

CAPABILITY ITERATION NETWORK FOR ROBOT PATH PLANNING

Buqing Nie¹, Yue Gao^{*2}, Yidong Mei³ and Feng Gao⁴

Abstract

Path planning is an important topic in robotics. Recently, value iteration based deep learning models have achieved good performance such as Value Iteration Network(VIN). However, previous methods suffer from slow convergence and low accuracy on large maps, hence restricted in path planning for agents with complex kinematics such as legged robots. Therefore, we propose a new value iteration based path planning method called Capability Iteration Network(CIN). CIN utilizes sparse reward maps and encodes the capability of the agent with state-action transition probability, rather than a convolution kernel in previous models. Furthermore, two training methods including end-to-end training and training capability module alone are proposed, both of which speed up convergence greatly. Several path planning experiments in various scenarios, including on 2D, 3D grid world and real robot with different map sizes are conducted. The results demonstrate that CIN has higher accuracy, faster convergence and lower sensitivity to random seed compared to previous VI-based models, hence more applicable for real robot path planning.

Key Words

path planning, value iteration, value iteration network, capability iteration network

I. INTRODUCTION

Path planning is an important direction for autonomous robot, whose objective is to compute a path to the goal region that not only avoids collisions with obstacles but also satisfies kinematic constraints of the robot [1], [2]. Various classic methods including A*

[3], Potential Fields [4] and Rapidly exploring Random Tree [5] are proposed. Recently, Deep Reinforcement Learning(DRL) has been utilized to solve path planning problems without handcraft parameters tuning according to the robot topological structures [6], for example on robotic manipulator [7], [8], drilling robot [9], planetary rover [10], and long-range navigation task [11].

DRL has the benefits of representing complex state space and learning a network which enables function approximation of the features and future rewards values [12]. However, most existing DRL based path planning models lack planning modules which makes them rely on large amounts of training data and difficult to generalize for unseen environments [13]. This makes DRL models difficult to be applied in real scenarios, like planning paths for real robots. A recent work, *Value Iteration Network*(VIN) [13] combines recurrent convolution neural networks and max-pooling to emulate Value Iteration(VI) [14] algorithm, which enables VIN to learn and plan paths for unseen mazes. However, VIN is particularly challenging on training when given large maps. In addition, how to utilize VIN to solve path planning for robots with complex kinematics remains a challenging task.

For RL models, sparse reward is an important property required for convergence and generalization [15]–[17]. Current VIN method utilizes convolution layers to compute reward map $R(s,a)$, which is not sparse and may decrease the accuracy of VIN because of its defects, especially on large maps. Besides, VIN utilizes concatenation and a global convolution kernel P_s^a (probability distribution for the next state only based on the action to take) instead of the state conditioned transition probability $P_{ss'}^a$ in the Value Iteration algorithm.

In this paper, a new VI-based model *Capability Iteration Network* (CIN) is proposed. In CIN, sparse reward maps instead of reward maps generated by CNN are utilized, which limit reward maps influencing the accuracy of the model, especially on large maps. Besides, as illustrated in fig.1, a capability module is utilized to encode each state s with a state-conditioned transition probability $P_{ss'}^a$, which is similar to the environment dynamics utilized in model based RL. These state-conditioned transition probabilities in CIN represents the capability of the agent, which differentiates itself

¹Buqing Nie is with Department of Automation, Shanghai Jiao Tong University, Shanghai, P.R. China, niebuqing@sjtu.edu.cn

²Yue Gao is with MoE Key Lab of Artificial Intelligence and Department of Automation, Shanghai Jiao Tong University, Shanghai, P.R. China, yuegao@sjtu.edu.cn

³Yidong Mei is with Department of Automation, Shanghai Jiao Tong University, Shanghai, P.R. China, meiyidong@sjtu.edu.cn

⁴Feng Gao is with State Key Laboratory of Mechanical System and Vibration, Shanghai Jiao Tong University, Shanghai, P.R. China, fengg@sjtu.edu.cn

*Corresponding author.

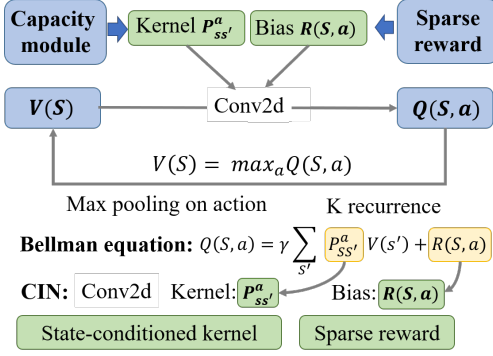


Fig. 1: Structure of the *Capability Iteration Network* (CIN). CIN utilizes sparse reward maps and state-conditioned transition probabilities $P_{ss'}^a$ instead of reward generated by CNN and global transition probabilities P_s^a in VIN.

from previous VI based models that utilize P_s^a without the knowledge of the current state. CIN is trained to learn the capability of the agent, which is relevant to the local region and independent with the global map. Therefore, CIN is able to be trained on local maps with faster learning speed. In the experiment section, several experiments demonstrate that CIN can be applied to scenarios including 2D mazes, large-scale 3D terrain maps and path planning for real hexapod robots.

The main contributions of this paper are as follows:

- A new VI-based model called CIN is proposed, which greatly improves performance on accuracy and learning speed compared with previous VI-based algorithms on various scenarios.
- We provide two training methods: end-to-end training and capability module training alone on small maps. Both of two methods greatly improves learning speed compared with the previous VI-based model.
- Experiments on 2D, 3D maps and hexapod robot demonstrate that CIN outperforms previous models on accuracy and learning speed, and can be utilized for path planning with real robots.

