

Article

Not peer-reviewed version

Path Planning for Robot Combined with Zero-Shot and Hierarchical Reinforcement Learning in Inexperienced Environments

[Liwei Mei](#) and [Pengjie Xu](#)*

Posted Date: 30 October 2024

doi: 10.20944/preprints202410.2427.v1

Keywords: path planning; zero-shot learning; hierarchical reinforcement learning; adaptive agents



Preprints.org is a free multidiscipline platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Article

Path Planning for Robot Combined with Zero-Shot and Hierarchical Reinforcement Learning in Inexperienced Environments

Mei Liwei ^{1,†} and Xu Pengjie ^{2,*,†}

¹ School of Information Science and Technology, East China University of Science and Technology

² School of Mechanical Engineering, Shanghai Jiaotong University

* Correspondence: xupengjie194105@sjtu.edu.cn

† These authors contributed equally to this work.

Abstract: Path planning for robots based on reinforcement learning encounters challenges in integrating semantic information about environments into the training process. In those unseen or complex environmental information, agents often perform sub-optimally and require more training time. In response to these challenges, this manuscript pioneers a framework integrating zero-shot learning combined with hierarchical reinforcement learning to enhance agent decision-making in complex environments. Zero-shot learning enables agents to infer correct actions for previously unseen objects or situations based on learned semantic associations. Subsequently, the path planning component utilizes hierarchical reinforcement learning with adaptive replay buffer, directed by the insights gained from zero-shot learning, to make decisions effectively. Two parts are trained separately, so zero-shot learning is available in different and unseen environments. Through simulation experiments, the proposed method proves that this structure can make full use of environmental information to generalize across unseen environments and plan collision-free paths.

Keywords: path planning; zero-shot learning; hierarchical reinforcement learning; adaptive agents

1. Introduction

Robotic path planning involves algorithms that instruct a robot to take reasonable steps to approach a user-specified location in an unknown environment. This task is essential for ensuring effective navigation and decision-making capabilities of robots and unmanned vehicles in the evolving landscape of autonomous systems. The ability to navigate through complex and dynamic environments is crucial for applications ranging from autonomous driving to warehouse automation and search and rescue operations.

Moreover, traditional path-planning algorithms, including A* and Dijkstra's algorithm [1,2], although effective in static environments, often struggle in some environments where full information is not available. For instance, in warehouse automation, the A* algorithm can be used to find the shortest path for a robot to navigate from a storage location to a packing area. However, if an unexpected obstacle, such as a misplaced package, blocks the planned route, the algorithm must re-plan the path. This re-planning can introduce significant delays, reducing operational efficiency. Additionally, A* can sometimes find paths that are theoretically optimal but practically infeasible due to narrow corridors or tight turns that the robot cannot navigate, necessitating manual intervention. Apart from this, heuristic methods like Particle Swarm Optimization (PSO) [3,4], though shows capacity to dynamically adjust to environmental changes without the need for complete map information, also present challenges. For example, in agricultural robotics, PSO has been applied to navigate drones for crop monitoring. However, the algorithm can suffer from premature convergence, where particles get trapped in local optima, leading to suboptimal paths. These algorithms generally require a pre-built map which is based on simultaneous localization and mapping (SLAM) [5], and do not adapt well to changes, necessitating frequent re-planning, which can be computationally expensive and inefficient. The accuracy of SLAM heavily depends on the precision of the sensors used, which can be compromised by factors such as sensor noise, range limitations, and resolution constraints.

Additionally, the high cost and instability of the equipment under various weather conditions are significant constraints that limit the practical deployment of these traditional methods in real-world scenarios [6].

Because of the mechanism of RL, a robot is able to make decisions by collecting interactions with the environment and then choosing an optimal policy by maximizing the collected rewards [7]. An existing problem in the reinforcement learning training process is how to provide appropriate rewards. Numerous methods can achieve this functionality, but they come with some compromises in some other performance. For instance, the approach of utilizing frequency and rewards to form gradient information guidance [8] may lead reinforcement learning to fall into local optima and decrease its exploration capabilities. To solve the mentioned challenges, this manuscript explores a method combining zero-shot learning and hierarchical reinforcement learning.

