

Deep Active Localization

Sai Krishna^{1*}, Keehong Seo^{2*}, Dhaivat Bhatt¹, Vincent Mai¹, Krishna Murthy¹, Liam Paull¹

Abstract— Active localization is the problem of generating robot actions that allow it to maximally disambiguate its pose within a reference map. Traditional approaches to this use an information-theoretic criterion for action selection and hand-crafted perceptual models. In this work we propose an end-to-end differentiable method for learning to take informative actions that is trainable entirely in simulation and then transferable to real robot hardware with zero refinement. The system is composed of two modules: a convolutional neural network for perception, and a deep reinforcement learned planning module. We introduce a multi-scale approach to the learned perceptual model since the accuracy needed to perform action selection with reinforcement learning is much less than the accuracy needed for robot control. We demonstrate that the resulting system outperforms using the traditional approach for either perception or planning. We also demonstrate our approaches robustness to different map configurations and other nuisance parameters through the use of domain randomization in training. The code (available [here](#)) is compatible with the OpenAI gym framework, as well as the Gazebo simulator.

I. INTRODUCTION

Localization against a global map, or *global* localization, is a prerequisite for any robotics task where a robot must know where it is (e.g. any task involving navigation). In such settings, it is beneficial for the agent to first select actions in order to disambiguate its location within the environment (map). This is referred to as “*active* localization”.

Traditional methods for localization, such as Markov Localization (ML) [4] and Adaptive Monte Carlo Localization (AMCL) [29] are “*passive*” (agnostic to how actions are selected). They provide a recursive framework for updating an approximation of the state belief posterior as new measurements arrive. In the case of ML, this is usually done by some type of discretization of the state space (i.e. a fixed grid) and in AMCL is achieved by maintaining a set of *particles* (state hypotheses). While these methods have seen widespread success in practice, they are still fundamentally limited since the map representation is hand-engineered and specifically tailored for the given on-board robot sensor. A common example is the pairing of the occupancy grid map [6] with the laser scanner since the measurement likelihood can be computed efficiently and in closed form with scan matching [23]. However, this choice of representation can be sub-optimal and inflexible. Furthermore, sensor parameters such as error covariances tend to be hand-tuned for performance, which is time consuming and error-inducing.

Active variants of ML and AMCL are classical ways to solve the active localization problem [2], [3], [9], [18], [27].

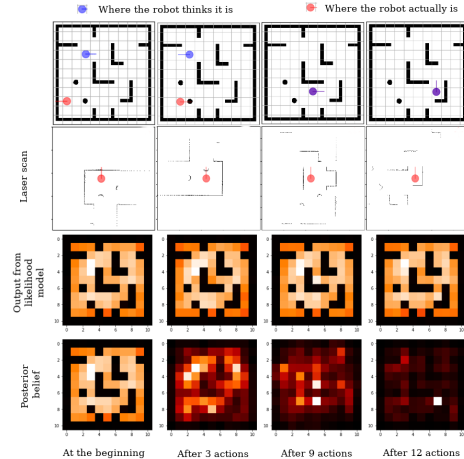


Fig. 1: DAL demonstration in simulation. At the beginning, the robot has no real clue of where it is relative to a known map, but as it executes actions and moves around, it converges onto a true location (see posterior belief improving in the fourth row). Furthermore, it transfers *zero-shot* onto a real robot, outperforming several baseline active localization methods. Refer to Sec. V for more details.

These methods typically leverage an information-theoretic measure to evaluate the benefit of visiting unseen poses. Again, these hand-tuned heuristics tends to lead to overfitting of a particular algorithm to a specific robot-sensor-environment setup.

Passive deep learning localization algorithms [1], [13]–[15] are not directly extendable to the active localization domain in the same fashion because they tend not to output calibrated measures of uncertainty.

However, other learning-based methods have been built that are dedicated specifically to this task, such as Active Neural Localization (ANL) [5]. In the ANL framework, a Deep Reinforcement Learning (DRL) policy model is trained on the outputs of a *fixed* measurement likelihood model (no learnable parameters), but still within a recursive Bayesian framework. However, several of the assumptions made in this work simply do not apply in real world robot deployments (for a more detailed analysis please see Sec. VI).

DRL has seen an incredible recent success on some robotics tasks such as manipulation [19] and navigation [33], albeit predominantly in simulation. In order for a DRL agent to be trained in simulation and deployed on a real robot either 1) the reality gap is small or is overcome with training techniques such as Domain Randomization (DR) such that no refinement of the policy is needed on the real robot (zero-

¹ Mila and DIRO, Universite de Montreal, Canada

² SAIT, Samsung Electronics Co. Ltd. Suwon, South Korea

* The first two authors contributed equally to this work.

shot transfer), or 2) the agent policy is fine-tuned on the real robot after primarily training in simulation. The latter option reduces the burden on the simulation fidelity and training regime, but fine-tuning on the robot can be difficult or impossible if the reward is determined by leveraging ground truth parameters (like localization) that are only available in the simulator. In this work, we argue that the framework of DRL holds potential for the task of active localization, and can be used to train policies that transfer to a robot in real-world indoor environments. Borrowing concepts from Bayes filtering [28], we define a reward as a function of the posterior belief. Specifically, our reward is defined as the probability of posterior belief at ground truth pose.

A. Contributions

In summary we claim the following contributions:

- We propose a multi-layer learned likelihood model which can be trained in simulation from automatically labeled data and then refined in an end-to-end manner,
- We show that this method works for *zero-shot* transfer onto a real robot and outperforms classical methods.
- We have integrated classical (Bayesian filtering, scan matching for ground truth likelihood) and learning-based (supervised and reinforcement learning) approaches appropriately to yield an active localization system that runs on a real robot.
- We developed automated processes for map generation, domain randomization and likelihood and policy model training.
- We provide an openAI-gym environment for the task of active localization to facilitate further research in this area. We have integrated it with SOTA deep RL algorithms to provide baselines
- We outline a set of domain randomization techniques and show that our learned likelihood model is more robust than classical hand-tuned techniques.

II. OVERVIEW OF APPROACH

In this section we will define our problem precisely and outline the structure of our proposed solution.

