

RL-PGO: Reinforcement Learning-based Planar Pose-Graph Optimization

Nikolaos Kourtzanidis and Sajad Saeedi

Abstract—The objective of pose SLAM or pose-graph optimization (PGO) is to estimate the trajectory of a robot given odometric and loop closing constraints. State-of-the-art iterative approaches typically involve the linearization of a non-convex objective function and then repeatedly solve a set of normal equations. Furthermore, these methods may converge to a local minima yielding sub-optimal results. In this work, we present to the best of our knowledge the first Deep Reinforcement Learning (DRL) based environment and proposed agent for 2D pose-graph optimization. We demonstrate that the pose-graph optimization problem can be modeled as a partially observable Markov Decision Process and evaluate performance on real-world and synthetic datasets. The proposed agent outperforms state-of-the-art solver g^2o on challenging instances where traditional nonlinear least-squares techniques may fail or converge to unsatisfactory solutions. Experimental results indicate that iterative-based solvers bootstrapped with the proposed approach allow for significantly higher quality estimations. We believe that reinforcement learning-based PGO is a promising avenue to further accelerate research towards globally optimal algorithms. Thus, our work paves the way to new optimization strategies in the 2D pose SLAM domain.

I. INTRODUCTION

Bundle adjustment (BA) is a process that plays a significant role in computer vision applications today, such as Structure from Motion (SfM) and on the back-end of Simultaneous Localization and Mapping (SLAM) algorithms [1]. This process involves estimating the 3D coordinates describing the scene and pose of the camera concurrently. This estimation is typically performed by optimally minimizing a cost function which minimizes re-projection error. In large complex scenes, however, the number of landmarks are significantly more prominent, which will, in turn increase the computation time and may affect real-time performance. When converting the problem of this type, one can either avoid taking into account all historical data from previous time steps and apply a sliding window technique or eliminate the optimization of the landmarks [2]. The latter is commonly referred to as pose-graph optimization. These relative pose measurements are typically obtained from cameras, laser sensors, IMUs, or wheel odometry. Given odometric or loop closing measurements at successive time intervals, the objective of PGO is to return the optimal configuration of pose estimates, which maximally explain the available measurements observed. The computation of the maximum likelihood estimate of robot poses results in a high-dimensional non-convex optimization problem with multiple local minima.

State-of-the-art methods in today’s modern era view SLAM and BA as Maximum a posteriori (MAP) or Maximum Likelihood inference. A specific class of graphs often analyzes both problems, known as factor graphs [3].

There exist instances of SLAM problems which consist of highly non-convex cost functions, attributed to system

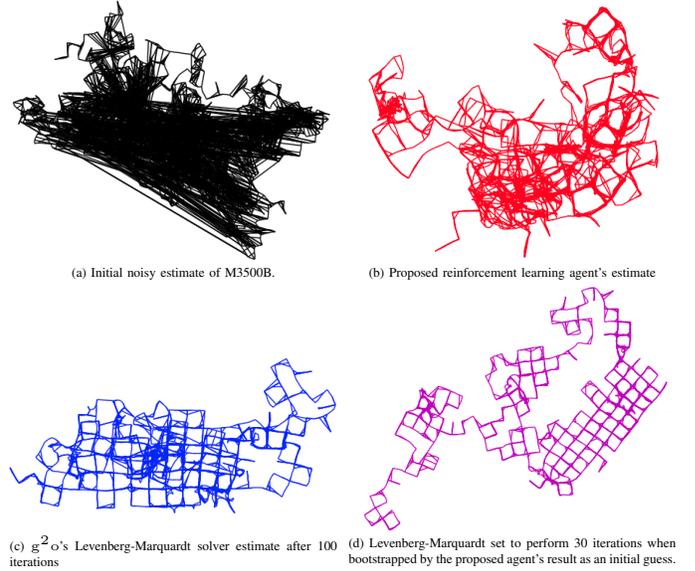


Fig. 1: Evaluation on the standard Manhattan world M3500B [6]. In this dataset, we notice that the estimate produced by Levenberg-Marquardt set to perform 100 iterations yields an unsatisfactory solution. When bootstrapped by our proposed agent’s initial guess, Levenberg-Marquardt set to complete just 30 iterations results in significantly improved accuracy.

nonlinearity, high levels of noise corruption, large inter-nodal distance spacing, and incorrect data associations from the front end [4]. This further results in cost functions with *large valleys* [5], which may cause classical gradient-based approaches to fail or converge at unsatisfactory solutions, as depicted in Fig. 1-(c).

This work, shows that our proposed reinforcement learning (RL) agent can perform exceptionally well on graphs with poor initial guesses. The agent can be utilized to search and explore the pose trajectory space configurations by applying the optimal retractions [7] directly to the desired poses, with the goal of returning and recovering the optimal global trajectory. In certain situations, obtaining ground truth trajectory labels in the SLAM domain may be laborious and expensive. The main advantage of leveraging a reinforcement learning-based agent is that existing supervised networks require to be trained on a large collection of labeled data. Suitable or empirical datasets used for label-based training are not always guaranteed [8].

Our contributions can be summarized as follows:

- The first DRL model to learn a policy that predicts the optimal orientation retractions [7] from pose-graph observations for the application of planar PGO. The proposed modular encoder-agent architecture is comprised of two main components. The first component is a neural network initially proposed in [9] for graph optimality classification. The second is a recurrent-based soft actor-critic (SAC) [10] agent whose policy predicts the optimal orientation retractions [7].
- A planar pose-graph environmental framework based on [11] and GUI for visualization during evaluation. The environment and the proposed agent will be made

publicly available¹ to encourage more research into developing and comparing other RL approaches in the pose-graph domain.

- Extensive experiments on simulated and real-world benchmarks illustrate the accuracy of the proposed agent, and further demonstrate the generalization capability of the agent to graphs never seen before during training.

The paper is organized as follows: Sec. II provides an overview of the related works. Sec. III includes preliminary definitions and related terminology. In Sec. IV, we present our proposed environment, and describe the encoder and agent architecture. Sec. V and VI present the experimental results, followed by the conclusions and future works.

II. RELATED WORK

Conventional Pose-Graph Optimization Approaches.

