jul. 2005.
- [49] M. J. Wainwright and M. I. Jordan, "Graphical models, exponential families, and variational inference," Found. Trends® Mach. Learn., vol. 1, no. 1-2, pp. 1–305, Nov. 2008.
- [50] F. Meyer *et al.*, "Message passing algorithms for scalable multitarget tracking," *Proc. IEEE*, vol. 106, no. 2, pp. 221–259, Feb. 2018.
- [51] G. Soldi, F. Meyer, P. Braca, and F. Hlawatsch, "Self-tuning algorithms for multisensor-multitarget tracking using belief propagation," *IEEE Trans. Signal Process.*, vol. 67, no. 15, pp. 3922–3937, Aug. 2019.
- [52] M. G. Kibria, K. Nguyen, G. P. Villardi, O. Zhao, K. Ishizu, and F. Kojima, "Big data analytics, machine learning, and artificial intelligence in next-generation wireless networks," *IEEE Access*, vol. 6, pp. 32328– 32338, May 2018.
- [53] B. Camajori Tedeschini, M. Brambilla, L. Barbieri, G. Balducci, and M. Nicoli, "Cooperative lidar sensing for pedestrian detection: Data association based on message passing neural networks," *IEEE Trans. Signal Process.*, vol. 71, pp. 3028–3042, Aug. 2023.
- [54] B. Camajori Tedeschini, M. Brambilla, and M. Nicoli, "Message passing neural network versus message passing algorithm for cooperative positioning," *IEEE Trans. Cogn. Commun. Netw.*, vol. 9, no. 6, pp. 1666–1676, Dec. 2023.
- [55] B. Camajori Tedeschini, M. Brambilla, and M. Nicoli, "Split consensus federated learning: an approach for distributed training and inference," *IEEE Access*, vol. 12, pp. 119 535–119 549, Aug. 2024.
- [56] R. S. Sutton and A. G. Barto, Reinforcement Learning: An Introduction. A Bradford Book, Oct. 2018.
- [57] D. P. Bertsekas, A Course in Reinforcement Learning, 1st ed. Athena Scientific, Jun. 2023.
- [58] D. P. Bertsekas, Reinforcement Learning and Optimal Control, 1st ed. Athena Scientific, Jul. 2019.
- [59] V. Mnih et al., "Human-level control through deep reinforcement learning," Nature, vol. 518, no. 7540, pp. 529–533, Feb. 2015.
- [60] D. P. Bertsekas, Lessons from AlphaZero for Optimal, Model Predictive, and Adaptive Control, 1st ed. Athena Scientific, Aug. 2021.
- [61] T. P. Lillicrap et al., "Continuous control with deep reinforcement learning," in Proc. 4th Int. Conf. Learn. Representations, Sep. 2016, pp. 1–14.
- [62] N. R. Ke et al., "Learning dynamics model in reinforcement learning by incorporating the long term future," in *Proc. 7th Int. Conf. Learn. Representations*, Mar. 2019, pp. 1–14.
- [63] D. P. Bertsekas, Dynamic Programming and Optimal Control, 4th ed. Athena Scientific, Oct. 2000.
- [64] D. P. Bertsekas, Dynamic Programming and Optimal Control, Volume II: Approximate Dynamic Programming, 4th ed. Athena Scientific, Jun. 2012.
- [65] N. C. Luong et al., "Applications of deep reinforcement learning in communications and networking: A survey," *IEEE Commun. Surveys* & Tuts., vol. 21, no. 4, pp. 3133–3174, Oct. 2019.
- [66] H. Zhang, N. Yang, W. Huangfu, K. Long, and V. C. M. Leung, "Power control based on deep reinforcement learning for spectrum sharing," *IEEE Trans. Wireless Commun.*, vol. 19, no. 6, pp. 4209–4219, Mar. 2020
- [67] F. Meng, P. Chen, L. Wu, and J. Cheng, "Power allocation in multiuser cellular networks: Deep reinforcement learning approaches," *IEEE Trans. Wireless Commun.*, vol. 19, no. 10, pp. 6255–6267, Jun. 2020.
- [68] X. Li, J. Fang, W. Cheng, H. Duan, Z. Chen, and H. Li, "Intelligent power control for spectrum sharing in cognitive radios: A deep reinforcement learning approach," *IEEE Access*, vol. 6, pp. 25463–25473, Apr. 2018.