A. Problem Setup

A mobile robot is equipped with a sensor that provides exteroceptive measurements z_t used for localization and proprioceptive measurements used for odometry o_t (this could be the same physical sensor). The robot moves through the map by sending control inputs u_t to its actuators. The agent is provided with a map \mathcal{M} .

Problem (Active Localization) *Assuming that the agent is placed at some point in the map, find the sequence of control inputs, $u_{1:T}$ that allow it to maximally disambiguate its pose within the map.*

Active localization solutions usually take the form:

$$u_{1:T}^* = \operatorname{argmax}_{u_{1:T}} f(\operatorname{bel}(x_T), x^*) \quad (1)$$

Algorithm 1 Deep Active Localization

```

1: procedure DAL( $\mathcal{M}, f^{(0)}, f^{(1)}, \pi$ )
2:   Input:
3:    $\mathcal{M}$ : Map of the environment
4:    $f^{(0)}$ : Trained likelihood model for level-0
5:    $f^{(1)}$ : Trained likelihood model for level-1
6:    $\pi$ : Trained policy model
7:   while True do
8:     Read LiDAR scan  $z_t$ 
9:     Obtain high and low res scan images  $S_h$  and  $S_l$ 
10:     $p^{(0)}(z_t|x) \leftarrow f^{(0)}(M_h, S_h)$   $\triangleright$  Low-res likelihood
11:     $\operatorname{bel}^{(0)} \leftarrow p^{(0)}(z_t|x) \odot \operatorname{bel}^{(0)}$   $\triangleright$  Update belief
12:     $p(a|s) \leftarrow \pi(\operatorname{bel}^{(0)}, M_l, S_l)$   $\triangleright$  Get action prob
13:     $a \leftarrow \operatorname{sample}(p(a|s))$   $\triangleright$  Sample an action
14:     $c_t \leftarrow \operatorname{argmax}(\operatorname{bel}^{(0)})$   $\triangleright$  Get current pose
15:     $g_t \leftarrow \operatorname{get-next-pose}(c_t, a)$ 
16:     $\operatorname{bel}^{(0)} \leftarrow \operatorname{transition-belief}(\operatorname{bel}^{(0)}, a)$ 
17:    find top  $s$  cells in  $p^{(0)}(z_t|x)$ 
18:    for  $i \in [0, s]$  do
19:       $M_c, S_c \leftarrow \operatorname{crop}(M_h, S_h, (n_i, m_i))$ 
20:       $p^{(1)}(z_t|x)[\operatorname{block}_i] \leftarrow f^{(1)}(M_c, S_c)$ 
21:       $\operatorname{bel}^{(1)} \leftarrow p^{(1)}(z_t|x) \odot \operatorname{bel}^{(1)}$   $\triangleright$  High-res belief
22:       $c_t^{(1)} \leftarrow \operatorname{argmax}(\operatorname{bel}^{(1)})$   $\triangleright$  Pose: max belief
23:       $u \leftarrow \operatorname{control}(c_t^{(1)}, g_t)$ 

```

where the function f in some way quantifies the weight in the state belief posterior at the end of the horizon T , $\operatorname{bel}(x_T)$ at the ground truth pose x^* .

Any algorithm used to solve the optimization in (1) must jointly consider both the *perception* task (i.e., how the belief posterior is being updated) as well as the planning and control problem that is used to generate the control actions conditioned on the belief.

B. System Overview

Our approach to solving the active localization problem is summarized in Fig. 2 and Algorithm 1¹. For simplicity, we have shown two hierarchical levels in Fig. 2 since it is appropriate for our setup, but further levels of hierarchical could be added following the same approach (i.e., by providing: a measurement likelihood model, a planner, and a transition function) At each level in the hierarchy, i , we define a grid representation of size $N^{(i)} \times M^{(i)} \times \Theta^{(i)}$ which represents the state space \mathcal{X} . The belief posterior is represented as a matrix where $\operatorname{bel}(x_t = [n, m, \theta])$ is the probability mass at location $[m, n, \theta]$.

In Fig. 2, we assume that the neural network models are already trained, for a description of the training procedure for the measurement likelihood models see Sec. III-A and for the Reinforcement Learning (RL) policy see Sec. IV-A.

¹Conventions: superscripts $x^{(i)}$ denote levels in the hierarchy, subscripts denote time indices, and preceding superscripts $\{^r\}x$ denote reference frames. In the absence of a frame it is assumed that the variable is in the map fixed (global) frame.

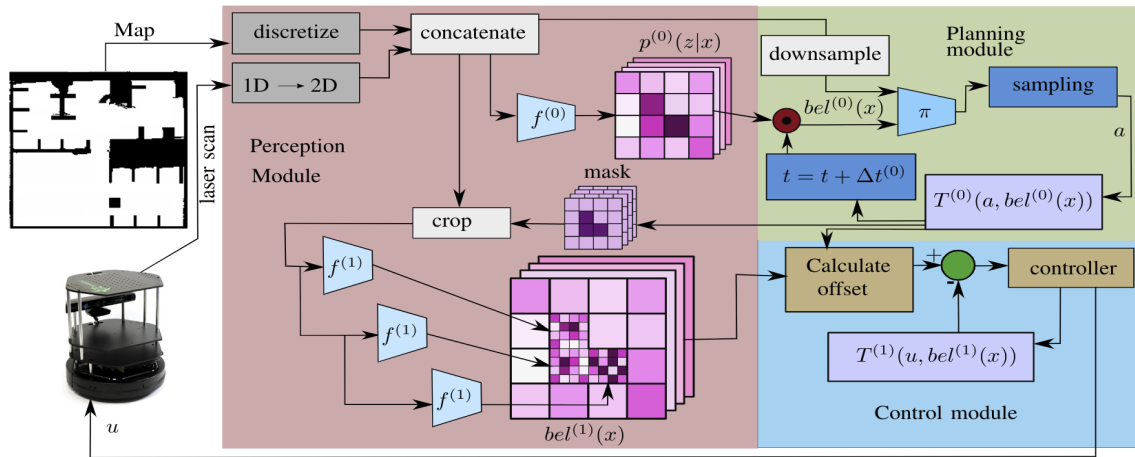


Fig. 2: **Deep Active Localization**: Our proposed method takes a map and sensor data (laser scan) and generates control actions in an end-to-end differentiable framework that includes learned perceptual modules (neural networks) at different scales and a learned policy network that is trained with reinforcement learning (RL). We maintain a coarse low dimensional pose estimate for RL, but then refine it to generate a higher precision pose estimate that is used for robot control.