II. RELATED WORK

A. Path Planning

Path planning is one of the essential tasks in the automation process of a system that moves in the environment while avoiding obstacles and respecting various constraints [18]. It has been widely applied to lots of scenarios including auto driving, autonomous underwater vehicle control and video games [19]–[22]. It has been widely studied by robot researchers, and plenty of methods have been proposed.

A^* algorithm developed by Hart et al. [3] is widely used in path planning because of the following properties: (1) A^* gives optimality when applied with the visible graph. (2) the path given by A^* is unique (3) computation cost is small with a good heuristic function. Many various variants of A^* are proposed such as D^* Lite [23] and any-angle A^* [21]. Artificial potential field (APF) introduced by Oussama Khatib [4] is another important method which creates an artificial potential field to attract the agent around the goal and repulse them around the obstacles [24]. APF is easy to be implemented with smooth paths in real-time. However, the agent may stay in local minimum regions, which is the main drawback of APF method. Probabilistic path planning algorithms such as Rapidly-exploring random trees [5] and probabilistic roadmaps [25] are also effective methods, which randomly select non-collision points in motion space and then connect them to find the best path. These algorithms don't require any environment modeling, which outperform previous algorithms such as A^* in terms of computation cost [22]. Genetic algorithm is another effective path planning problem with high robustness in various scenarios such as robot manipulators and unmanned surface vehicle [26], [27].

Recently, with the development of reinforcement learning and deep learning, deep reinforcement learning (DRL) has been applied to solve path planning problems in many scenarios such as robotic manipulator [7], [8], drilling robot [9] and planetary rover [10]. DRL based methods try to summarize patterns through numerous attempts to generate a suitable path in a new environment, which doesn't require modeling of environment and robot before planning paths. However, current DRL methods rely heavily on training data, which limits the application scenarios.

B. Reinforcement Learning

Reinforcement learning (RL) is a powerful paradigm to solve the sequential decision making problem for Markov Decision Process (MDP), whose objective is to find an optimal policy for the agent through interaction with the environment and maximizing the rewards. RL methods can be categorized into model-free and model-based methods depend on whether the policy has access to the underlying model of the environment.

Model-free methods learn the policy directly from interactions with the environment [12]. For example, Mnih et al. proposed Deep Q network [6] which approximates Q function with a neural network to find the best policy; Lillicrap et al. proposed Deep Deterministic Policy Gradient [28] which utilizes policy networks to generate actions while using Q-networks to criticize the policy. Model-free methods have achieved great success in many application scenarios such as video games [29]

and robot control [30]. However, they are suffering from low sample efficiency for training data and random seed sensitivity, which limits the real application of model free methods.

Model-based RL methods can simulate the transition probability of the environment, resulting in increased sample efficiency, which is particularly important in domains where each interaction with the environment is expensive [12]. For example, PILCO [31] and its variants utilize Gaussian process to learn the environment dynamics and achieve good performance in nonlinear control problems [32]. World Model [33] utilizes MDN-RNN to simulate the environment model with raw images as input and outperforms previous methods in gym tasks. However, learning the environment model may introduces extra complexities and the errors in the environment model may in turn influence the policy accuracy. The state-conditioned transition probabilities utilized in CIN is similar to the environment dynamics in model based RL. In this paper, we combine the idea of the environment dynamics with the Value Iteration based deep learning methods to improve the performance in the path planning tasks.

C. Value Iteration Based Methods

For a markov decision process (MDP) problem with known transition probability $P_{ss'}^a$, a dynamic programming method called Value Iteration [14] can give an optimal policy π^* . The following equations (1) and (2) are the basis of VI algorithm called the Bellman equations:

$$Q(s, a) = R(s, a) + \gamma \sum_{s'} P_{s, s'}^a V(s') \quad (1)$$

$$V'(s) = \max_a Q(s, a) \quad (2)$$

The optimal value map $V^*(s)$ and Q value map $Q^*(s, a)$ can be obtained by calculating between $V(s)$ and $Q(s)$ iteratively until convergence, after which the optimal policy π^* can be obtained utilizing equation (3).

$$\pi^* = \arg \max_a Q^*(s, a) \quad (3)$$

As is shown in Fig. 2, in order to implement planning computation in existing RL models, Tamar et al. [13] proposed a VI-based RL model called *Value Iteration Network* (VIN). VIN combines value iteration with convolution and max-pooling operators to plan paths with unknown environment dynamics. Niu et al. proposed generalized value iteration network (GVIN) [34] which is able to learn and plan paths on irregular spatial graphs using a novel graph convolution operator. Lee et al. proposed Gated Path Planning Network(GPPN) [35] which replaces the unconventional recurrent update

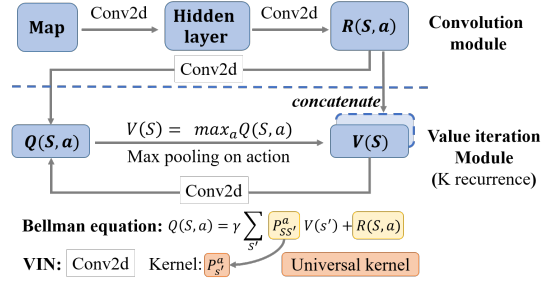


Fig. 2: Structure of the *Value Iteration Network* (VIN) [13] VIN combines the VI algorithm with CNN to implement a RL model with planning capability.

in VIN with a gated LSTM recurrent operator to alleviate optimization issues like random seed sensitivity. Experiment results on 2D and 3D path planning demonstrate that GPPN outperforms VIN on accuracy, learning speed, training instability and sensitivity to random seeds. Pflueger et al. proposed Rover-IRL [10] based on VIN and inverse reinforcement learning(IRL) and achieves good performance on planetary rover path planning problems.