- [69] J. Vlachogiannis and N. Hatziargyriou, "Reinforcement learning for reactive power control," *IEEE Trans. Power Syst.*, vol. 19, no. 3, pp. 1317–1325, Aug. 2004.
- [70] L. Liang, H. Ye, and G. Y. Li, "Spectrum sharing in vehicular networks based on multi-agent reinforcement learning," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 10, pp. 2282–2292, Aug. 2019.
- [71] O. Naparstek and K. Cohen, "Deep multi-user reinforcement learning for distributed dynamic spectrum access," *IEEE Trans. Wireless Commun.*, vol. 18, no. 1, pp. 310–323, Jan. 2019.
- [72] H. Song, L. Liu, J. Ashdown, and Y. Yi, "A deep reinforcement learning framework for spectrum management in dynamic spectrum access," *IEEE Internet Things J.*, vol. 8, no. 14, pp. 11 208–11 218, Jan. 2021.
- [73] Q. Luo, C. Li, T. H. Luan, and W. Shi, "Collaborative data scheduling for vehicular edge computing via deep reinforcement learning," *IEEE Internet Things J.*, vol. 7, no. 10, pp. 9637–9650, Mar. 2020.
- [74] T. T. Nguyen, N. D. Nguyen, and S. Nahavandi, "Deep reinforcement learning for multiagent systems: A review of challenges, solutions, and applications," *IEEE Trans. Cybern.*, vol. 50, no. 9, pp. 3826–3839, Mar. 2020.
- [75] F. A. Oliehoek and C. Amato, A Concise Introduction to Decentralized POMDPs, ser. SpringerBriefs in Intelligent Systems. Springer International Publishing, Jun. 2016.
- [76] L. Kraemer and B. Banerjee, "Multi-agent reinforcement learning as a rehearsal for decentralized planning," *Neurocomputing*, vol. 190, no. C, p. 82–94, May 2016.
- [77] S. Bhattacharya, S. Kailas, S. Badyal, S. Gil, and D. Bertsekas, "Multiagent reinforcement learning: Rollout and policy iteration for POMDP with application to multirobot problems," *IEEE Trans. Robot.*, vol. 40, pp. 2003–2023, Dec. 2024.
- [78] S. Omidshafiei, J. Pazis, C. Amato, J. P. How, and J. Vian, "Deep decentralized multi-task multi-agent reinforcement learning under partial observability," in *Proc. 34th Int. Conf. Mach. Learn.*, Aug. 2017, p. 2681–2690
- [79] S. Gronauer and K. Diepold, "Multi-agent deep reinforcement learning: a survey," *Artif. Intell. Rev.*, vol. 55, no. 2, pp. 895–943, Apr. 2021.