The robot is provided with a map, \mathcal{M} as input. Each time a sensor input, z_t is received, it is converted to a 2D image and, combined with the \mathcal{M} to form the input to the network at level 0, $f^{(0)}(\mathcal{M}, z_t)$, which produces a coarse measurement likelihood, $p^{(0)}(z_t|x)$. The measurement likelihood is combined with the prior belief (via element-wise product) to produce the belief posterior. The belief posterior is fed as input to the RL model. The other inputs to the RL model are low dimensional map and low dimensional scan. It is important in order for this procedure to train efficiently that the input is relatively low dimensional. The RL policy is sampled to generate an action, a_t from the set $\{\text{left}, \text{right}, \text{straight}\}$. This action produces a goal pose that is either the centroid of an adjacent cell at the same orientation, or a pure rotation of $\pm 2\pi/|\Theta^{(0)}|$, where $\Theta^{(0)}$ is the number of uniformly spaced discrete angles being considered at level 0. The action and belief posterior are fed through a noisy transition function to generate the belief prior at the next timestep:

$$\begin{aligned} \bar{bel}^{(0)}(x_{t+1} = [n, m, \theta]) &= T^{(0)}(a_k, bel(x_t)) \\ &= \begin{cases} bel^{(0)}(x_t = [n, m, (\theta - 2\pi/|\Theta^{(0)}|)]) & a_k = \text{left} \\ bel^{(0)}(x_t = [n, m, (\theta + 2\pi/|\Theta^{(0)}|)]) & a_k = \text{right} \\ bel^{(0)}(x_t = [n + l_n \cos(\theta), m + l_m \sin(\theta), \theta]) & a_k = \text{str} \end{cases} \end{aligned} \quad (2)$$

where l_n and l_m are the distance between adjacent cell centroids in the $N^{(0)}$ and $M^{(0)}$ directions respectively. For increased robustness, noise is injected into this transition model to represent the fact that the cell transition is not deterministic.

Since each time the robot executes an action it will not arrive exactly at the centroid of the adjacent cell due to dead reckoning error, we must account for the error accrued and compensate for it. We refine our coarse measurement likelihood to predict where *within* the cell our agent actually is to compute a correction. This offset within the cell is used

to calculate an offset for the next action so that the error with respect to the grid centroids is *bounded over time*. To achieve this, we chose the “likely” cells from the coarse belief $bel^{(0)}(x)$ and refine the measurement likelihood to determine a more precise estimate of the location of the agent within the cell. This is used to calculate an offset to add to the relative pose transformations between adjacent cells to calculate an actual reference location in the robot frame, $r: {}^r x_{ref}$. We use a standard tracking controller to generate control inputs, u_t , which are sent to the robot and a standard kino-dynamic transition model $T^{(1)}(u, bel^{(1)}(x))$ to dead-reckon towards the reference location.

III. PERCEPTION

The objectives of the perception system are:

- 1) To provide a measurement likelihood to be used by RL agent, which is trainable in an end-to-end manner,
- 2) To provide a refined estimate of pose to the inner loop controller so that the error induced by dead reckoning may be bounded,
- 3) To be fully trainable in simulation,
- 4) Not to require any information other than the (possibly noisy) map at test time.

A. Learning Measurement Models

In typical robotics pipelines, the measurement likelihood is constructed using custom metrics, based on models of the sensors in question. Inevitably, some elements of the model are imprecise. For example, covariances in visual odometry models are typically tuned manually since closed-form solutions are difficult to obtain. Additionally, the Gaussian assumption (e.g., in an EKF) is clearly inappropriate for the task of global localization, as evidenced by the prevalence of Monte Carlo based solutions [29].

In our setting, given a map of the environment \mathcal{M} and scan input z_t , we want to *learn* the likelihood of the robot’s pose at all candidate points on a multi-resolution grid.

Following the notation defined in Sec II-B, the output of the likelihood model at time step t is given by:

$$p^{(0)}(z_t|x) = f_{\phi}^{(0)}(\mathcal{M}, z_t) \quad (3)$$

where ϕ are the parameters of the neural network model.

1) *Data Preprocessing*: It is difficult for a neural network to learn from different dimensional inputs. Though we can use different embeddings and concatenate at deeper layers, our experiments revealed that this is not very efficient. So, as a pre-processing step, we convert the scan input z_t into a scan image S_h (h denotes high resolution) of the same size as the grid map. We concatenate both the 2D scan image and the map of the environment and use it as input to our neural network. During the training phase, we obtain the ground truth likelihood $\tilde{p}^{(0)}(z_t|x)$ by taking the cosine similarity or correlation of the current scan with the scans at all other possible positions.

2) *Ground Truth Data Generation*: We generate and store triplets of (scan-image S_h , map of the environment M_h , ground truth likelihood $\tilde{p}^{(0)}(z_t|x)$) while randomly moving the robot in the simulator and randomly resetting the initial pose after every T time steps and randomly resetting the environment after every e such episodes. We train it on m such environments (maps). Training on these triplets using Resnet-152 [11] or Densenet-121 [12] gave very good results in simulation but did not transfer when these models are transferred to real robot in a real environment. This demonstrates the need for domain randomization to achieve zero-shot transfer. See III-D for further details.

B. Bayes' Filter

We can now update the belief by taking element wise product of likelihood and the previous belief over all the $N \times M \times O$ grid cells

$$bel^{(0)} = \bar{bel}^{(0)} \odot p^{(0)}(z_t|x)$$

C. Hierarchical Likelihood Model

We estimate the robot pose on a coarse grid, which is sufficient for RL, but we require a refined estimate for two reasons:

- 1) It's possible that the coarse estimate is simply not precise enough, particularly in very large maps
- 2) Without a more precise estimate of the robot's pose within a coarse grid cell, the robot will gradually drift from the grid centroids as a result of dead reckoning.