Popular MAP iterative solvers such as g^2o [12], Ceres [13], and GTSAM [14] can be explicitly utilized to tackle problems that can be represented as a graph using nonlinear optimization methods. Line search methods such as Stochastic Gradient Descent and Gauss-Newton (GN) are typically used to perform nonlinear optimization. Given an initial guess, second order methods repeatedly linearize a nonlinear least-squares objective function and then solve the unique least-squares equations until convergence [3]. Depending on the quality of this initial guess, there is no assurance for convergence to global optimal solutions [4]. As the non-convexities of the problem increase, iterative methods such as Powell’s dogleg, Levenberg-Marquardt (LM), and Gauss-Newton, are subject higher computational efforts required due to recurring matrix computations involved in the linearization step. This in turn, led to an increased interest in regards to direct and indirect linear solvers as well as factorization techniques to be introduced. Direct and indirect sparse linear solver methods include the preconditioned conjugate gradient (PCG), SuiteSparseQR [15], and CHOLMOD [16]. For larger scaled problems, bundle adjustment methods were introduced to perform at a much more accelerated rate. This is done by performing computations synchronously parallel, distributing the workload, and making use of the recent developments in hardware, i.e., IPUs. An example of this was implemented in [17] and [18].

Bootstrapping or initialization strategies were also shown to improve convergence of iterative methods. Multi Ancestor Spatial Approximation Tree for 2D and 3D pose-graph optimization [19] demonstrate a computationally efficient and light-weight initialization method, which when followed by a gradient-based solver, achieves good results on classical 2D and 3D pose-graph optimization datasets.

In recent years, the semidefinite relaxation (SDR) technique has been involved in exciting developments in the area of SLAM, and has also shown great significance on a variety of applications. In general, it can be applied to many non-convex quadratically constrained quadratic programs (QCQPs) such as the pose SLAM and Landmark SLAM problem instances.

SE-Sync [20] is a certifiably correct algorithm for performing synchronization over the special Euclidean group.

Leveraging duality to devise an algorithm was also shown to enable global verification of a given estimate [21]. The sparse bounded degree sum-of-squares (Sparse-BSOS) optimization method [22], formulate SLAM problems as polynomial optimization programs and demonstrate the ability to achieve global minimum solutions without initialization. A deep learning approach for pose-graph global optimality classification was also recently proposed in [9].

Learning-Based Optimization Approaches.

There has been a recent surge of interest on incorporating deep neural networks into state estimation and SLAM pipelines. Others have approached alternative solutions which are learned variations of the traditional approaches. LS-Net [23] introduced the first approach to a learned optimizer for minimizing photometric residuals. Several learning-based methods involve reparametrization to allow for gradient back propagation and end-to-end learning. An example of this was illustrated in [24] and [25], which involved a reparametrization of the damping mechanism to enable differentiability in Levenberg-Marquardt solvers. Recent line of works, present an end-to-end approach for learning estimators modeled as factor-graph smoothers for state estimation applications. A surrogate loss was proposed in [26] for end-to-end training of smoothing-based estimators, and was evaluated on Visual Tracking, and Visual Odometry tasks.

Graph Neural Networks.

There has also been a plethora of works involving Graph Neural Networks for multiple rotation averaging (MRA). NeuRoRA [27] were the first to apply a graph neural network to regress the rotations in a view graph. The model consisted of a cleaning network which was responsible for mitigating outlier measurements, followed by an additional network used for fine tuning and regressing the absolute rotations given the measurements. Using a single Message Passing Neural Network (MPNN) with edge attention, [28] further improved the robustness, training time, and execution speed of the previous deep learning model and eliminated the requirement for an additional network. Recently, PoGO-Net [29] introduce a novel joint loss MRA formulation which was shown to outperform state-of-the-art and operate in real-time.

Reinforcement Learning.

All the aforementioned learning techniques either require ground truth labels or apply techniques to enable differentiability for end-to-end training. Reinforcement learning revolves around the concept of learning from interaction [30]. It is different from the supervised and unsupervised machine learning paradigms in that it does not have an external supervisor to pair situations with labeled actions and does not try to propose a valid structure hidden in unlabelled situations. Instead, RL mainly focuses on maximizing a reward signal through interactions between the agent and an environment. It is essentially based on associating transitional experiences to actions in such a manner that would maximize a reward signal, and is used in several applications, e.g. healthcare [31], robotics [32], finance [33], and many more [34].

¹<https://sites.google.com/view/rl-pgo>

III. PROBLEM STATEMENT

In Sec. III-A, we provide relevant background knowledge and preliminary definitions in regards to PGO and Lie Theory. Related RL terminology is also further provided in Sec. III-B.

A. 2D Pose-Graph Optimization and Lie Theory

Assuming independent Gaussian measurement noise for relative orientation and translation measurements, i.e. $\Sigma_{ij} = \text{diag}(\sigma_t, \sigma_R)$, we formulate the problem for planar scenarios as follows: Given m relative pose measurements encoded by edges $(i, j) \in \mathcal{E}$ on a directed graph shared between n robot poses, and co-variance matrix for all measurement pairs, the goal of pose-graph optimization is to return the optimal estimate or configuration of poses which best fit the measurements observed. This can be more formally denoted as seeking the minimum of the objective cost function $F(\mathbf{x})$ [4]:

$$F(\mathbf{x}) = \sum_{(i,j) \in \mathcal{E}} \left\| \mathbf{R}_i^\top (\mathbf{t}_j - \mathbf{t}_i) - \tilde{\mathbf{t}}_{ij} \right\|_{\Sigma_{\sigma_t}}^2 + \sum_{(i,j) \in \mathcal{E}} \left\| (\mathbf{R}_j - \mathbf{R}_i) - \tilde{\mathbf{R}}_{ij} \right\|_{\Sigma_{\sigma_R}}^2, \quad (1)$$

where $\|\cdot\|_{\Sigma}^2$ symbolizes the Mahalanobis distance, and $F(\mathbf{x})$ is a unitless weighted squared sum of all error residuals used to determine the accuracy of the state estimation (smaller the better). $\mathbf{x} \in \mathbb{S}\mathbb{E}(2)$ denotes the set of all incremental pose coordinates. For each pose in the set, $\mathbf{x}_i = [\mathbf{t}_i^\top \mathbf{R}_i]^\top$, $\mathbf{t}_i \in \mathbb{R}^2$ and $\mathbf{R}_i \in \mathbb{S}\mathbb{O}(2), \forall i = 1 \dots n$ denote the absolute translation and orientation configurations. The relative translation and orientation measurements observed between neighboring poses i and j are represented as $\tilde{\mathbf{t}}_{ij} \in \mathbb{R}^2$ and $\tilde{\mathbf{R}}_{ij} \in \mathbb{S}\mathbb{O}(2)$. The level of difficulty in solving this optimization task lies in the estimation of the absolute orientations. It was shown in [35] and [4] that estimating rotations first in the case of pose-graph initialization, allows for greater convergence guarantees and closed form solutions in the planar case. Further proven in [4], the pose-graph optimization would be a linear least-squares problem if orientation states are known which allows for significant computational advantage. This principal is illustrated in our proposed method where a reinforcement learning agent seeks to recover the optimal orientation configurations for each robot pose in $\mathbb{S}\mathbb{O}(2)$, followed by a linear least-squares translation estimation step shown in [4]. When applying perturbations on the orientation states, one cannot just simply apply matrix addition, as this operation does take into consideration the wrap-around in the $\mathbb{S}\mathbb{O}(2)$ Lie Group. Instead, we apply updates via local reparametrization more commonly known as the exponential mapping, followed by the retraction operation which is denoted by the symbol \oplus [7]. In $\mathbb{S}\mathbb{O}(2)$, this is defined as