Currently, existing Value Iteration based methods implement planning computation without access to the environment dynamics $P_{ss'}^a$, which is not consistent with the Bellman equations. In this paper, a new Value Iteration based path planning model called CIN is proposed, which utilizes $P_{ss'}^a$ and the sparse reward to implement VI algorithm with convolution neural networks.

III. CAPABILITY ITERATION NETWORK

A. Framework

The framework of CIN is illustrated in Fig. 3 which consists of two important modules: the capability module and value iteration module. The input of CIN is a grid map $\mathcal{M}(m \times m)$, current state s_t and the goal s_g . A sparse reward map $R(s)$ is generated with positive reward r_p at goal state and small negative living cost r_n otherwise, i.e.

$$R(s) = \begin{cases} r_p, & s = s_g, \\ r_n, & s \neq s_g. \end{cases}$$

The core component of CIN is a capability module where environment states are processed as state-conditioned transition probabilities utilizing a neural network classifier. The value iteration module computes the Q-value map $Q(s, a)$ and value map $V(s)$ iteratively utilizing convolution layers with kernels given by the capability module and max-pooling layers respectively. The optimal Q-value map $Q^*(s, a)$ obtained after K iterations is utilized to find the best action a^* for current state s_t by maximization through the action space.

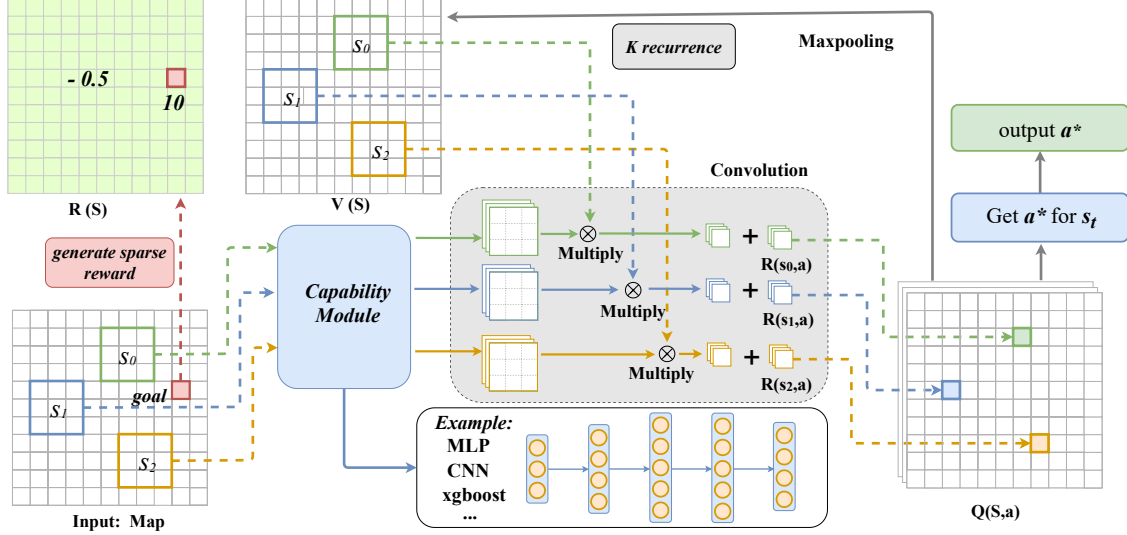


Fig. 3: Structure of *Capability Iteration Network* (CIN). The input of CIN is the original map, current state s_t and the goal position, which is used to generate a sparse reward map firstly. For each state s_i , the algorithm extracts its neighborhood and generates the state-conditioned transition probability $P_{ss'}^a$ utilizing the capability module. Then the algorithm computes the Q-value map $Q(s, a)$ and the value map $V(s)$ utilizing convolution and max-pooling layers iteratively combined with the Bellman equations. After K times iteration, the best action a^* is obtained by max-pooling on the Q-value map: $a^* = \arg \max_a Q(s_t, a)$.

B. Capability Module

Capability module is the key component of CIN, which utilizes neural network classifiers and environment information to predict probability distribution of the next state s_{t+1} , i.e. $P(s_{t+1}|s_t, a)$, often denoted as $P_{ss'}^a$. Actually, it is similar to the idea of the environment dynamics utilized in the model based RL methods.

The input of the capability module is an $F \times F$ local map, which contains information of state $s_{i,j}$ at the center and states can be reached in one step from $s_{i,j}$. F is the kernel size and a hyper-parameter for CIN, which is related to the action capability of the agent. In maze environment, $F = 3$ is utilized. The output is the probability distribution of the next state from the $s_{i,j}$ with all the available actions, i.e.

$$f(\mathcal{M}_{[i,j,F]})_{k,l,a} = P(s_{i+k,j+l}|s_{i,j},a), \\ -\frac{2F-1}{2} \leq k, l \leq \frac{2F-1}{2}, a \in A$$

where $f(\cdot)$ is the capability module, $\mathcal{M}_{[i,j,F]}$ denotes the image patch centered at position (i, j) with kernel size F . Capability module learns and encodes the capability of the agent in the learned neural network, which is the only component needed to train in CIN and determines the performance of the whole model.

For path planning tasks like 2D and 3D navigation, a multilayer perceptron (MLP) classifier is enough; more complex classifiers like CNN [36] and xgboost [37] are necessary for more difficult tasks. More details of

the capability module are provided in the experiment section.