Traditional approaches try to solve this problem using scan matching technique, but it has the drawback of being computationally very expensive as the current scan information has to be compared with scan information at multiple poses.

Recent approaches like ANL [5] try to overcome this problem by using neural networks to predict the likelihood in a 3D grid of dimensions $N^{(1)} \times M^{(1)} \times \Theta^{(1)}$, where the whole environment is divided into $N^{(1)}$ rows, $M^{(1)}$ columns and $\Theta^{(1)}$ orientations. Each of these grid cells represent the likelihood of robot in that pose. But, this doesn't alleviate the problem of decoupling the localization precision from

the size of the map. With input being $N \times M \times 2$, and the output being $N \times M \times \Theta$, it is usually difficult and inefficient for a CNN to learn this mapping.

In HLE (Hierarchical Likelihood Estimation), likelihood is first estimated at coarse resolution ($N^{(0)} \times M^{(0)} \times O$) and then each of these grid cells is expanded to further finer grid cells ($k \times k$) to get a full resolution map of size ($N^{(1)} \times M^{(1)} \times \Theta$).

For the case of 2-level hierarchy, we have 2 neural networks $f^{(0)}(M_h, S_h)$ and $f^{(1)}(M_c, S_c)$ at levels 0 and 1 respectively. The input for $f^{(0)}$ is high the resolution map M_h (which is a processed version of map \mathcal{M} of the environment) ($N^0 \times M^0$), and scan image S_h ($N^0 \times M^0$) which is obtained from the current LiDAR scan z_t at the current pose. The likelihood output $p^{(0)}(z_t|x) = f^{(0)}(M_h, S_h)$ is of shape $N^{(0)} \times M^{(0)} \times \Theta^{(0)}$. The parameters of $f^{(0)}$ are optimized by minimizing the mean squared error (MSE) between $p^{(0)}(z_t|x)$ and ground truth likelihood $\tilde{p}^{(0)}(z_t|x)$ (which is at the same resolution of $N^{(0)} \times M^{(0)} \times \Theta^{(0)}$).

Grid cells corresponding to a maximum of c values from the coarse likelihood $p^{(0)}(z_t|x)$ are selected. For each of these grid cells, we crop a square patch from the high resolution map and high resolution scan (optional) around the grid cell. Concatenation of the cropped map and cropped scan is used as input for $f^{(1)}$. The network outputs a $k \times k$ block where each cell in the block represents likelihood at finer resolution. each of the cell in this block is multiplied with corresponding likelihood of the grid cell at previous level $p^{(0)}(z_t|x)$. For the cells which are not in the top c values, the likelihood value at level-0 $p^{(0)}(z_t|x)$ is normalized and directly copied to each cell in the corresponding block in $p^{(0)}(z_t|x)$. The parameters of $f^{(1)}$ are optimized by minimizing the mean squared error (MSE) between $p^{(1)}(z_t|x)$ and ground truth likelihood $\tilde{p}^{(1)}(z_t|x)$ (which is at the same resolution of $N^{(1)} \times M^{(1)} \times \Theta^{(1)}$).

D. Domain Randomization

DR has become a popular method to promote generalization, particularly when training an agent in a (necessarily imperfect) simulator, and then deploying it on real hardware [30]. While the LiDAR data modality generally transfers better than vision (camera images) from simulation to the real world scenario it is not without challenges since no sensor model is perfect. Just training the likelihood models without any variability in the simulator results in over-fitting and poor transfer. Since we do not have an exact recreation of the real world environment within our simulator, our learned agent will have to generalize to a new environment while simultaneously bridging the reality gap. Hence, it is important to account for various real world irregularities while training on the simulator. In our data collection pipeline, we randomized the following parameters:

- Thickness and length of obstacles to account for different types of obstacles in real environment.
- Error in robot pose to account for the possibility that it is not exactly at the centroid of a given cell.

- Temperature of softmax. It is important to normalize the ground truth likelihood to keep it bounded for a machine learning system to be able to learn. We use softmax with temperature β defined as:

$$\sigma(p^{(0)}(z_t|x))_{n,m,\theta} = \frac{e^{\beta p^{(0)}(z_t|x)_{n,m,\theta}}}{\sum_{n,m,\theta} e^{\beta p^{(0)}(z_t|x)_{n,m,\theta}}} \quad (4)$$

We vary the temperature parameter β randomly within limits to make sure that it is not biased towards overly-uniform or overly-sharp likelihood outputs

- Noise in the LiDAR scan: We train our network by adding Gaussian noise to every LiDAR scan point. Additionally, for every scan, we randomly set some of the incoming data points within the scan to have the value of $+\infty$.
- There are errors in the actual map of the environment created using gmapping [10]. This map is preprocessed and used as an input to the likelihood model and RL model for experiments on the real robot. Hence, it is important to ensure that some noise (more erosion and dilation of the map) is added to the map during training phase (on simulation) as well. However, note that the LiDAR readings will still come from the unperturbed map.

IV. PLANNING AND CONTROL

The multi-scale localization estimates are used for planning and control.

A. Reinforcement Learning

Following the traditional conventions in a Markov Decision Process [26], we denote the state $s \in \mathbb{R}^d$ and the actions $a \in \mathbb{R}^{d_a}$, the reward function $r(s, a)$ and a deterministic transition model $T = p(s'|s, a)$. In our case, the state (input to the RL model) is a concatenation of belief map $bel(x_t)$ ($\mathcal{N}^{(0)} \times M^{(0)} \times \Theta^{(0)}$) and the input map of the environment \mathcal{M} . We formulate the MDP over high-level actions (`left`, `right`, and `straight`). The goal of any RL algorithm is to optimize its policy $\pi(a|s)$ to maximize the discounted return defined as: $G_t = R_t + \gamma R_{t+1} + \dots$ over initial distribution of states (which is assumed uniform here). We use advantage actor critic (A2C) [22] algorithm to accomplish this task. There are various choices for a reward function. We used belief at true pose as our reward function since it is dense, well-behaved, and benefits from the availability of true pose in the simulator. The empirical evidence for the choice is given in Fig-3.