- [80] K. Zhang, Z. Yang, and T. Başar, Multi-Agent Reinforcement Learning: A Selective Overview of Theories and Algorithms. Springer International Publishing, Jun. 2021, pp. 321–384.
- [81] L. Busoniu, R. Babuska, and B. De Schutter, "A comprehensive survey of multiagent reinforcement learning," *IEEE Trans. Syst., Man, Cybern., Part C (Appl. Rev.)*, vol. 38, no. 2, p. 156–172, Mar. 2008.
- [82] T. Rashid, M. Samvelyan, C. S. de Witt, G. Farquhar, J. Foerster, and S. Whiteson, "Monotonic value function factorisation for deep multi-agent reinforcement learning," J. Mach. Learn. Res. 21(178):1-51, 2020, Mar. 2020.
- [83] K. Son, D. Kim, W. J. Kang, D. E. Hostallero, and Y. Yi, "QTRAN: Learning to factorize with transformation for cooperative multi-agent reinforcement learning," in *Proc. 36th Int. Conf. Mach. Learn.*, May 2019, pp. 5887–5896.
- [84] J. Wang, Z. Ren, T. Liu, Y. Yu, and C. Zhang, "QPLEX: Duplex dueling multi-agent Q-learning," in *Proc. 38th Int. Conf. Mach. Learn.*, Aug. 2021, pp. 1–27.
- [85] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," in *Proc.* 31th Int. Conf. Neural Inf. Process. Syst., Dec. 2017, pp. 6382–6393.
- [86] J. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson, "Counterfactual multi-agent policy gradients," in *Proc. 32nd AAAI Conf. Artif. Intell.*, Feb. 2018, pp. 2974–2982.
- [87] C. Yu et al., "The surprising effectiveness of PPO in cooperative, multiagent games," in Proc. 36th Int. Conf. Neural Inf. Process. Syst., Mar. 2021, pp. 1–14.
- [88] V. Mnih et al., "Asynchronous methods for deep reinforcement learning," in Proc. 33rd Int. Conf. Mach. Learn., Feb. 2016, pp. 1928–1937.
- [89] J. K. Gupta, M. Egorov, and M. Kochenderfer, "Cooperative multiagent control using deep reinforcement learning," in *Auton. Agents Multiagent Syst.*, ser. Lecture Notes in Computer Science, G. Sukthankar and J. A. Rodriguez-Aguilar, Eds. Springer International Publishing, Nov. 2017, pp. 66–83.
- [90] F. Christianos, G. Papoudakis, M. A. Rahman, and S. V. Albrecht, "Scaling multi-agent reinforcement learning with selective parameter sharing," in *Proc. 38th Int. Conf. Mach. Learn.*, Feb. 2021, pp. 1989– 1998.
- [91] J. K. Terry, N. Grammel, S. Son, and B. Black, "Parameter sharing for heterogeneous agents in multi-agent reinforcement learning," *ArXiv*, pp. 1–16, May 2020.