C. Value Iteration Module

The value iteration module utilizes convolution and max-pooling operators to simulate the process of Value Iteration algorithm, which computes the Q-value map $Q(s, a)$ and value map $V(s)$ iteratively. The following equations are the Bellman equations, the basis of VI algorithm:

$$Q(s, a) = \gamma \sum_{s'} P_{ss'}^a V(s') + R(s, a) \quad (4)$$

$$V(s) = \max_a Q(s, a) \quad (5)$$

As is shown in Fig. 3, CIN utilizes convolution to simulate the equation (4). The transition probabilities $P_{ss'}^a$ provided by the capability module are the convolution kernels, and the sparse reward $R(s)$ are the bias, i.e.

$$Q^{(k+1)}(s_{i,j}, a) = \gamma f(\mathcal{M}_{[i,j,F]})_{i,j} \cdot V^{(k)}(s_{[i,j,F]}) + R(s_{i,j}), \quad 1 \leq i, j \leq m$$

Then the new value map $V(s)$ is obtained utilizing max-pooling illustrated in equation (5) and the Q-value map $Q(s, a)$ given by the previous convolution operator:

$$V^{(k+1)}(s_{i,j}) = \max_{a \in A} Q^{(k)}(s_{i,j}, a), \\ 1 \leq i, j \leq m$$

The VI algorithm tell us that the value iteration module will output the optimal Q-value map $Q^*(s, a)$ after convergence:

$$\lim_{k \rightarrow \infty} Q^{(k)}(s, a) = Q^*(s, a)$$

In practical, the number of iteration K is a hyper-parameter, which needs to be large enough to ensure convergence. Then the best action a^* for current state s_t is obtained utilizing maximization on the Q value map through the action space:

$$a^* = \arg \max_{a \in A} Q^*(s_t, a)$$

Compared to VIN model described in Fig. 2, CIN utilizes state-conditioned transition probabilities $P_{ss'}^a$ in the value iteration module instead of P_s^a in VIN, which is completely consistent with the Bellman equation. This makes what CIN learns closer to the nature of the task compared with VIN. The sparse reward map used in CIN consists with positive reward at the goal and small living cost otherwise, which is simple to generate and avoids the accuracy of CIN being influenced by the defects in reward maps generated by the neural network, especially on large maps.

D. Training Methods

There is no parameter to train in the value iteration module and the capability module is just a simple classifier to predict the probability distribution of the next state in MDP. Thus, we can train capability module alone by supervised learning on local maps. Besides, training CIN end-to-end is also a good choice.

End-to-end: The whole network of CIN is differentiable, which means we can train CIN end-to-end using any RL or Imitation learning (IL) [38] algorithms. The Cross Entropy Loss is utilized for IL in our experiment:

$$\mathcal{L}_{CE} = - \sum_{i=1}^{|A|} p_i \log \hat{p}_i, \quad p_i = \begin{cases} 1, & a_i = a^* \\ 0, & a_i \neq a^* \end{cases}$$

where \hat{p}_i is the probability of the model to choose the action a_i .

CIN utilizes sparse reward maps immediately, rather than complex reward maps generated by CNN in previous models. This makes CIN need less computation cost and has faster learning speed compared to VIN, especially on large maps. The experiment results show that CIN can achieve better accuracy and faster learning speed compared to existing VI based models and reactive policies. More details are described in the experiment section.

Train Capability Module Alone: We can also train capability module alone by supervised learning. Take 2D grid world task as an example: the agent is controlled with a random policy to obtain random trajectories. For

each state s_t on the trajectories, 3×3 small local map with center s_t is sampled, which is utilized as the training data combined with the action a_t taken by the agent. The corresponding position of the next state s_{t+1} is the training label. The capability module is seen as a classifier and trained utilizing supervised learning with the Mean Squared Error(MSE) loss:

$$\mathcal{L}_{MSE} = \frac{1}{F^2} \sum_{k=1}^F \sum_{l=1}^F (f(\mathcal{M}_{[i,j,F]})_{k,l,a} - P_{i,j,k,l})^2$$

, where $P_{i,j,k,l} = P(s_{i+k-\frac{2F-1}{2}}, s_{j+l-\frac{2F-1}{2}} | s_{i,j}, a)$.

To speed up convergence of the capability module, curriculum learning [39] is needed. Take 3D terrain map task as an example, our training goal is to teach capability module the capability of the agent to overcome obstacles (i.e., the max height difference Δh^* the agent can overcome). Thus, we start training capability module by “simple” local maps with small or big height difference. Gradually, “harder” local maps with height difference approach to Δh^* are needed for training. Training with curriculum learning can speed up learning and achieves better accuracy.

IV. EXPERIMENT AND DISCUSSION

In this section, we evaluate the performance of CIN in three types of experiment environments(2D, 3D grid world and real hexapod robot) and make a comparison with existing methods including VIN, GPPN and ResNet-based reactive policy to demonstrate the effectiveness and improvement of CIN.

A. 2D Grid-World Domain

Grid world path planning with obstacles randomly located in the map is one of the most common experiment environments for path planning problem. In 2D grid-world domain, the input map contains binary numbers 0 (obstacles) and 1 (free-space) only. The 2D grid world dataset is created with a maze generation algorithm which uses Depth-First Search and the Recursive Backtracker algorithm [40] with random starting and ending points. The size of the training dataset $N_{tr} = 10K$; validation dataset size $N_{val} = 1K$; test dataset size $N_{te} = 1K$.

CIN: In this experiment, we train the capability module alone with supervised learning for simplicity. The sparse reward map consists of a high reward +10 at the goal state and living cost -0.5 at other states. A multilayer perceptron(MLP) consisting of four letent layers, one softmax layer and ReLU activation function is utilized for the capability module in this task. We utilize Mean Square Error (MSE) loss with Adam optimizer [41] for optimization.