$$\mathbf{R}_i \oplus \xi \triangleq \mathbf{R}_i \cdot \exp \hat{\xi}, \quad (2)$$

where $\xi \in \mathbb{R}$, and the hat operator is responsible for the transformation from Cartesian vector space to the Lie algebra, which is a 2×2 matrix. The inverse of the exponential mapping is referred to as the logarithmic mapping, which transfers

elements from the Lie group to Cartesian coordinates denoted as: $\text{Log}(\mathbf{R}_i)$.

B. Reinforcement Learning Definitions

The basic components of reinforcement learning consist of an agent and an environment. The main objective of a reinforcement learning algorithm involves training an agent to learn an optimal policy $\pi(a|s)$, such that the agent can achieve a high cumulative reward under a user defined evaluation metric, known as the reward function. A policy is defined as a probabilistic mapping from states to actions. For every discrete time step of interaction t , the agent is responsible for making optimal sequential decisions by applying actions $\mathbf{a}_t \in \mathcal{A}$, observing a state $\mathbf{s}_t \in \mathcal{S}$, and then receiving a feedback signal or reward $R(\mathbf{s}_t, \mathbf{a}_t)$ for the corresponding state-action pair. We also define a Q value function which returns the expected sum of rewards assuming the agent is in state \mathbf{s} and performs actions \mathbf{a} following policy π .

The standard RL objective seeks to maximize the expected sum of rewards: $\mathbb{E}_\pi [\sum_t \gamma^t R(\mathbf{s}_t, \mathbf{a}_t)]$, where $\gamma \in [0, 1)$ is the weighting which defines how much importance we give for future rewards. It is also assumed that the state must include information about all aspects of the past interactions which enables informative decisions for the future [30]. When an observation does not contain the complete state information, the environment is said to be partially observable.

IV. PROPOSED APPROACH

The end-to-end agent-environment interaction is illustrated in Fig. 2. For each episode step, the orientation residual components of the input graph are first passed into the encoder network. The encoder then returns a highly expressive low dimensional input representation of the state. Once the state is then provided to the recurrent-based SAC policy network, the optimal retractions on neighbouring poses are applied. In Sec. IV-A and IV-B, we review the environment details and the framework used to encode our observations into a descriptive latent embedding. Sec. IV-C overviews the recurrent SAC policy network used to apply the optimal retractions on the graph orientation state space.

A. Environment Details

To frame the pose-graph optimization problem in terms of reinforcement learning, we need to first define an environment. The environmental framework was inspired from the synthetically generated pose-graphs provided in [4]. Each pose-graph instance is categorized by five environmental parameters: the number of poses n , orientation measurement uncertainty σ_R , translation measurement uncertainty σ_t , inter-nodal distance spacing d , and probability of loop closures lc .

As shown in Fig. 2, the agent's location along the graph is highlighted in red, where steps occur between each edge. The number of steps per episode is dependant on the number of edges in the graph and a user defined value referred to as the number of cycles. This value determines the number of times the agent traverses across each edge, and termination of

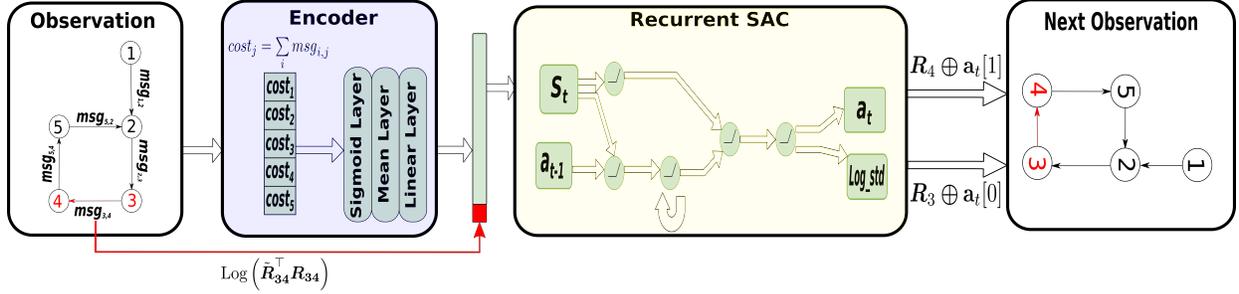


Fig. 2: Agent (red), cycling through every edge in the pose-graph and applying optimal retractions to the orientation components of each neighbouring pose.

an episode. The definition of one cycle consists of when an agent traverses every edge in the pose-graph once.

For each step-by-step transition along the edges, the output from the policy is passed through the \tanh activation layer and then multiplied by a user defined action range factor. The resultant vector is utilized to apply a retraction on the neighbouring poses connected by the edge, in which the agent resides at that particular instance in time: $R_i \oplus a_t[0]$ and $R_j \oplus a_t[1]$. In one cycle, each pose is perturbed twice assuming no loop closures. Following the completion of an episode, the final orientation estimate returned by the agent is then utilized in the concluding linear least-squares translation estimation provided in [4].

The state in this case is comprised of two components. The first component includes the encoded history of the orientation residuals at each time step. The second component is essentially the angular difference between the observed and measured orientations corresponding to the edge in which the agent is located at time t :

$$\text{Log} \left(\tilde{\mathbf{R}}_{ij}^\top \mathbf{R}_{ij} \right), \quad (3)$$

where $\mathbf{R}_{ij} \doteq \mathbf{R}_i^\top \mathbf{R}_j$, $(i, j) \in \mathcal{E}$. The graph observation encoding details are further discussed in Sec. IV-B.