B. Closed-Loop Control

After the high level actions (`left`, `right`, `straight`) are chosen by the RL model, it is important to ensure that the robot reaches its goal position without deviating from the path. Even a minor deviation in each time step results in compounding errors. So, the low level actions (linear and angular velocities of the mobile robot) are given based on the current pose and the goal pose. The current pose is obtained at much finer resolution using our hierarchical likelihood

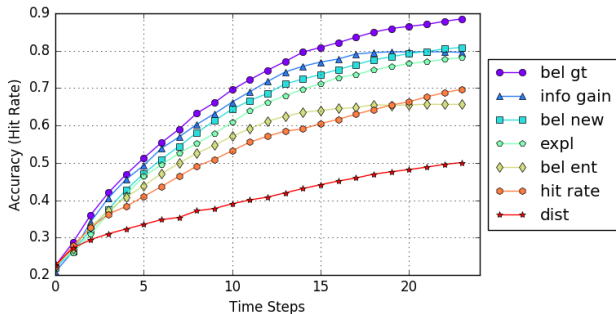


Fig. 3: We trained 7 RL policy models with different rewards for each: probability mass of belief at true pose (`bel gt`), decrease in the entropy of the belief (`info gain`), reward of +1 for a new pose in belief (`bel new`), reward of +1 for a new true pose (`expl`), negative entropy of belief (`bel ent`), reward of +1 if Manhattan distance error is equal to 0 (`hit rate`), Manhattan distance error (`dist`).

model. See the Sec. III-C for more details of hierarchical models. Optionally, the HLE model can be used after every fixed number of time steps to correct for the accumulated drift.

C. Zero-shot Transfer

The transferability of the perceptual model is enabled through the use of domain randomization (DR) as discussed in Sec. III-D. Another common use of DR is to randomize over physical properties of the robot (i.e., dynamics). This is not necessary in our case since we are performing RL at the *planning* level of abstraction and using a more traditional feedback controller to execute the plans.

V. EXPERIMENTS

We demonstrate the efficacy of DAL by conducting several experiments in simulation and on real robots and varying environments. This section describes the experimental setup and presents the results obtained, which demonstrate that DAL outperforms state-of-the-art active localization approaches in terms of localization performance and robustness.

The real environments used for testing were simple and re-configurable. Using a SLAM software we then obtained maps as shown in Fig. 6.

A. Dataset Generation

Using the following algorithm, we generated a dataset to train the likelihood model to estimate the measurement likelihood $p(z|x)$. The algorithm generates random maps, sensor data at some random poses as well as the target distributions at those poses.

- 1) We use Kruskal’s algorithm to generate random maze-like environments at low resolution. The output is a square matrix with 1s representing obstacles and 0s representing free cells. The algorithm guarantees that all the open spaces are connected.
- 2) Based on the grid map at low resolution generated above, we place obstacles (of different sizes) at the

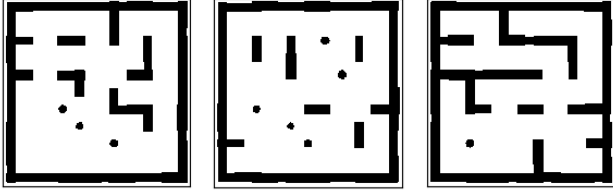


Fig. 4: Random maps generated from Kruskal’s maze algorithm with some of the walls pruned

cells in the high resolution. The algorithm then randomly dilates and erodes the surfaces of obstacles in the map so that its texture becomes more realistic.

- 3) We place the robot at all the grid centers of low resolution map at all $\Theta^{(0)}$ directions. From each of those locations, we simulate a laser scan. So, we obtain $M^{(0)} \times N^{(0)} \times \Theta^{(0)}$ laser scans. We call this as scan matrix and it’s of dimension $M^{(0)} \times N^{(0)} \times \Theta^{(0)} \times 360$. Given a map, we precompute a range vector at each location with respect to the low-resolution grids. A range vector in our experiment has the length of 360, representing distances at each of 360 degrees.
- 4) We now spawn the robot at 10 random locations and compute the cosine similarity of the scan at the current pose with the scan at all possible poses to get likelihood map. With each of those likelihoods, we form a triplet of map, current scan, likelihood and use these as training set

The procedure was repeated for 10000 different maps. The noise added to this process is described in III-D and V-B

B. Training Likelihood and Policy Models

We first describe how we trained our likelihood model (LM) which was used in our real robot experiments. We generated a dataset containing 10,000 maps with 10 scan inputs for each map by following the procedure explained above. The dimension of ground truth likelihood (GTL) was set to $33 \times 33 \times 8$, i.e, 8 headings with 33×33 x and y locations.

Likelihood model was trained with a densenet201 model on this dataset. We applied the following randomization for training LM:

- the temperature for softmax function was manually decreased over epochs from 1.0 to 0.1,
- input scan was randomly rotated in the range of $\pm 2^\circ$,
- 100 randomly selected pixels were flipped ($0 \leftrightarrow 1$),
- and drop-out was applied in densenet with the rate of 0.1.

Likelihood model was trained first with 100,000 data instances over 16 epochs and adapted to the map of the real environment in a simulator for 10,000 inputs.

The policy model π was trained on a simulator with randomly generated 2,000 maps, with 20 episodes (each of length 24) for each map. Reward (equal to probability mass of belief at true pose) was given at each time step.

C. Experimental Setup

To demonstrate the feasibility and robustness of our approach, we tested our trained likelihood and policy models on 2 mobile robots: JAY and Turtlebot and successfully localized in two different environments.

1) *Experiments on JAY*: JAY (Rainbow Robotics) is cylinder-shaped with the radius of 0.30m, the height of 0.50m and the weight of 46 kg before modification, whose CPU is IntelCore i7-8809G Processor. We mounted an extra LiDAR A2M8 (Slamtec) on the top of the robot in the center to secure 360-degree view. The LiDAR provides 0 to 360 degree angular range, and 0.15 to 8.0 m distance range, with 0.9 degree angular resolution, at the rate of 10 Hz.