2. Related Work

2.1. Zero-Shot Generalization

Zero-shot learning (ZSL) [9], a class of algorithms capable of performing well on new tasks without additional data collection, holds immense potential across various domains. In this manuscript, two concepts are mainly discussed, ZSL in image processing and zero-shot generalization(ZSG) in path planning.

ZSL in image processing: It enables models to recognize and classify previously unseen objects by leveraging semantic information encoded in class attribute descriptions and embeddings [10]. By bridging the gap through semantic knowledge transfer, zero-shot learning generalizes beyond the confines of training data to novel instances. The Contrastive Language-Image Pre-Trained (CLIP) model is a masterpiece which has been implemented in multiple areas. For instance, in feature extraction, Zhang et al. demonstrated that ZSL could be enhanced by training multi-modal embeddings using a deep learning model [11]. In environment detection, it combines 3D point clouds and RGB images to enhance the capability to identify 3D objects. Zhu et al. put forth PointCLIP V2 [12]. It projects 3D cloud points to a 2D plane, gets the embedding by Bidirectional Encoder Representations from Transformers, and then compares the similarity of embedding of projecting the image and semantic information to form an identifying result. Additionally, in the multimedia area, Song et al. introduced MeshCLIP [13], which processes cross-modal information for 3D mesh data using zero-shot learning, thereby improving reconstruction quality. In this manuscript, CLIP is used to generate action guidance according to visual information.

Zero-shot generalization: It is an essential metric that evaluates the capability of a model to perform effectively on new, unseen tasks without additional training [14]. As shown in Figure 1, the data distribution of this problem can be divided into three situations e.g. the training set and test set are the same, they follow the same distribution, or they follow different distributions. This capability is particularly crucial in inexperienced environments where pre-defined training datasets cannot cover all possible scenarios. ZSG enables models to leverage prior knowledge and apply it to novel situations by relying on semantic information and transfer learning mechanisms [15]. In the context of RL, zero-shot generalization is formalized as the ability of a policy, trained in one set of contexts, to perform well in entirely new, unseen contexts. This involves specifying which subset of ZSG problems is being addressed, as it encompasses a range of scenarios rather than a single specific problem. For instance, a robot trained to navigate in a particular type of environment should be able to adapt its navigation strategies to different environments it has never encountered before, relying on the transfer of learned semantic relationships and policy structures.

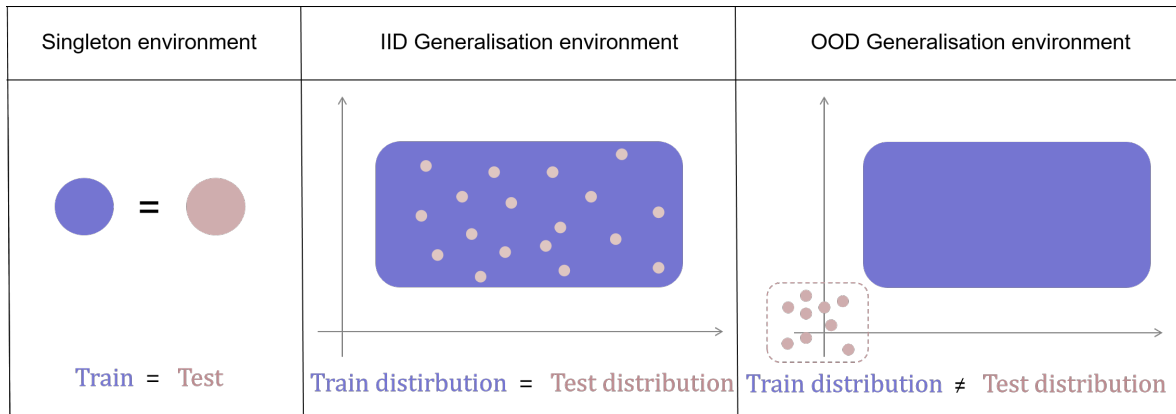


Figure 1. Zero-shot generalization data distribution, the singleton environment means that the train set is identical to the test set. The independent and identical(IID) generalization environment means that the two sets follow the same distribution while the Out-Of-Distribution(OOD) means they follow different distribution.