The reward function for each time step is hybrid (dense/sparse) and is calculated as

$$\text{Reward} = \frac{100}{OC + 1}, \quad (4)$$

with an additive constant of +25 for every step the function output experiences a relative decrease by factor of 10, i.e. (0.001, 0.0001, 0.00001). The orientation cost (OC) is

$$OC = \sqrt{\sum_{(i,j) \in \mathcal{E}} \left\| \tilde{\mathbf{R}}_{ij} - \mathbf{R}_{ij} \right\|_F^2}, \quad (5)$$

where $\|B\|_F$, denotes the Frobenius norm of matrix B .

B. Graph Encoder Architecture

We adopt the architecture proposed in [9], formally presented for the task of global optimality prediction. The poses or in other words, nodes of each graph input, store a cost feature and the absolute orientation \mathbf{R}_i of the node itself. We utilize an augmented message passing function which is dependant on orientation components only, as opposed to both translation and orientation as depicted in Eq. (6). Further

provided by [9], once the input pose-graph observation passes forward, a two step process occurs. The first step is the message passing step which involves computing the Frobenius norm of the absolute orientations shared by each edge, where β is a learnable weight:

$$\text{msg}_{i,j} = \beta \times \left\| \mathbf{R}_i \mathbf{R}_{ij} - \mathbf{R}_j \right\|_F. \quad (6)$$

Proceeding the computation of messages for each of the connected nodes, an aggregate sum is then stored in the cost feature associated with the corresponding node itself,

$$\text{cost}_j = \sum_i \text{msg}_{i,j}. \quad (7)$$

The mean of all cost features from graph observation at time t is then passed through a linear layer and concatenated with the angular difference between the observed and measured orientations, corresponding to the agent's edge location computed from Eq. (3). The output dimension of the linear layer is user-defined and for all of our evaluations we set this value to 20. This resultant state vector \mathbf{s}_t at time t is passed into the recurrent SAC policy for an action to be applied.

C. Recurrent SAC

Soft actor-critic (SAC) [36], [37] is a state-of-the-art continuous control RL algorithm with an augmented objective:

$$\mathbb{E}_\pi \left[\sum_t \gamma^t [R(\mathbf{s}_t, \mathbf{a}_t) + \alpha \mathcal{H}(\pi(\cdot | \mathbf{s}_t))] \right], \quad (8)$$

where the temperature parameter α dictates the relative importance between the entropy \mathcal{H} and reward R . Thus, the agent is encouraged to explore the state space and unseen trajectories which can speed up learning and prevent sub optimal convergence. We utilize a recurrent version of the original implementation as shown in [10]. After observing the current state and history of previous state-action pairs, the recurrent SAC policy network applies the optimal retractions to each of the poses shared by the edge in which the agent is located at for every time step.

The recurrent SAC algorithm utilizes two networks, the Q-function and the policy. We will consider the policy π_ϕ , parameterized by ϕ modelled as a Gaussian with mean \mathbf{a}_t , and covariance as depicted in Fig. 2. The two dimensional action vector at time t is given by:

$$\mathbf{a}_t = \pi_\phi(\epsilon_t; \mathbf{s}_t, \mathbf{a}_{t-1}, \mathbf{z}_t), \quad (9)$$

where ϵ_t is an input noise vector sampled from a normal distribution with mean 0, and standard deviation of 1. Both policy and Q value architectures follow the structure of the recurrent-based DDPG algorithm proposed in [38].

The recurrent branch further allows the agent to make informed optimal decisions on the next actions from previous state and action pairs, by the internal state representation [38]: $z_t = z(h_t)$ modelled by an LSTM. We denote $h_t = [a_{t-1}, s_{t-1}, a_{t-2}, s_{t-2}, \dots]$ to represent the history of the previous states and actions.

The objective loss functions are essentially the same as presented in [37], with slight modification to the gradient estimators, accounting for back propagation through time (BPPT). Details can be referred to in [38] and [39].

V. EXPERIMENTAL RESULTS

In this section, we evaluate our RL approach and provide comparisons against the gradient-based g^2o [12] Gauss Newton, and Levenberg-Marquardt iterative solvers. As proven in [4], the non-convexity of a given problem instance are related to the ratio of orientation to translation uncertainty $\frac{\sigma_R}{\sigma_t}$, and the squared sum of all measurement distances. In Sec. V-B and V-C, we conduct analytical tests similarly done and inspired by [4], where comparisons were made against Gauss Newton 10 (GN10) and 5 (GN5) iterations in their evaluations. For our comparisons we set the maximum to 100. The quality of estimation provided by our approach is assessed under the influence of the adjustable environmental parameters, and further illustrate the effectiveness of the proposed approach under challenging scenarios. A head-to-head comparison with the Levenberg-Marquardt solver is also provided on standard real-world and synthetic benchmarks in Sec. V-D.

We train five separate agents, and for all evaluations depict results for the agent which performed best on the testing environment. Each of the agents are trained on small pose-graph environments of size $n = 20$ randomly sampled from their assigned environmental noise distributions every episode. The agents are then evaluated on synthetic and real-world graphs, much larger in size and unseen during training to demonstrate generalizability. The five training environment details are indicated in Table I, and cumulative reward plots for these agents during training are depicted in Fig. 3.

Env.	n	σ_R [rad]	σ_t [m]	d [m]	lc	Cycles	Action Range
1	20	0.3	0.01	3	0.5	7	0.25
2	20	0.3	0.01	3	0.5	8	0.25
3	20	0.2	0.1495	1	0.5	5	0.4
4	20	0.2	0.1495	1	0.5	6	0.25
5	20	0.1	0.01	10	0.5	6	0.25

TABLE I: Training environment parameters include the number of poses n , orientation measurement uncertainty σ_R , translation measurement uncertainty σ_t , inter-nodal distance spacing d , and probability of loop closures lc . See Sec. IV-A for more details.

A. Implementation and Training Details

For the policy and Q-value network, all related layers consist of 512 fully connected units followed by 512 LSTM units in the recurrent branch. Kaiming Initialization [40] was employed on both Q-function and policy network weights. Rectified non-linearity [41] (ReLU) were used for all hidden layers.