Two baselines are utilized in our experiments: **VIN** [13]: VIN is a classic VI-based model combined with VI

m	VIN		GPPN		CIN		ResNet18	
	%Opt	%Suc	%Opt	%Suc	%Opt	%Suc	%Opt	%Suc
8	98.9	99.3	98.6	99.2	99.9	99.9	99.7	99.9
15	89.1	90.5	97.6	98.4	99.7	99.8	97.4	98.2
28	80.9	81.0	96.3	96.7	99.5	99.5	86.5	91.8
45	63.3	72.1	92.3	95.2	99.0	99.1	54.0	66.3

TABLE I: Experiment result in 2D grid world with varying map size ($m \times m$).

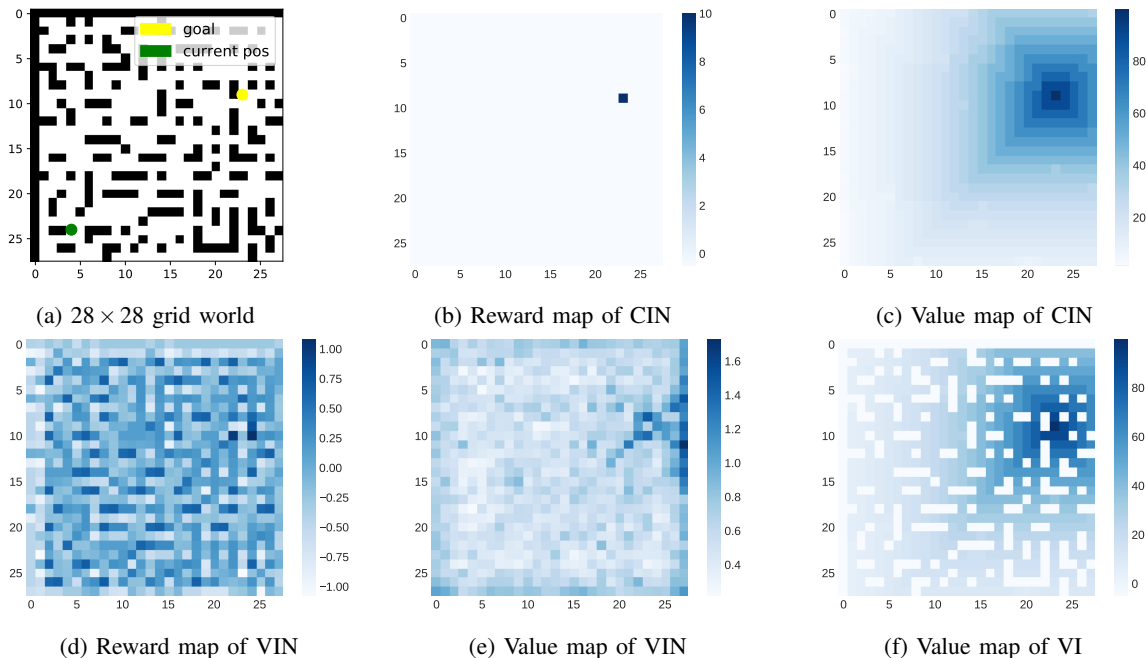


Fig. 4: Visualization of the 2D maze, reward maps and value maps.

algorithm and CNN. VIN achieves good performance in path planning problems because of its planning computation. **ResNet18** [42]: The problem setting for this task is similar to the image segmentation, in which problem each pixel in the given image needs to be assigned a label (the best action choice in our case). Thus, image segmentation model ResNet18 can be seen as a reactive path planning model and utilized as a baseline in our experiments.

We conduct experiments on maps with four different sizes: 8×8 , 15×15 , 28×28 and 45×45 . Two metrics are utilized to evaluate the performance of algorithms, including **%Optimal(%Opt)**—the percentage of states whose predicted paths under the policy estimated has optimal length, **%Success(%Suc)**— the percentage of states whose predicted paths under the policy estimated reach the goal state.

Table I shows the performance of models on each metric. Two conclusions can be concluded: (1) CIN outperforms other models in 2D domain on two metrics, especially on maps with large sizes. (2) Models with planning computation has better performance on big

maps, compared to reactive policies (ResNet18 in this experiment). The conclusion (2) has been proposed in [13].

B. Model Result Visualization

The visualization of the model results are depicted in Fig. 4. The reward map used in CIN (Fig. 4b) is sparse, which is consistent with the value iteration algorithm. Reward map of VIN (Fig. 4d) is generated by CNN, which is more messy, sensitive to the map settings and may influence the accuracy of VIN on unseen or large mazes. Compared to VIN (Fig. 4e), value map of CIN (Fig. 4c) is sharp at the edges of obstacles with obvious tendency towards goal point, and is generally consistent with optimal value map (Fig. 4f) generated by the value iteration algorithm.