2.2. Artificial Potential Field

The Artificial Potential Field (APF) method is a widely used approach in robotic path planning [16]. It conceptualizes the environment as a potential field, where the robot is treated as a particle moving under the influence of this field. The method employs two primary components: an attractive potential field that pulls the robot towards the goal and a repulsive potential field that pushes it away from obstacles.

In robotic path planning, the APF method is advantageous due to its simplicity and real-time applicability, especially in partially known environments. Since this method can lead to local minima where the net force on the robot is zero, preventing further movement, we make some modifications to suit the hierarchical reinforcement learning process.

2.3. Hierarchical Reinforcement Learning

HRL draws inspiration from the way humans approach complex problems—by decomposing them into smaller, more manageable tasks [9]. HRL operates on multiple levels of decision-making, where high-level policies guide the overall direction toward the goal, and low-level policies handle the specific action. Such a structure mirrors the cognitive process of setting a general objective and executing detailed actions to achieve it [17]. This approach improves learning efficiency and enhances adaptability in dynamic environments. Recent advancements have seen HRL successfully applied across various domains, including robotic navigation, gaming, and autonomous driving [18,19]. These studies highlight the effectiveness of HRL in managing high-dimensional state spaces and executing complex sequential decisions, showcasing its potential in facilitating long-term planning and precision control. Christen et al. put forth an explicit task decomposition method [20], which can conduct a zero-shot of the planning layer across different low-level agents without retraining. Chen et al. propose a soft actor-critic structure with a prioritized experience replay [21], solving the problem of low sample utilization. Ye et al. propose a hierarchical policy learning with intrinsic-extrinsic modeling [22]. As analyzed above, information about environments, especially unseen ones, can impact the efficiency of an agent. Additionally, the hierarchical structure has shown great potential for the improvement of efficiency and generalization ability.

In the robotic visual servoing path planning task, the perception of vision sensors can be incomplete due to issues such as occlusion. To address this problem, we employ a hierarchical reinforcement learning framework. The high-level decision-making module is guided by visual information, and the visual guidance is constrained within a specific area. This approach ensures that the path planning process can effectively handle challenges arising from incomplete perception, improving the robustness and accuracy of the task.

To enhance agent performances, we explored a path planning method that integrates ZSL with HRL. The contributions of our work are described as follows:

- 1) The proposed method first fuses ZSL into the reinforcement learning process for robot path planning. After integrating the unseen semantic information, the agent becomes more intelligent regarding path selection and training cost control.
- 2) A reasonable fusion architecture is proposed specifically. ZSL is utilized on a high level to infer correct macro-actions for previously unseen objects or situations based on learned semantic relationships. Then, the HRL follows cues offered by zero-shot learning to make decisions effectively in specific path selections.
- 3) Detailed performance analysis is provided for the proposed combined learning framework. The simulation results corroborate the proposed method's capability of leveraging visual cues for decisions and modifying agents' actions.

3. Proposed Method

The system architecture of the proposed method is shown in Figure 2. It is a two-stage learning process with a combination of zero-shot learning and hierarchical reinforcement learning. Part A and Part B belong to the zero-shot learning methods. Part C is the main structure of hierarchical reinforcement learning. The designed visual-act model utilizes a zero-shot learning image encoder to transfer images into embedding. Then, image embeddings and correspondent labels are sent to a multi-head self-attention network. This work is shown in Figure 3. The outputs of this stage are action instructions according to input images. It identifies the focusing parts of image embedding related to the exact actions in each direction. The latter stage is a hierarchical reinforcement learning process. It utilizes a high-level and a low-level policy to form paths. The high-level policy determines the plausible area while the low-level policy provides exact action. In the meantime, images captured by a visual sensor are sent to the pre-trained model. Then, the pre-trained model provides action instructions to correct decisions made by the reinforcement learning part.