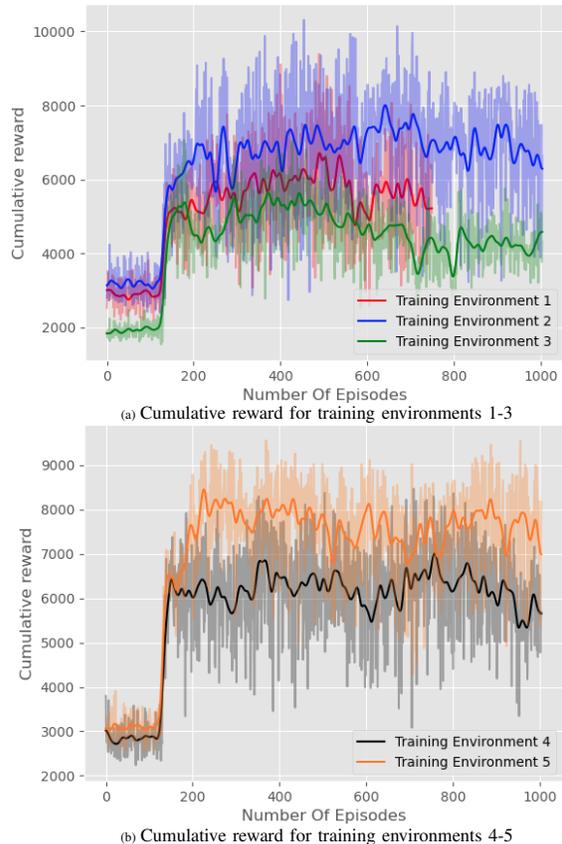


Fig. 3: (a)-(b) Cumulative reward plots for Training Environments 1-5.

We utilized Adam optimizer [42] with a learning rate of $3.0E - 04$, minibatch size of 128 and a discount factor of $\gamma = 1.0$. All Q-function and policy networks are updated after every episode with target networks updated by an exponential moving average, with $\tau = 1.0E - 02$. The networks are trained on a single Nvidia GeForce RTX 3090 GPU with 24GB memory. For best performance the final policy was set to deterministic and all objective cost function values $F(x)$, as well as optimization times recorded, are averaged over 10 runs or episodes.

B. Effect on Measurement Uncertainty Ratio

In this analytical case study, we allow Gauss Newton to perform 100 iterations (GN100) which can be interpreted as applying 100 pose retractions or *actions* to each pose in the graph. We utilize the agent trained in environment number 2 where the number of evaluation cycles was set to 8, and therefore performs far fewer retractions per pose than GN100. This analysis involves evaluation of the performances on a test pose-graph environment with decreasing translation uncertainty while all other parameters are kept fixed. In this case, $n = 300$, $\sigma_R = 0.3\text{rad}$, $d = 3\text{m}$, $lc = 0.5$, while σ_t ranges from $[0.2, 0.1, 0.05, 0.03, 0.01]\text{m}$. Objective cost function values and total elapsed optimization times are presented in Table II and an example of the $\sigma_t = 0.05\text{m}$ instance is shown in Fig. 4. It is observed that as the ratio of orientation to translation uncertainty increases, estimates produced by GN100 are extremely inaccurate as the non-convexities of the pose-graph optimization problem also increase. Standalone

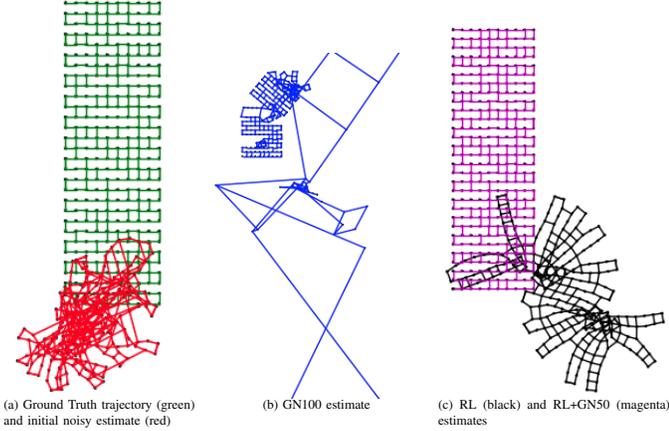


Fig. 4: Analysis on the influence of measurement uncertainty ratio for the instance $\sigma_t = 0.05m$. The RL+GN50 estimate (magenta) as shown in (c), only required 8 iterations for convergence and successfully recovered the global minimum solution. As the ratio of uncertainty increases, the RL standalone estimate (black) provides a much accurate estimation when compared to GN100 (blue).

RL was shown to outperform GN100 in situations where $\sigma_t = 0.05m$, $0.03m$, and $0.01m$. It is also noticed when utilizing the RL estimate as an initial guess, the ground truth trajectory and hence, global minimum, is returned in almost all cases in far fewer retractions per pose. Regardless of GN being set to perform 50 iterations, once bootstrapped by the RL estimate graph instances $\sigma_t = 0.2m$, $0.1m$, $0.05m$, and $0.03m$ converged in less than 11.

σ_t [m]	Metric	GN100	RL	RL+GN50
0.2	$F(x)$	0.00	6.04E+03	0.00
	\bar{time} [s]	0.04	10.08	10.12
0.1	$F(x)$	0.00	2.68E+04	0.00
	\bar{time} [s]	0.04	9.87	9.91
0.05	$F(x)$	1.59E+08	5.31E+04	0.00
	\bar{time} [s]	0.04	9.69	9.73
0.03	$F(x)$	5.76E+06	1.89E+05	0.00
	\bar{time} [s]	0.04	10.01	10.05
0.01	$F(x)$	1.47E+10	5.52E+06	1.21E+04
	\bar{time} [s]	0.04	9.97	10.01

TABLE II: Effect on measurement uncertainty ratio while other environmental parameters are held fixed.

C. Effect on Inter-nodal distance

In this section we conduct our evaluations on test pose-graph environments synthetically generated with various inter-nodal distance spacing d , while all other parameters are kept fixed. The test environment parameters in this case are, $n = 300$, $\sigma_R = 0.1rad$, $\sigma_t = 0.01m$, $lc = 0.5$, while d ranges from $[1, 3, 5, 8, 10]m$. Agent trained on environment 5 was utilized for evaluations and the final metrics are provided in Table III.

One may notice that standalone RL outperforms GN100 on instances with larger inter-nodal distance spacing. These results also conform to the analysis performed in [4] where it was further indicated that the inter-nodal distances are directly related to the sum of square measurements distances, which further increase the non-convexities of the pose-graph problem. GN10 bootstrapped by the RL estimate was also capable of achieving the global minimum cost in almost all instances. An example of the most challenging case $d = 10m$ is illustrated in Fig. 5.

d [m]	Metric	GN100	RL	RL+GN10
1	$F(x)$	0.00	6.82E+04	0.00
	\bar{time} [s]	0.04	6.96	7.0
3	$F(x)$	0.00	7.80E+05	0.00
	\bar{time} [s]	0.04	6.97	7.01
5	$F(x)$	1.17E+06	1.03E+06	0.00
	\bar{time} [s]	0.04	6.96	7.0
8	$F(x)$	1.58E+07	5.82E+06	0.00
	\bar{time} [s]	0.04	6.96	7.0
10	$F(x)$	3.35E+14	1.24E+07	4.31E+05
	\bar{time} [s]	0.04	6.96	7.0

TABLE III: Effect on inter-nodal distance spacing while other environmental parameters are held fixed.