C. Learning Speed and Stability

The learning curve of each model on the 15×15 2D mazes is depicted in Fig. 5. In order to compare the learning speed of each model fairly, CIN is trained end-to-end utilizing Imitation Learning(IL) in this experiment. All the models are trained utilizing at least 5

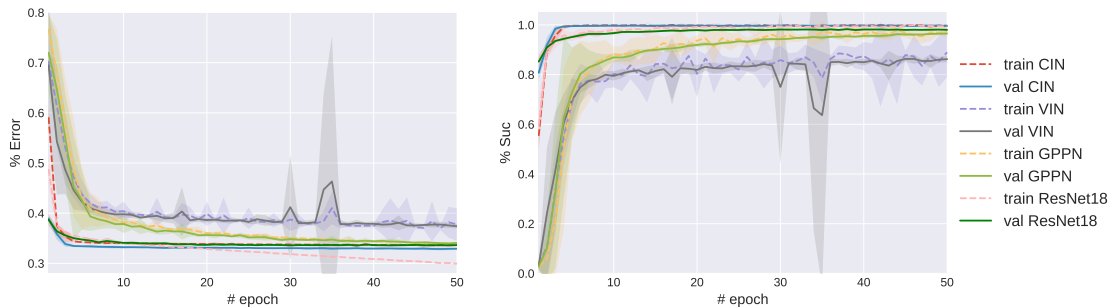


Fig. 5: The learning curve of each model on the 15×15 2D mazes with different random seeds. All the models are trained end-to-end using imitation learning for fairness. CIN outperforms other models on learning speed, training stability and sensitivity to random seeds.

different random seeds with training dataset size $N_{tr} = 10K$ in 50 epochs.

The %Error in Fig. 5 is the percentage of states whose predicted action by the model estimated is different with the action given by the expert. As shown in the figure, CIN and ResNet18 have faster learning speed compared to GPPN and VIN, because CIN utilizes sparse reward immediately while VIN and GPPN uses CNN to generate reward maps with more parameters to train. However, ResNet18 appears overfitted because it belongs to reactive policy without planning computation. Besides, GPPN utilizes gated LSTM recurrent operator which improves the training stability and random seed sensitivity compared to VIN. As shown in the figure, CIN has better random seed sensitivity and converges more stably than other models.

D. 3D Terrain Map

Unlike maps in 2D domain only contain 0 and 1, matrix elements in 3D terrain maps can be any float value, representing height in a specific position. The agent can only walk to its neighborhoods if the height difference is smaller than a certain value. Other experiment settings are same to the 2D grid world experiment. Compared to 2D mazes, 3D terrain maps are more complex and challenging for path planning models.

Two metrics used in this experiment are same to which utilized in the 2D environment. CIN is trained by training capability module alone while other models are trained end-to-end with imitation learning. As shown in the Table. II, CIN still achieves better performance than VIN and GPPN on various metrics and map settings, especially on big maps.

E. Real Hexapod Robot

Hexapod robot. The hexapod robot has the characteristics of good stability, large carrying capability and flexibility, and can adapt to complex environments and

special terrains. The robot used in this experiment is a parallel structure robot, also known as Parallel-Parallel (PP) structure six-legged robot. Each leg of the robot consists of three connecting rods $[l_1, l_2, l_3]$. The lifting, lowering and stepping of the legs are realized by the expansion and contraction of the three rods, and each leg can be taken in multiple directions.

Adaption to hexapod robot. Currently, VI-based algorithms are only effective for discrete state and action space. Thus, simplification and discretization are needed to be utilized. Given a real environment, we construct a 3D map through point cloud data of Lidar with SLAM algorithm. Then the 2D map with small squares of $0.2m \times 0.2m$ is obtained by projecting the 3D map to 2D planes. Each small square is given a constant float number as the average height. The 2D coordinate of the square which the robot mass center lies in is utilized to represent the current state of the robot.

To simplify the large continuous action space of the hexapod robot, we design a triangle-gait for the robot. In triangle-gait, six legs are divided into two group with three non-adjacent feet as a group. The robot take three feet first and then the other ones for a walk cycle. The gait is defined as (step length, step height, rising length) with eight walking directions. Available combinations are $(0.1m, 0.3m, 0.3m)$, $(0.2m, 0.2m, 0.2m)$, $(0.3m, 0.1m, 0.1m)$. For simplicity, We train capability module alone to reduce the time cost.

We simulate the robot in CoppeliaSim (V-rep) which is a versatile, scalable and powerful general-purpose robot simulation framework [43]. As is shown in Fig. 6, the robot can step up stairs, make U-turn on the corner and handle uneven surface, utilizing the path planned by CIN. Experiment in the real environment is also conducted. As is shown in Fig. 7, we apply CIN trained in the simulation to the real hexapod robot which also achieves good performance. More details are provided in the attached video.

m	VIN		GPPN		CIN		ResNet18	
	%Opt	%Suc	%Opt	%Suc	%Opt	%Suc	%Opt	%Suc
8	92.0	93.7	93.6	94.0	95.2	96.5	94.0	94.2
15	83.7	86.8	89.3	89.7	93.7	94.3	79.1	81.3
28	74.1	76.5	85.3	87.5	89.6	89.9	56.2	60.0
45	60.9	63.1	80.3	81.5	85.3	87.6	47.3	51.3

TABLE II: Experiment result in 3D grid world with varying map size ($m \times m$).

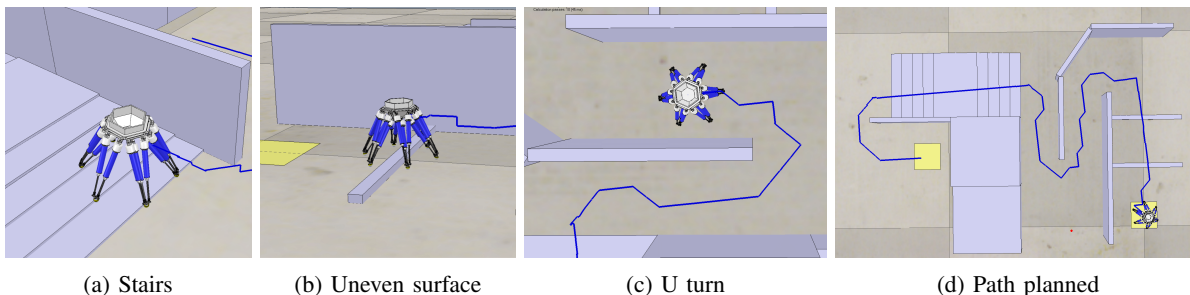


Fig. 6: Simulation in CoppeliaSim (V-rep) with a hexapod robot. CIN achieves good performance for the path planning on stairs, uneven terrain and U turn.