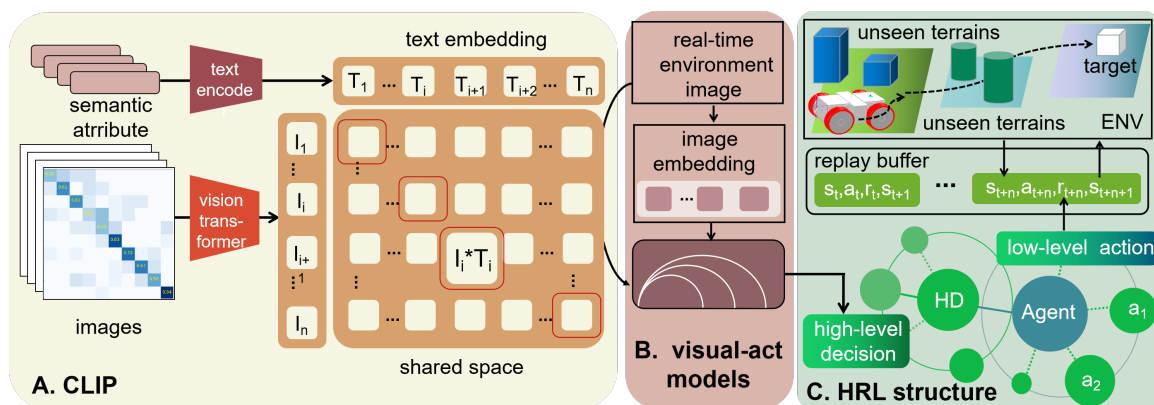


Figure 2. Main structure of the proposed method.

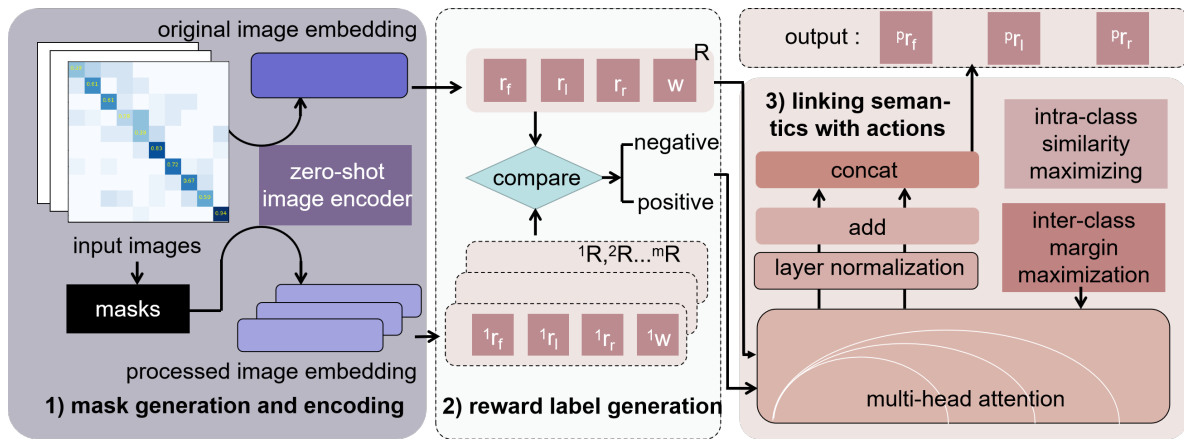


Figure 3. pre-trained visual-act decision model.

3.1. Contrastive Language-Image Pre-Trained Model

The CLIP model leverages the relationships between images and their corresponding text descriptions to achieve zero-shot learning capabilities. In this section, we describe the core components and mechanisms of the CLIP model and how it contributes to our proposed method. As in Figure 1, the method can be described as below:

Images I : train set images with the corresponding text descriptions.

Semantic attributes T : different text descriptions, including those that match the images and those that do not.

Vision transformer $f(I)$: [23] interpret image information into vectors v_I .

Text encoder $g(T)$: [24] transform semantic attributes into text embedding v_T .