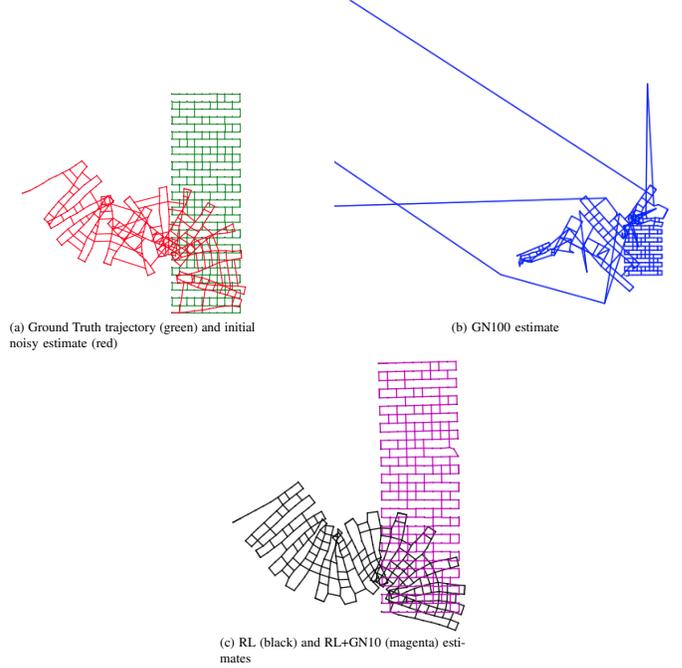


Fig. 5: Analysis on the influence of inter-nodal distance spacing for the instance $d = 10m$. The RL+GN10 estimate (magenta) as shown in (c), only required 8 iterations for convergence, to a visually meaningful solution. As the sum of total measurement distances increase, the proposed RL agent (black) provides a much accurate estimation when compared to GN100 (blue).

D. Standard Benchmark Datasets

To further assess the efficacy of the proposed approach, we evaluate performance on standard real-world and synthetic benchmarks provided by [6] and [4], never seen by our agent. Comparisons are made with g^2o 's LM solver fixed at 30 (LM30), and 100 (LM100) iterations, with stopping criteria based on reaching a relative error decrease or number of iterations. Additionally, we provide the objective cost function values and optimization time required for the LM30 estimate initially bootstrapped by our approach. The datasets include Manhattan world M3500 along with its three variants A, B, and C, which by default have standard deviations of 0.1rad, 0.2rad, and 0.3rad added to the relative orientation measurements of the original. Further, we evaluate on City10K [43], as well as Intel and MIT, which were obtained by processing raw measurements from wheel odometry and laser range finder measurements obtained at the Intel Research Lab in Seattle and Killian Court. In this assessment, we double the number of cycles the agent is set to perform during test time except for the evaluation on City10K. Increasing the number of cycles for evaluations was found to further improve the quality of the RL estimate due to a larger number of pose

Dataset (# of poses, # of edges)	Metric	RL	LM30	LM100	RL+LM30
M3500 (3500, 5453)	$F(x)$	1.97E+03	1.38E+02	1.38E+02	1.38E+02
	time [s]	870.02	-0.26	-0.46	870.28
M3500A (3500, 5453)	$F(x)$	2.94E+04	1.46E+05	9.68E+04	5.72E+03
	time [s]	745.29	-0.38	-1.51	745.67
M3500B (3500, 5453)	$F(x)$	4.66E+04	1.52E+04	1.43E+04	9.84E+03
	time [s]	744.70	-0.33	-1.48	745.03
M3500C (3500, 5453)	$F(x)$	5.64E+04	6.01E+04	3.74E+04	7.12E+03
	time [s]	744.01	-0.41	-1.63	744.42
City10K (10000, 20687)	$F(x)$	4.50E+04	4.53E+04	3.19E+04	5.12E+02
	time [s]	3914.25	-2.72	-9.85	3916.97
Intel (1228, 1483)	$F(x)$	9.96E+05	1.30E+05	5.24E+04	4.65E+02
	time [s]	122.07	0.09	0.27	122.16
MIT (808, 827)	$F(x)$	5.91E+03	1.11E+04	5.26E+02	7.71E+02
	time [s]	34.49	0.05	0.13	34.54

TABLE IV: Comparison Amongst Standard Benchmark Datasets.

refinements provided by the agent. Allowing the agent to perform more cycles however, come at the cost of longer optimization times required. This is apparent for graphs much larger in size such as City10K. In regards to evaluations on M3500A, B, C, City10k, and MIT, we employ the agent trained on environment number 3 which was found to perform best on those datasets. For the M3500 and Intel datasets, agents trained on environments 4 and 1 were utilized for evaluation. Objective cost function values, and total elapsed optimization times are given in Table IV.

It is observed that our proposed agent is effective in exploring the state space subject to highly non-convex cost functions. This is illustrated in instances M3500A, and C for example, which have a larger ratio of orientation to translation uncertainty due to the added standard deviation on the orientation measurements. In these graph instances the factors that influence non-convexity have a larger impact. In such cases, RL standalone was able to outperform LM30 in terms of state estimation accuracy. The RL estimate was also capable of achieving a smaller objective cost function value than LM100 on M3500A, further illustrated in Fig. 6. Due to the sequential step-by-step nature of the agent applying only two actions at once, far much more computation time is required for the optimization episode to complete in all cases.

In regards to dataset City10K, the number of edges in this graph and therefore, total sum of squared distance measurements, are also much larger than other datasets. Although RL standalone was able to outperform LM30 in City10K and MIT as well, by allowing more iterations, LM100 was eventually able to attain much accurate estimations. Nonetheless, when the RL estimate is utilized as a initial guess LM30 was capable of achieving even higher quality estimations than LM100 on all datasets except for MIT. We hypothesize that this is attributed to the fact that MIT has the fewest number of edges and loop closures, which provide the agent with less information on the residual state and decrease the number of actions applied. Further analysis which involved LM to run until convergence, indicate that the LM solver bootstrapped by our agent’s estimate was still of higher quality (see website of the project: <https://sites.google.com/view/rl-pgo>).