The ROS navigation package including AMCL was manually initialized and used for ground truth localization. DAL was running on a separate server equipped with 4 Nvidia GTX-1080Ti GPUs and an Intel CPU E5-2630 v4 (2.1 GHz, 10 cores) to send velocity commands to JAY over a 2.4 GHz WIFI channel.



(a) JAY with top-mounted 360-degree view LiDAR (b) Turtlebot with top-mounted 260-degree view LiDAR

Fig. 5: Robots used in our experiments

With JAY deployed in the indoor environment shown as ‘env0’ in Fig. 8(a), we tested our likelihood model. To compare the localization performance with some of existing methods that run on grids, the experiments included the following methods.

- Trained likelihood model (LM): use the trained model for the sensor measurement likelihood at each pose.
- Scan Matching (SM): use cosine similarity between the precomputed scan matrix and the input scan for the sensor measurement likelihood.
- Reinforcement Learning (RL): use the trained policy model π for sampling next action.
- Active Markov Localization (AML): the robot virtually goes one step ahead along each of the possible actions and select the one with the largest reward, defined as the decrease in entropy.

- Random Action (RA): sample next action randomly with uniform probability.

For each test condition, it ran for 10 episodes of length 11. Fig. 6(a) shows the 10 initial poses that were randomly sampled to be applied to all the tests.

A custom control law computed velocity command at each step to move JAY to the next target pose by using the odometry feedback, where the target pose was computed from the believed pose $belx_t$ and the action sampled from the policy.

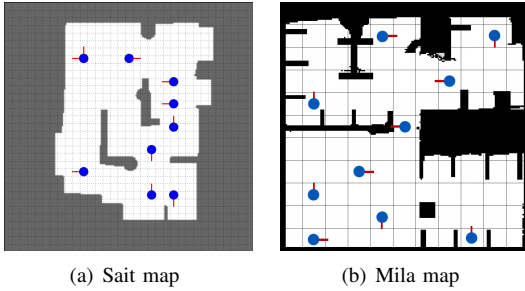


Fig. 6: The maps of the real environments in SAIT and Mila were generated by using a SLAM algorithm, *gmapping* [10]. The map size was adjusted to 224×224 pixels each of size $0.04m \times 0.04m$. We randomly selected 10 poses on the map and used them repeatedly as the starting poses for the 10 episodes of each test.

2) *Results from JAY Experiment:* We evaluated the test results using the following metrics:

- Earth Mover’s or Wasserstein’s distance W : It quantifies the error between our belief map and true belief map (which has probability 1.0 at true pose and 0 everywhere else)

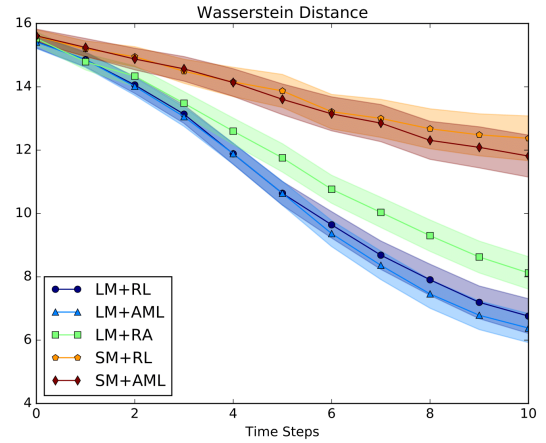
$$W = \sum_x p(x)D(x, x^*), \quad (5)$$

where we define D to be the Manhattan distance between two poses.

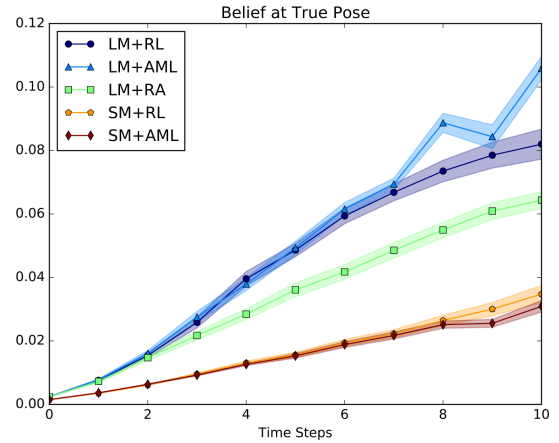
- estimated belief at true pose (This is also the reward that was used for training) $p(x^*)$
- hit rate, which is defined as the number of times the DAL pose is exactly same as true pose.

Each test consists of 10 episodes; the mean and standard deviation over the 10 episodes were then computed for each of 11 steps as plotted in 7 to illustrate how the metrics changed during an episode in average.

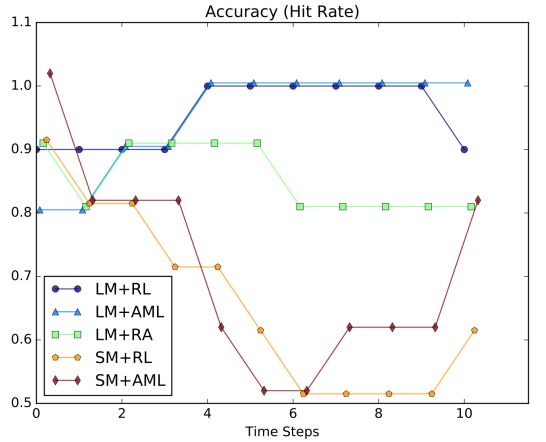
From the results plotted in Fig. 7, we can clearly see that our trained likelihood model performed much better than scan matching in all the metrics. This could be attributed to the susceptibility of SM to noises such as inaccurate map and the robot being off the center of a grid cell. Since we accounted for all these noises during training, our LM proved to be robust. Our RL model performed as good as AML. In terms of time complexity, our RL model is however multiple times faster than AML as we show in the time complexity analysis (cf. Table I)