Shared space: a mapping relationship formed by contrastive learning.

$$v_I = f(I), \quad v_T = g(T) \quad (1)$$

The similarity between an image and a text description is then quantified by a similarity metric S , typically the cosine similarity, which is computed as:

$$S(v_I, v_T) = \frac{v_I \cdot v_T}{\|v_I\| \|v_T\|} \quad (2)$$

Elements in $S(v_I, v_T)$ are separated into corresponding pairs and non-corresponding ones. During training, the model is optimized to maximize the similarity between corresponding images and texts while minimizing the similarity between non-corresponding pairs [25]. Such a semantic-based zero-shot learning approach allows the model to acquire the correct embedding of images when faced with previously unseen objects by leveraging the relationship between semantic information and image features. Because this encoder is trained using a contrastive learning strategy, the model can only focus on the relative relationships between images and semantic information and cannot attend to absolute information, such as the spatial coordinates of objects. This limitation is addressed by the training approach of the proposed method, which utilizes masks to specify exact objects related to decision-making.

3.2. Visual-Act Model

This part utilizes the zero-shot image encoder to transform images into embedding vectors. Then, the model figures out the relationship between visual information and action selection, which can extract semantic information related to decision-making, thus optimizing the path-planning process.

Algorithm 1 Optimized model for high-level macro-action decision making in action selection

```

1: Input: Image dataset  $\mathcal{D}$  with annotations, Set of masks  $\{M_i\}$ , Number of epochs  $N$ 
2: Output: Optimized model for action decision  $R_p[r_p(f), r_p(l), r_p(r)]$ , processed image embeddings  $I_1, I_2, I_3$ 
3: Initialize the zero-shot image encoder  $Z$ 
4: Initialize action-specific networks for each action  $O_a$ 
5: Initialize parameters for contrastive loss
6: Define total loss function  $L_{total}$  with weighting factor  $\lambda$ 
7: for each epoch do
8:   for each (Image,  $R$ ,  $\{R_i\}$ ) in dataset  $\mathcal{D}$  do
9:      $E \leftarrow Z(\text{Image})$ 
10:    for each mask  $M_i$  in Set of masks do
11:       $E_i \leftarrow Z(\text{Image} \oplus M_i)$ 
12:    end for
13:    Initialize task-specific losses  $L_a$  for this batch
14:    for each action  $a$  in {forward, left, right} do
15:       $R_p, I_1, I_2, I_3 \leftarrow \text{action-specific network}(E, a)$ 
16:    end for
17:     $L_a \leftarrow \|R_p - R\|_2$ 
18:     $L_c \leftarrow \text{contrastive loss}(E, I_1, I_2, I_3, \{E_i\}, \{R_i\})$ 
19:     $L_{total} \leftarrow \sum L_a + \lambda \times L_c$ 
20:    Backpropagate  $L_{total}$  and update parameters
21:  end for
22:  Optionally evaluate the model on the validation set
23: end for
24: Save the optimized model parameters

```

(1) Mask generation and encoding

For each input image I , a series of object masks [26] are applied to segment the image into different regions. The final number of masks is capped at sixteen by disregarding far-away and trivial objects. The original image and mask-processed images are represented as E and E_i . By putting masks on images, the zero-shot image encoder can enlarge the features of the mask-applied part, and it can be reflected in the output embeddings in a way that allows attention mechanism and contrast learning.

For each action $a \in \{a_f, a_l, a_r\}$, we calculate a reward label r_a considering several factors in this direction: the distance to obstacles, the dynamic/static nature of these obstacles, and the distance to the goal. The backward action is excluded from our model due to its focus on forward-moving scenarios. Mathematically, the reward label for action a can be expressed as:

$$r_a = f(d_{\text{obs}}, d_g, \text{type}_{\text{obs}}) \quad (3)$$

where d_{obs} denotes the distance to the nearest obstacle, d_g represents the distance to the goal, and type_{obs} indicates whether the obstacle is static or dynamic. The function f computes the reward label, encapsulating the trade-offs between navigating safely around obstacles and efficiently moving toward the goal.

The original image E and processed images E_i are labelled as R and R_i . For R_i , we focus on whether the masked object is in the corresponding direction, then denote this with a difference as a mark for contrastive loss. A weight λ is stitched to R_i to decrease the loss value if the masked object belongs to the background.