VI. CONCLUSION AND FUTURE WORKS

In this work, we have demonstrated a proposed RL agent to effectively explore the orientation state space in chal-

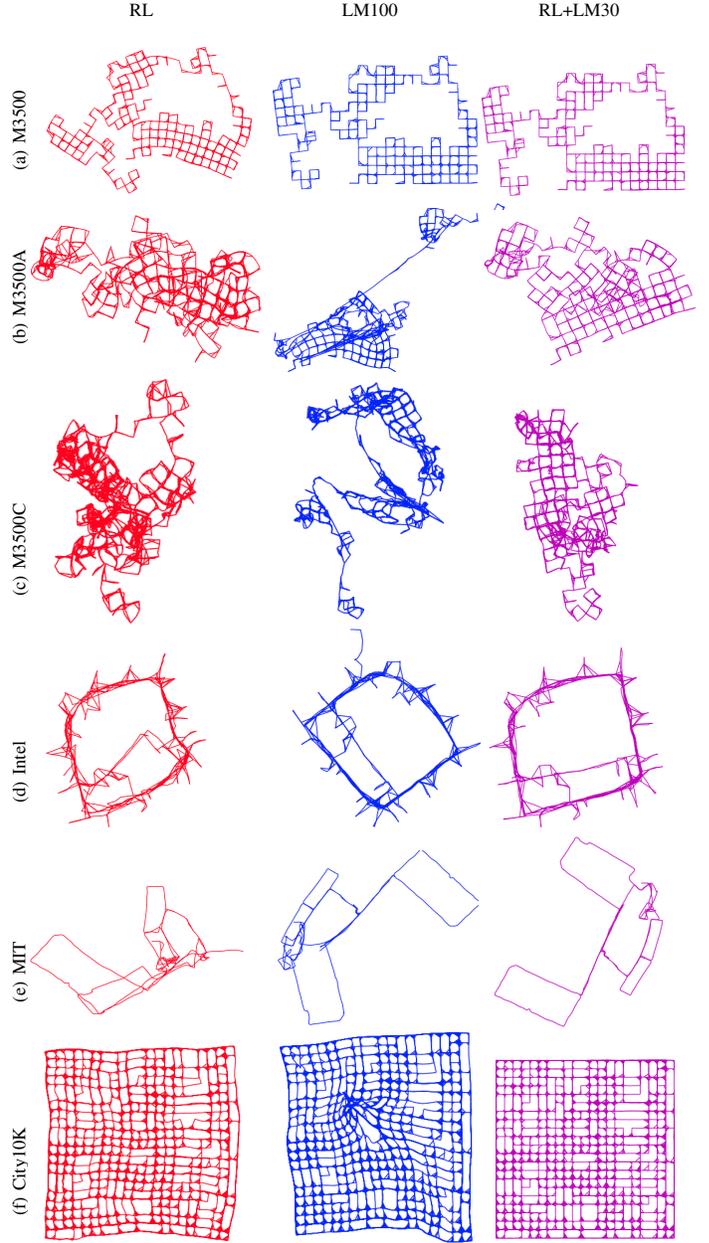


Fig. 6: Comparison amongst the standard real-world and synthetic graphs. Performance metrics are depicted from Table IV. From left to right: Standalone RL best estimate from 10 evaluations (red), LM100 estimate (blue), LM30 estimate when bootstrapped by the RL result as an initial guess (magenta). It is observed that RL standalone outperforms LM100 on M3500A. In all datasets except for MIT, RL+LM30 produced estimations with the highest quality (lowest objective cost function value). Interestingly, the standalone proposed RL agent was able to achieve solutions of adequate structure, despite having never seen any of these test graphs throughout training.

lenging pose-graph instances subject to highly non-convex cost functions for the application of planar pose-graph optimization. In particular scenarios our agent was capable of outperforming state-of-the-art Gauss Newton and Levenberg-Marquardt gradient-based solvers, in fewer retractions per pose. Nevertheless, although the RL-based approach was able to attain promising results, due to the small action space and inability to apply actions to all poses simultaneously the proposed agent does not seem to computationally scale well for graphs much larger in size. Thus, the methods presented in this work are more well-suited for offline Batch-SLAM applications. Our agent was shown to tolerate higher than usual

levels of noise and inter-nodal distance spacing, which further degrades the quality of estimation. Comparisons on simulated environments further exploit the factors which influence non-convexity [4], where state-of-the-art gradient descent-based solvers may catastrophically fail or return poor quality estimations. As shown in Fig. 3, the pose-graph optimization problem can be modelled as a partially observable MDP.

In situations where obtaining ground truth trajectories or labels may be a laborious and expensive process, we have demonstrated and highlighted the agent’s ability to perform well on much larger real-world and standard synthetic graphs, unseen during training. RL agents trained on smaller toy pose-graph environments of size $n = 20$ were able to generalize well to graphs with dissimilar noise distributions and odometric trajectories at test time.

Future works would involve extensions to pose-graph optimization instances in the 3D domain. Increasing action space size or a multi-agent framework may significantly reduce computational time. Moreover, a robust encoder architecture that is capable of mitigating the influence on outliers and false positive loop closures would allow for better standalone performance on real-world datasets. This may involve integrating robust kernels into the message passing function or de-noising layers adopted from [29].