- [92] Z. Xia et al., "Multi-agent reinforcement learning aided intelligent UAV swarm for target tracking," *IEEE Trans. Veh. Technol.*, vol. 71, no. 1, pp. 931–945. Nov. 2021.
- [93] B. Peng, G. Seco-Granados, E. Steinmetz, M. Frohle, and H. W. Wymeersch, "Decentralized scheduling for cooperative localization with deep reinforcement learning," *IEEE Trans. Veh. Technol.*, vol. 68, no. 5, pp. 4295–4305, May 2019.
- [94] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "CARLA: An open urban driving simulator," in *Proc. 1st Conf. Robot Learning*, Nov. 2017, pp. 1–16.
- [95] M. Arulampalam, S. Maskell, N. Gordon, and T. Clapp, "A tutorial on particle filters for online nonlinear/non-Gaussian Bayesian tracking," *IEEE Trans. Signal Process.*, vol. 50, no. 2, pp. 174–188, Feb. 2002.
- [96] A. Tampuu et al., "Multiagent cooperation and competition with deep reinforcement learning," PLOS ONE, vol. 12, Apr. 2017.
- [97] N. Mehta, P. Tadepalli, and A. Fern, "Multi-agent shared hierarchy reinforcement learning," in *ICML Workshop Rich Representations* Reinforcement Learn., Aug. 2005, pp. 45–50.
- [98] P. Sunehag et al., "Value-decomposition networks for cooperative multi-agent learning," ArXiv, pp. 1–17, Jun. 2017.
- [99] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Mach. Learn.*, vol. 8, no. 3-4, pp. 229–256, Jan. 2005.
- [100] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel, "High-dimensional continuous control using generalized advantage estimation," in *Proc. 4th Int. Conf. Mach. Learn.*, Jun. 2016, pp. 1–14.
- [101] J. Schulman, S. Levine, P. Moritz, M. I. Jordan, and P. Abbeel, "Trust region policy optimization," in *Proc. 31st Int. Conf. Mach. Learn.*, Feb. 2015, pp. 1889–1897.
- [102] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," ArXiv, pp. 1–12, Jul. 2017.
- [103] V. Feinberg, A. Wan, I. Stoica, M. I. Jordan, J. E. Gonzalez, and S. Levine, "Model-based value estimation for efficient model-free reinforcement learning," in *Proc. 35st Int. Conf. Mach. Learn.*, Feb. 2018, pp. 1–12.
- [104] J. Buckman, D. Hafner, G. Tucker, E. Brevdo, and H. Lee, "Sample-efficient reinforcement learning with stochastic ensemble value expansion," in *Proc. 31th Int. Conf. Neural Inf. Process. Syst.*, Dec. 2018, pp. 8234–8244.
- [105] D. P. Bertsekas, Rollout, Policy Iteration, and Distributed Reinforcement Learning, 1st ed. Athena Scientific, Aug. 2020.
- [106] A. Ilyas et al., "A closer look at deep policy gradients," in Proc. 8th Int. Conf. Learn. Representations, Nov. 2020, pp. 1–27.
- [107] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proc. 3rd Int. Conf. Learn. Representations*, Dec. 2015, pp. 1–15.
- [108] J. Jhung, H. Suk, H. Park, and S. Kim, "Hardware accelerators for autonomous vehicles," in *Artificial Intelligence and Hardware Accelerators*, A. Mishra, J. Cha, H. Park, and S. Kim, Eds. Cham: Springer International Publishing, 2023, pp. 269–317.
- [109] M. Mikami, K. Serizawa, Y. Ishida, H. Nishiyori, K. Moto, and H. Yoshino, "Field experimental evaluation on latency and reliability performance of 5G NR V2V direct communication in real express highway environment," in 2020 IEEE 91st Veh. Technol. Conf. (VTC2020-Spring). IEEE, May 2020, pp. 1–5.
- [110] B. Coll-Perales et al., "End-to-end V2X latency modeling and analysis in 5G networks," *IEEE Trans. Veh. Technol.*, vol. 72, no. 4, pp. 5094– 5109, Apr. 2023.
- [111] Study on enhancement of 3GPP Support for 5G V2X Services, TR 22.886 Version 16.2.0, 3rd Generation Partnership Project (3GPP), Sophia Antipolis, France, 2018, Dec.
- [112] D. Chen, K. Zhang, Y. Wang, X. Yin, Z. Li, and D. Filev, "Communication-efficient decentralized multi-agent reinforcement learning for cooperative adaptive cruise control," *IEEE Trans. Intell.* Veh., pp. 1–14, Feb. 2024.
- [113] L. Zhang, R. Zhang, T. Wu, R. Weng, M. Han, and Y. Zhao, "Safe reinforcement learning with stability guarantee for motion planning of autonomous vehicles," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 32, no. 12, pp. 5435–5444, Dec. 2021.
- [114] M. Han, Y. Tian, L. Zhang, J. Wang, and W. Pan, "Reinforcement learning control of constrained dynamic systems with uniformly ultimate boundedness stability guarantee," *Autom.*, vol. 129, p. 109689, May 2021.