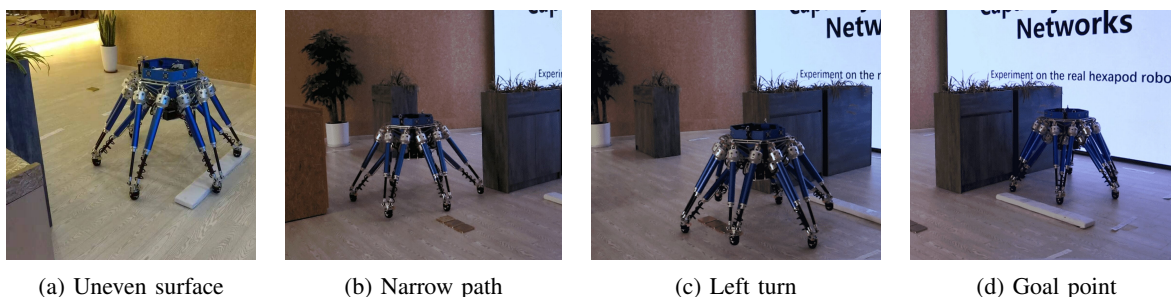


Fig. 7: Experiment on the real hexapod robot. The same hexapod robot and environment adaptation methods are utilized as in the simulation. CIN trained in the simulation still achieves good performance in the real world. More information is provided in the attached video.

V. CONCLUSION

In this paper, we propose a new VI-based path planning model, *Capability Iteration Network*, which improves accuracy and learning speed compared with existing VI-based algorithms, especially on large-scale maps. CIN utilizes sparse reward maps and state-conditioned transition probabilities $P_{ss'}^a$ instead of reward maps generated by CNN and global convolution kernel P_s^a in the previous algorithms. This makes CIN consistent with VI algorithm, only need to learn the capability of the agent with higher accuracy and faster learning speed. Several experiments are conducted, including path planning on 2D, 3D grid world and a real hexapod robot. The results demonstrate that CIN outperforms existing algorithms and is able to be utilized in path planning tasks for real robots.

REFERENCES

- [1] T. T. Mac, C. Copot, D. T. Tran, and R. De Keyser, "Heuristic approaches in robot path planning: A survey," *Robotics and Autonomous Systems*, vol. 86, pp. 13–28, 2016.
- [2] E. Plaku, L. E. Kavraki, and M. Y. Vardi, "Motion planning with dynamics by a synergistic combination of layers of planning," *IEEE Transactions on Robotics*, vol. 26, no. 3, pp. 469–482, 2010.
- [3] P. E. Hart, N. J. Nilsson, and B. Raphael, "A formal basis for the heuristic determination of minimum cost paths," *IEEE transactions on Systems Science and Cybernetics*, vol. 4, no. 2, pp. 100–107, 1968.
- [4] O. Khatib, "Real-time obstacle avoidance for manipulators and mobile robots," in *Autonomous robot vehicles*. Springer, 1986, pp. 396–404.
- [5] S. M. LaValle, "Rapidly-exploring random trees: A new tool for path planning," 1998.
- [6] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, p. 529, 2015.
- [7] M. Sadeghzadeh, D. Calvert, and H. A. Abdullah, "Autonomous visual servoing of a robot manipulator using reinforcement