REFERENCES

- [1] Y. Chen and G. Wang, “Bundle adjustment revisited,” *ArXiv*, vol. abs/1912.03858, 2019.
- [2] X. Gao, T. Zhang, Y. Liu, and Q. Yan, *14 Lectures on Visual SLAM: From Theory to Practice*. Publishing House of Electronics Industry, 2017.
- [3] F. Dellaert and M. Kaess, *Factor Graphs for Robot Perception*. Now Publishers Inc., 2017.
- [4] L. Carlone, R. Aragues, J. Castellanos, and B. Bona, “A fast and accurate approximation for planar pose graph optimization,” *International Journal of Robotics Research*, vol. 33, pp. 965–987, 2014.
- [5] E. Olson, J. J. Leonard, and S. J. Teller, “Fast iterative alignment of pose graphs with poor initial estimates,” *International Conference on Robotics and Automation*, 2006.
- [6] L. Carlone and A. Censi, “From angular manifolds to the integer lattice: Guaranteed orientation estimation with application to pose graph optimization,” *IEEE Transactions on Robotics*, vol. 30, no. 2, pp. 475–492, 2014.
- [7] J. Deray and J. Solà, “Manif: A micro Lie theory library for state estimation in robotics applications,” *Journal of Open Source Software*, vol. 5, p. 1371, 2020.
- [8] C. Cadena, L. Carlone, H. Carrillo, Y. Latif, D. Scaramuzza, J. Neira, I. Reid, and J. J. Leonard, “Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age,” *IEEE Transactions on Robotics*, vol. 32, p. 1309–1332, 2016.
- [9] R. Azzam, F. H. Kong, T. Taha, and Y. Zweiri, “Pose-graph neural network classifier for global optimality prediction in 2d slam,” *IEEE Access*, vol. 9, pp. 80466–80477, 2021.
- [10] Z. Ding, “Popular RL algorithms,” <https://github.com/quantumiracle/Popular-RL-Algorithms>, 2019.
- [11] J. Dong and Z. Lv, “miniSAM: A flexible factor graph non-linear least squares optimization framework,” *ArXiv*, vol. abs/1909.00903, 2019.
- [12] R. Kümmerle, G. Grisetti, H. Strasdat, K. Konolige, and W. Burgard, “G2O: A general framework for graph optimization,” *IEEE International Conference on Robotics and Automation*, 2011.
- [13] S. Agarwal, K. Mierle, and Others, “Ceres solver,” <http://ceres-solver.org>.
- [14] F. Dellaert, “Factor graphs and GTSAM: A hands-on introduction,” *GT RIM, GT-RIM-CP&R-2012-002*, 2012.
- [15] T. A. Davis, “Algorithm 915, SuiteSparseQR: Multifrontal multithreaded rank-revealing sparse QR factorization,” *ACM Trans. Math. Softw.*, vol. 38, no. 1, 2011.
- [16] Y. Chen, T. A. Davis, W. W. Hager, and S. Rajamanickam, “Algorithm 887: CHOLMOD, supernodal sparse Cholesky factorization and update/downdate,” *ACM Trans. Math. Softw.*, vol. 35, no. 3, 2008.
- [17] A. J. Davison and J. Ortiz, “Futuremapping 2: Gaussian belief propagation for spatial AI,” *arXiv*, vol. abs/1910.14139, 2019.
- [18] J. Ortiz, M. Pupilli, S. Leutenegger, and A. J. Davison, “Bundle adjustment on a graph processor,” in *Conference on Computer Vision and Pattern Recognition*, 2020.
- [19] K. Harsányi, A. Kiss, T. Szirányi, and A. Majdik, “Masat: A fast and robust algorithm for pose-graph initialization,” *Pattern Recognition Letters*, vol. 129, 11 2019.
- [20] D. M. Rosen, L. Carlone, A. S. Bandeira, and J. J. Leonard, “SE-Sync: A certifiably correct algorithm for synchronization over the special euclidean group,” *The International Journal of Robotics Research*, vol. 38, no. 2-3, pp. 95–125, 2019.
- [21] L. Carlone, G. C. Calafiore, C. Tommolillo, and F. Dellaert, “Planar pose graph optimization: Duality, optimal solutions, and verification,” *IEEE Transactions on Robotics*, vol. 32, no. 3, pp. 545–565, 2016.
- [22] J. G. Mangelson, J. Liu, R. M. Eustice, and R. Vasudevan, “Guaranteed globally optimal planar pose graph and landmark SLAM via sparse-bounded sums-of-squares programming,” in *International Conference on Robotics and Automation*, 2019.
- [23] R. Clark, M. Bloesch, J. Czarnowski, S. Leutenegger, and A. J. Davison, “LS-Net: Learning to solve nonlinear least squares for monocular stereo,” in *European Conference on Computer Vision*, 2018.
- [24] C. Tang and P. Tan, “BA-Net: Dense bundle adjustment network,” *ArXiv*, vol. abs/1806.04807, 2018.
- [25] K. M. Jatavallabhula, S. Saryazdi, G. Iyer, and L. Paull, “gradSLAM: Automagically differentiable SLAM,” in *International Conference on Robotics and Automation*, 2020.
- [26] B. Yi, M. A. Lee, A. Kloss, R. Martín-Martín, and J. Bohg, “Differentiable factor graph optimization for learning smoothers,” in *International Conference on Intelligent Robots and Systems*, 2021.
- [27] P. Purkait, T.-J. Chin, and I. Reid, “NeuRoRA: Neural robust rotation averaging,” in *European Conference on Computer Vision*, 2020.
- [28] J. Thorpe, R. Tennakoon, and A. Bab-Hadiashar, “Rotation averaging with attention graph neural networks,” in *Conference on Computer Vision and Pattern Recognition*, 2020.
- [29] X. Li and H. Ling, “PoGO-Net: Pose graph optimization with graph neural networks,” *International Conference on Computer Vision*, 2021.
- [30] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. The MIT Press, 2018.
- [31] C. Yu, J. Liu, and S. Nemati, “Reinforcement learning in healthcare: A survey,” *ArXiv*, vol. abs/1908.08796, 2020.
- [32] J. Kober, J. A. Bagnell, and J. Peters, “Reinforcement learning in robotics: A survey,” *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1238–1274, 2013.
- [33] T. G. Fischer, “Reinforcement learning in financial markets - a survey,” in *IDEAS*, 2018.
- [34] V. François-Lavet, P. Henderson, R. Islam, M. G. Bellemare, and J. Pineau, *An Introduction to Deep Reinforcement Learning*, 2018.
- [35] L. Carlone, R. Tron, K. Daniilidis, and F. Dellaert, “Initialization techniques for 3D SLAM: A survey on rotation estimation and its use in pose graph optimization,” in *International Conference on Robotics and Automation*, 2015.
- [36] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, “Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor,” in *International Conference on Machine Learning*, 2018.
- [37] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, and S. Levine, “Soft actor-critic algorithms and applications,” *ArXiv*, vol. abs/1812.05905, 2018.
- [38] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, “Sim-to-Real transfer of robotic control with dynamics randomization,” *International Conference on Robotics and Automation*, 2018.
- [39] N. Heess, J. J. Hunt, T. P. Lillicrap, and D. Silver, “Memory-based control with recurrent neural networks,” *ArXiv*, vol. abs/1512.04455, 2015.
- [40] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification,” *ArXiv*, vol. abs/1502.01852, 2015.
- [41] A. F. Agarap, “Deep learning using rectified linear units (ReLU),” *ArXiv*, vol. abs/1803.08375, 2018.
- [42] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *International Conference on Learning Representations*, 2015.
- [43] M. Kaess, A. Ranganathan, and F. Dellaert, “iSAM: Incremental smoothing and mapping,” *IEEE Transactions on Robotics*, vol. 24, pp. 1365–1378, 2008.