Bernardo Camajori Tedeschini (Graduate Student Member, IEEE) is pursuing the Ph.D. degree in Information Technology at the Dipartimento di Elettronica, Informazione e Bioingegneria (DEIB), Politecnico di Milano, Milan, Italy, since November 2021. He received his M.Sc. (Hons.) degree in Telecommunications Engineering and B.Sc. (Hons.) degree in Computer Science from the Politecnico di Milano, Milan, Italy, in 2021 and 2019, respectively.

Currently, he is a Visiting PhD Researcher at the Wireless Information and Network Sciences Labo-

ratory, the Massachusetts Institute of Technology (MIT), Cambridge, MA, USA. In 2021, he has served as a Visiting Research Scientist at CERN, Geneva, Switzerland, where he worked on the CAFEIN project, focusing on the development and deployment of a Federated network platform. His research interests encompass federated learning, machine learning for signal processing and sensing over networks, and localization methods.

Mr. Camajori Tedeschini is a recipient of a Ph.D. grant from Italy's Ministero dell'Istruzione, dell'Università e della Ricerca (MIUR) and the Roberto Rocca Doctoral Fellowship, which was jointly awarded by MIT and Politecnico di Milano. He earned both his Bachelor's and Master's degrees with highest honors and he was honored with the best freshmen prize from Politecnico di Milano in 2017.



Moe Z. Win (Fellow, IEEE) is a Professor at the Massachusetts Institute of Technology (MIT) and the founding director of the Wireless Information and Network Sciences Laboratory. Prior to joining MIT, he was with AT&T Research Laboratories and with NASA Jet Propulsion Laboratory.

His research encompasses fundamental theories, algorithm design, and network experimentation for a broad range of real-world problems. His current research topics include ultra-wideband systems, network localization and navigation, network interfer-

ence exploitation, and quantum information science. He has served the IEEE Communications Society as an elected Member-at-Large on the Board of Governors, as elected Chair of the Radio Communications Committee, and as an IEEE Distinguished Lecturer. Over the last two decades, he held various editorial positions for IEEE journals and organized numerous international conferences. He has served on the SIAM Diversity Advisory Committee.

Dr. Win is an elected Fellow of the AAAS, the EURASIP, the IEEE, and the IET. He was honored with two IEEE Technical Field Awards: the IEEE Kiyo Tomiyasu Award (2011) and the IEEE Eric E. Sumner Award (2006, jointly with R. A. Scholtz). His publications, co-authored with students and colleagues, have received several awards. Other recognitions include the MIT Frank E. Perkins Award (2024), the MIT Everett Moore Baker Award (2022), the IEEE Vehicular Technology Society James Evans Avant Garde Award (2022), the IEEE Communications Society Edwin H. Armstrong Achievement Award (2016), the Cristoforo Colombo International Prize for Communications (2013), the Copernicus Fellowship (2011) and the *Laurea Honoris Causa* (2008) from the Università degli Studi di Ferrara, and the U.S. Presidential Early Career Award for Scientists and Engineers (2004).



Mattia Brambilla (Member, IEEE) received the B.Sc. and M.Sc. degrees in telecommunication engineering and the Ph.D. degree (cum laude) in information technology from the Politecnico di Milano, in 2015, 2017, and 2021, respectively.

He was a Visiting Researcher with the NATO Centre for Maritime Research and Experimentation (CMRE), La Spezia, Italy, in 2019. In 2021 he joined the faculty of Dipartimento di Elettronica, Informazione e Bioingegneria (DEIB) at the Politecnico di Milano as Research Fellow. His research interests

include signal processing, statistical learning, and data fusion for cooperative localization and communication.

Dr. Brambilla was the recipient of the Best Student Paper Award at the 2018 IEEE Statistical Signal Processing Workshop.



Monica Nicoli (Senior Member, IEEE) received the M.Sc. (Hons.) and Ph.D. degrees in communication engineering from Politecnico di Milano, Milan, Italy, in 1998 and 2002, respectively. She was a Visiting Researcher with ENI Agip, from 1998 to 1999, and Uppsala University, in 2001. In 2002, she joined Politecnico di Milano as a Faculty Member. She is currently an Associate Professor in telecommunications with the Department of Management, Economics and Industrial Engineering.

Her research interests include signal processing, machine learning, and wireless communications, with emphasis on smart mobility and Internet of Things (IoT). She was a recipient of the Marisa Bellisario Award, in 1999, and a co-recipient of the best paper awards of the EuMA Mediterranean Microwave Symposium, in 2022, the IEEE Symposium on Joint Communications and Sensing, in 2021, the IEEE Statistical Signal Processing Workshop, in 2018, and the IET Intelligent Transport Systems journal, in 2014. She is an Associate Editor of the IEEE Transactions on Intelligent Transportation Systems. She has also served as an Associate Editor for the EURASIP Journal on Wireless Communications and Networking, from 2010 to 2017, and a Lead Guest Editor for the Special Issue on Localization