(a) Wasserstein distance



(b) Belief at ground truth

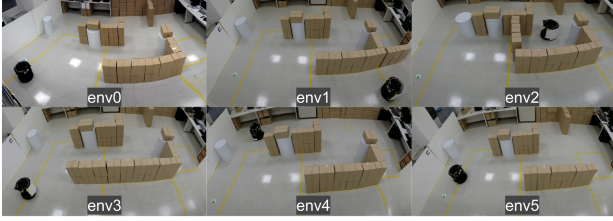


(c) hit rate

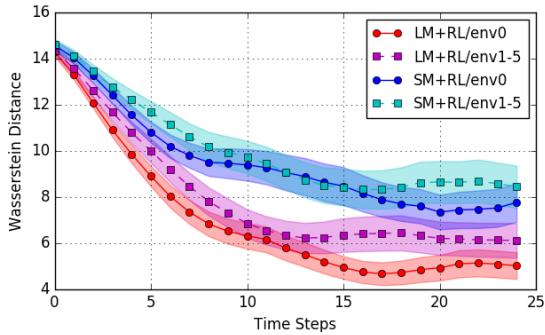
Fig. 7: Localization performance was evaluated in 8 headings and 33×33 grids with JAY in the real environment. We compared our proposed likelihood model for 33×33 (LM), scan-matching (SM), active Markov localization (AML), random action (RA), and the trained policy π with reinforcement learning (RL).

	SM	LM	RL	AML(+LM)
4x11x11	0.202	0.0531	0.00122	1.23
4x33x33	0.316	0.0518	0.00172	1.29
8x33x33	0.492	0.0866	0.00203	1.38
24x33x33	0.770	0.123	0.00194	1.47

TABLE I: Time complexity analysis (in seconds)



(a) The original and the 5 modified environments



(b) Wasserstein distance between belief and true belief

Fig. 8: LM is tested in modified environments env1 to env5 and compared with the error from the original environment env0. Error is measured in terms of Wasserstein distance between belief and true belief.

To further justify our robustness claims of our trained LM, we modified the environment by adding or moving some obstacles and ran the same experiment again. We tested it for LM+RL versus SM+RL on 5 differently modified environments from ‘env1’ to ‘env5’ shown in Fig. 8(a). We also tested again at the original environment ‘env0’. Each test was done with the same 10 initial poses as above while the length of an episode was increased to 25. Number of headings was fixed at 4.

3) *Experiments on Turtlebot*: We used Hukoyo Laser UST-20LX which outputs 1040 measurement ranges with detection angle of 260° and angular resolution of 0.25° . The netbook runs on Intel Core i3-4010U processor with 4GB RAM. Our test environment is of size 9×9 meters² and is combination of 2 rooms as shown in 6(b). Similar results as that of JAY were observed (omitted for brevity).

D. gym-dal

Most of the robotics research involves the extensive use of ROS and Gazebo environments which are convenient but they are not ideal platform for implementing learning algorithms especially reinforcement learning. To alleviate these problems, we are open sourcing our gym environment.

We hope that this will act as a testbed for active localization research. Once trained on our simulator, the policies and the likelihood model learned are directly transferable to real robot. We have also integrated the environment with various state of the art RL algorithms [17] like PPO [25], A2C [22], ACKTR [32].

E. Hierarchical Likelihood models

We present the experimental results of our 2 novel hierarchical likelihood models. From Fig. 9(a), we can observe that both weighted and un-weighted variants of Hierarchical likelihood model performed equally well and obtained convergence. However, higher levels of hierarchy are very sensitive to hyper-parameters because they are dependent on the performance of previous layers. The robustification studies of the hierarchical model, a thorough empirical evaluation and making them insensible to hyper parameters are possible avenues for future research.

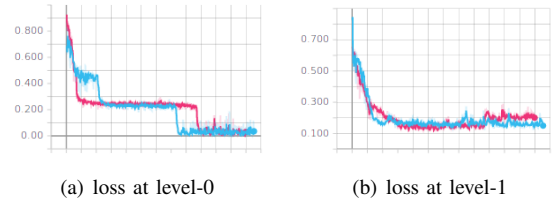


Fig. 9: blue: weighted HLE. red: unweighted HLE

VI. RELATED WORK

Localization is an extremely well-studied robotics problem. We present a treatment of relevant *active* and *passive* approaches to robot localization, drawing parallels and contrast to our method.

A. Passive localization

Passive localization methods can broadly be classified into two representative categories: *Bayesian* filtering methods and learning-based methods. Although related, we do not cover visual place recognition approaches as their primary intention is to provide *coarse* location estimates. We refer the interested reader to [20] for an excellent survey.

Bayesian filtering-based methods: Off-the-shelf localization modules such as Extended Kalman Filter (EKF)-based localization [28], Markov localization [4] and Adaptive Monte-Carlo localization (AMCL) [29] are based on a Bayesian filtering viewpoint of localization. All these methods begin with a prior (usually uniform) *belief* of a robot’s pose and recursively update the belief as new (uncertain) measurements arrive. The representation of the belief is what fundamentally differentiates each approach. EKF-based methods assume a Gaussian distribution to characterize the belief, while Markov localization [4] and AMCL [29] use non-parameteric distributions to characterize the belief (histograms and particles respectively). Consequently, EKF localization lacks the capability to characterize multiple belief hypotheses (owing to the unimodal Gaussian prior),

while Markov localization and AMCL suffer computational overhead as the population of the non-parametric distribution increases.

Learning-based methods: Recently there has been a resurgence of learning-based methods for localization. Approaches such as PoseNet [15] and VLocNet [1] perform visual localization by training a convolutional neural network to regress to *scene coordinates*, given an image [15] or a sequence of images [1]. Recently, differentiable particle filters (DPFs) have been proposed for global localization [13], [14]. However, such approaches need precisely annotated data to train a *PoseNet*, *VLocnet*, or a *DPF* for each new deployed environment. In contrast, DAL transfers essentially *zero-shot* across simulated environments, and from a simulator to the real-world.

B. Active localization

In terms of *active* localization methods, there are again two broad categories: information-theoretic approaches and learning-based approaches.

Information theoretic methods

Burgard et al. introduced *active* localization in their seminal work [3]. They demonstrated that, rather than passively driving a robot around, picking actions that reduce the expected localization uncertainty results in a more efficient and robust solution. Using an entropy measure characterized as a mixture-of-Gaussians, they demonstrate that the framework of Markov localization [4] can be extended to action selection. This was more a *proof-of-concept*, as the applicability of this method is confined to low-dimensional state-spaces, where entropy computation can be carried out efficiently. Since then, several other approaches [7], [8], [16], [21], [24], [31] use a similar, information gain maximization cost for the task of active localization or SLAM.