- learning,” *International Journal of Robotics and Automation*, vol. 31, no. 1, 2016.
- [8] T. Yan, W. Zhang, S. X. Yang, and L. Yu, “Soft actor-critic reinforcement learning for robotic manipulator with hindsight experience replay,” *International Journal of Robotics & Automation*, vol. 34, no. 5, pp. 536–543, 2019.
- [9] Y. Liu, M. Cong, H. Dong, and D. Liu, “Reinforcement learning and ega-based trajectory planning for dual robots,” *International Journal of Robotics & Automation*, vol. 33, no. 4, pp. 367–378, 2018.
- [10] M. Pflueger, A. Agha, and G. S. Sukhatme, “Rover-irl: Inverse reinforcement learning with soft value iteration networks for planetary rover path planning,” *IEEE Robotics and Automation Letters*, vol. 4, no. 2, pp. 1387–1394, 2019.
- [11] A. Faust, K. Oslund, O. Ramirez, A. Francis, L. Tapia, M. Fiser, and J. Davidson, “Prm-rl: Long-range robotic navigation tasks by combining reinforcement learning and sampling-based planning,” in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 5113–5120.
- [12] K. Arulkumar, M. P. Deisenroth, M. Brundage, and A. A. Bharath, “A brief survey of deep reinforcement learning,” *arXiv preprint arXiv:1708.05866*, 2017.
- [13] A. Tamar, Y. Wu, G. Thomas, S. Levine, and P. Abbeel, “Value iteration networks,” in *Advances in Neural Information Processing Systems*, 2016, pp. 2154–2162.
- [14] R. Bellman, “Dynamic programming,” *Science*, vol. 153, no. 3731, pp. 34–37, 1966.
- [15] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin, O. P. Abbeel, and W. Zaremba, “Hindsight experience replay,” in *Advances in neural information processing systems*, 2017, pp. 5048–5058.
- [16] M. Vecerik, T. Hester, J. Scholz, F. Wang, O. Pietquin, B. Piot, N. Heess, T. Rothörl, T. Lampe, and M. Riedmiller, “Leveraging demonstrations for deep reinforcement learning on robotics problems with sparse rewards,” *arXiv preprint arXiv:1707.08817*, 2017.
- [17] M. Riedmiller, R. Hafner, T. Lampe, M. Neunert, J. Degraeve, T. Wiele, V. Mnih, N. Heess, and J. T. Springenberg, “Learning by playing solving sparse reward tasks from scratch,” in *International Conference on Machine Learning*, 2018, pp. 4344–4353.
- [18] O. Souissi, R. Benatallah, D. Duvivier, A. Artiba, N. Belanger, and P. Feyszeau, “Path planning: A 2013 survey,” in *Proceedings of 2013 International Conference on Industrial Engineering and Systems Management (IESM)*. IEEE, 2013, pp. 1–8.
- [19] B. Paden, M. Cáp, S. Z. Yong, D. Yershov, and E. Frazzoli, “A survey of motion planning and control techniques for self-driving urban vehicles,” *IEEE Transactions on intelligent vehicles*, vol. 1, no. 1, pp. 33–55, 2016.
- [20] G. Che, L. Liu, and Z. Yu, “An improved ant colony optimization algorithm based on particle swarm optimization algorithm for path planning of autonomous underwater vehicle,” *Journal of Ambient Intelligence and Humanized Computing*, vol. 11, no. 8, pp. 3349–3354, 2020.
- [21] P. K. Y. Yap, N. Burch, R. C. Holte, and J. Schaeffer, “Any-angle path planning for computer games,” in *Seventh Artificial Intelligence and Interactive Digital Entertainment Conference*, 2011.
- [22] W. Gong, X. Xie, and Y.-J. Liu, “Human experience-inspired path planning for robots,” *International Journal of Advanced Robotic Systems*, vol. 15, no. 1, p. 1729881418757046, 2018.
- [23] S. Koenig and M. Likhachev, “D* lite,” in *Eighteenth national conference on Artificial intelligence*, 2002, pp. 476–483.
- [24] J. Wang, X.-j. Lin, H.-y. Zhang, G.-d. Lu, Q.-l. Pan, and H. Li, “Path planning of manipulator using potential field combined with sphere tree model,” *International Journal of Robotics and Automation*, vol. 35, no. 2, 2020.
- [25] L. E. Kavraki, P. Svestka, J.-C. Latombe, and M. H. Overmars, “Probabilistic roadmaps for path planning in high-dimensional configuration spaces,” *IEEE transactions on Robotics and Automation*, vol. 12, no. 4, pp. 566–580, 1996.
- [26] Y. Liu, M. Cong, H. Dong, D. Liu, and Y. Du, “Time-optimal motion planning for robot manipulators based on elitist genetic algorithm,” *International Journal of Robotics & Automation*, vol. 32, no. 4, pp. 396–405, 2017.
- [27] J. Xin, J. Zhong, J. Sheng, P. Li, and Y. Cui, “Improved genetic algorithms based on data-driven operators for path planning of unmanned surface vehicle,” *International Journal of Robotics and Automation*, vol. 34, no. 6, 2019.
- [28] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, “Continuous control with deep reinforcement learning,” *arXiv preprint arXiv:1509.02971*, 2015.
- [29] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, “Playing atari with deep reinforcement learning,” *arXiv preprint arXiv:1312.5602*, 2013.
- [30] Y. Duan, X. Chen, R. Houthoofd, J. Schulman, and P. Abbeel, “Benchmarking deep reinforcement learning for continuous control,” in *International Conference on Machine Learning*, 2016, pp. 1329–1338.
- [31] M. Deisenroth and C. E. Rasmussen, “Pilco: A model-based and data-efficient approach to policy search,” in *Proceedings of the 28th International Conference on machine learning (ICML-11)*, 2011, pp. 465–472.
- [32] R. McAllister and C. E. Rasmussen, “Data-efficient reinforcement learning in continuous state-action gaussian-pomdps,” in *Advances in Neural Information Processing Systems*, 2017, pp. 2040–2049.
- [33] D. Ha and J. Schmidhuber, “Recurrent world models facilitate policy evolution,” in *Advances in Neural Information Processing Systems*, 2018, pp. 2450–2462.
- [34] S. Niu, S. Chen, H. Guo, C. Targonski, M. C. Smith, and J. Kovačević, “Generalized value iteration networks: Life beyond lattices,” in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [35] L. Lee, E. Parisotto, D. S. Chaplot, E. Xing, and R. Salakhutdinov, “Gated path planning networks,” in *International Conference on Machine Learning*, 2018, pp. 2947–2955.
- [36] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, “Backpropagation applied to handwritten zip code recognition,” *Neural computation*, vol. 1, no. 4, pp. 541–551, 1989.
- [37] T. Chen and C. Guestrin, “Xgboost: A scalable tree boosting system,” in *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, 2016, pp. 785–794.
- [38] P. Abbeel and A. Y. Ng, “Apprenticeship learning via inverse reinforcement learning,” in *Proceedings of the twenty-first international conference on Machine learning*, 2004, p. 1.
- [39] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, “Curriculum learning,” in *Proceedings of the 26th annual international conference on machine learning*. ACM, 2009, pp. 41–48.
- [40] Wikipedia, “Maze generation algorithm — Wikipedia, the free encyclopedia,” 2020, [Online; accessed 01-September-2020].
- [41] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [42] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [43] E. Rohmer, S. P. Singh, and M. Freese, “V-rep: A versatile and scalable robot simulation framework,” in *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2013, pp. 1321–1326.