Learning-based methods

Active Neural Localization (ANL) [5] is the first known work to employ a learned model for active localization from images. Their approach comprises two modules: a perceptual model and a policy model. The perceptual model is not completely learned. The transition functions here are assumed to be deterministic, and this prohibits the policy model to be transferred onto a real robot. On the other-hand, DAL comprises a likelihood and a policy model, both of which are learned from data, which outperform their traditional counterparts. The hierarchical likelihood estimation also allows for scalability to much bigger environments compared to ANL.

VII. CONCLUSION

We have presented a learning-based approach to active localization from a known map. We propose multi-scale learned perceptual models that are connected with an RL planner and inner loop controller in an end-to-end fashion. We have demonstrated the effectiveness of the approach on real robot settings in two completely different setups.

In future work, we would like to explore the removal of the reliance of the system on a high fidelity map such as is generated by *gmapping*. Much more appealing would be, for example, if the robot could localize on a hand drawn sketch or some much more easily obtained representation.

REFERENCES

- [1] Noha Radwan Abhinav Valada and Wolfram Burgard. Deep auxiliary learning for visual localization and odometry. In *ICRA*, May 2018.
- [2] Tal Arbel and Frank P Ferrie. Viewpoint selection by navigation through entropy maps. In *ICCV*, 1999.
- [3] Wolfram Burgard, Dieter Fox, and Sebastian Thrun. Active mobile robot localization. In *IJCAI*, 1997.
- [4] Wolfram Burgard, Dieter Fox, and Sebastian Thrun. Markov localization for mobile robots in dynamic environments. *CoRR*, 2011.
- [5] Devendra Singh Chaplot, Emilio Parisotto, and Ruslan Salakhutdinov. Active neural localization. *CoRR*, 2018.
- [6] A. Elfes. Using occupancy grids for mobile robot perception and navigation. *Computer*, 22(6):46–57, 1989.
- [7] Hans Jacob S. Feder, John J. Leonard, and Christopher M. Smith. Adaptive mobile robot navigation and mapping. *IJRR*, 1999.
- [8] C. Forster, M. Pizzoli, and D. Scaramuzza. Appearance-based active, monocular, dense reconstruction for micro aerial vehicle. *RSS*, 2014.
- [9] Dieter Fox, Wolfram Burgard, and Sebastian Thrun. Active markov localization for mobile robots. *Robotics and Autonomous Systems*, 1998.
- [10] G. Grisetti, C. Stachniss, and W. Burgard. Improved techniques for grid mapping with rao-blackwellized particle filters. *IEEE Transactions on Robotics*, 23(1):34–46, Feb 2007.
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, 2015.
- [12] Gao Huang, Zhuang Liu, and Kilian Q. Weinberger. Densely connected convolutional networks. *CoRR*, 2016.
- [13] Rico Jonschkowski, Divyam Rastogi, and Oliver Brock. Differentiable particle filters: End-to-end learning with algorithmic priors. *CoRR*, 2018.
- [14] Péter Karkus, David Hsu, and Wee Sun Lee. Particle filter networks: End-to-end probabilistic localization from visual observations. *CoRR*, 2018.
- [15] Alex Kendall, Matthew Grimes, and Roberto Cipolla. Convolutional networks for real-time 6-dof camera relocalization. *ICCV*, 2015.
- [16] Ayoung Kim and Ryan M. Eustice. Active visual slam for robotic area coverage: Theory and experiment. *IJRR*, 2015.
- [17] Ilya Kostrikov. Pytorch implementations of reinforcement learning algorithms. GitHub repository, 2018.
- [18] Rainer Kummerle, Rudolph Triebel, Patrick Pfaff, and Wolfram Burgard. Monte carlo localization in outdoor terrains using multi-level surface maps. *J. Field Robotics*, 2008.
- [19] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *CoRR*, 2015.
- [20] Stephanie Lowry, Niko Sünderhauf, Paul Newman, John J Leonard, David Cox, Peter Corke, and Michael J Milford. Visual place recognition: A survey. *IEEE Transactions on Robotics*, 2016.
- [21] G. L. Mariottini and S. I. Roumeliotis. Active vision-based robot localization and navigation in a visual memory. In *ICRA*, 2011.
- [22] Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. *CoRR*, 2016.
- [23] Edwin B. Olson. Real-time correlative scan matching. In *ICRA*, 2009.
- [24] N. Roy, W. Burgard, D. Fox, and S. Thrun. Coastal navigation-mobile robot navigation with uncertainty in dynamic environments. In *ICRA*, 1999.
- [25] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *CoRR*, 2017.
- [26] Richard S. Sutton. On the significance of markov decision processes. In *ICANN*, 1997.
- [27] Sebastian Thrun. Finding landmarks for mobile robot navigation. In *ICRA*, 1998.
- [28] Sebastian Thrun, Wolfram Burgard, and Dieter Fox. *Probabilistic Robotics (Intelligent Robotics and Autonomous Agents)*. The MIT Press, 2005.

- [29] Sebastian Thrun, Dieter Fox, Wolfram Burgard, and Frank Dellaert. Monte carlo localization for mobile robots. In *ICRA*, 1999.
- [30] Joshua Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. *CoRR*, 2017.
- [31] R. Valencia, J. Valls Miro, G. Dissanayake, and J. Andrade-Cetto. Active pose slam. In *IROS*, 2012.
- [32] Yuhuai Wu, Elman Mansimov, Shun Liao, Roger B. Grosse, and Jimmy Ba. Scalable trust-region method for deep reinforcement learning using kronecker-factored approximation. *CoRR*, 2017.
- [33] Yuke Zhu, Roozbeh Mottaghi, Eric Kolve, Joseph J Lim, Abhinav Gupta, Li Fei-Fei, and Ali Farhadi. Target-driven visual navigation in indoor scenes using deep reinforcement learning. In *ICRA*, 2017.