(2) Linking semantics of images with actions

Based on the differences between the image labels R_i after mask processing and the original label R at each label element, the images that have undergone mask processing are labeled as positive

samples x^+ for the positions with differences, while those without any changes are labeled as negative samples x^- .

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (4)$$

where Q , K , and V represent queries, keys, and values separately.

Then, these samples are sent to a masked multi-head attention network [27], which aims to forge a shared feature representation pivotal for ensuing action decisions. For every pre-defined action, action-specific networks O_a process this shared representation, engendering action determinations. The self-attention mechanism allows the model to focus on specific elements in semantic embeddings. This dynamic is orchestrated by incorporating each action's loss L_a into a cumulative loss function, where the contrastive loss L_c plays a crucial role in guiding the model towards discerning distinct actions.

$$L_{\text{total}} = \sum_a L_a + \lambda \times L_c \quad (5)$$

$$L_c = -\log \frac{\exp(\text{sim}(x, x^+)/\tau)}{\exp(\text{sim}(x, x^+)/\tau) + \sum_{x^-} \exp(\text{sim}(x, x^-)/\tau)} \quad (6)$$

where $\text{sim}(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$ denotes the cosine similarity between two vectors x and y . x is the anchor sample, x^+ is a positive sample similar to the anchor, and x^- represents a negative sample dissimilar to the anchor. τ is a temperature scaling parameter that controls the separation between positive and negative pairs.

The model outputs serve as supplemental guidance for high-level policy in hierarchical reinforcement learning and provide instruction for high-level decisions.

Algorithm 2 Hierarchical path planning

- 1: Initialize high-level and low-level policies
 - 2: Initialize policy network $Q(s, a; \theta)$ and target network $Q'(s, a; \theta^-)$
 - 3: Initialize experience replay buffer D with capacity N
 - 4: **for** each episode **do**
 - 5: **for** each step in episode **do**
 - 6: Observe visual input I_t and current state s_t
 - 7: Generate low-level action $a_t = \pi(s_{low} | I_t)$
 - 8: **if** high-level block is unexplored **then**
 - 9: Generate direction $\vec{F} = -\nabla U$ using APF as equation (14-16)
 - 10: Decompose direction into actions for unexplored blocks
 - 11: **else**
 - 12: Follow policy $\pi_{explored}$ for explored blocks
 - 13: **end if**
 - 14: Execute action a_t , observe reward r_t and next state s_{t+1}
 - 15: Store transition (s_t, a_t, r_t, s_{t+1}) in replay buffer D
 - 16: **end for**
 - 17: Synthesize global strategy $\pi_{global} = \pi_{explored} \cup \pi_{unexplored}$
 - 18: Sample random batch from replay buffer D
 - 19: Compute target Q-value: $y = r_t + \gamma \max_{a'} Q'(s_{t+1}, a'; \theta^-)$
 - 20: Perform gradient descent step on $(y - Q(s_t, a_t; \theta))^2$
 - 21: Update policy network $Q(s, a; \theta)$
 - 22: **if** episode // 100 == 0 **then**
 - 23: Compare π_{global} with historical strategies
 - 24: Update target network $Q'(s, a; \theta^-) = Q(s, a; \theta_{recent})$
 - 25: **end if**
 - 26: Update policy network $Q(s, a; \theta)$ based on target network
 - 27: Dynamically adjust B and b based on R_{bar}
 - 28: **end for**
-

3.3. Hierarchical Policy

In this section, we utilize HRL to complete path-planning tasks under the guidance of leveraging visual cues for the action selection model. A two-level policy layer is conducted to increase efficiency without sacrificing exploring capability.

(1) Environment representation

The environments are represented in two forms, e.g. high-level grid form and low-level ones. The high-level network is designed to achieve global decision-making based on above work, while the lower-level is used for better local planning.

High-level grid map: It selects the next macro-action, directing the agent towards a specific region on the grid. The grid map is divided into larger blocks, and each block is treated as a high-level state. The high-level policy determines the sequence of blocks to be explored based on the current state and the goal location. The settings of the high-level grid help agents avoid areas where visual information has not been fully detected and accelerate the training process in the task of robotics path planning.

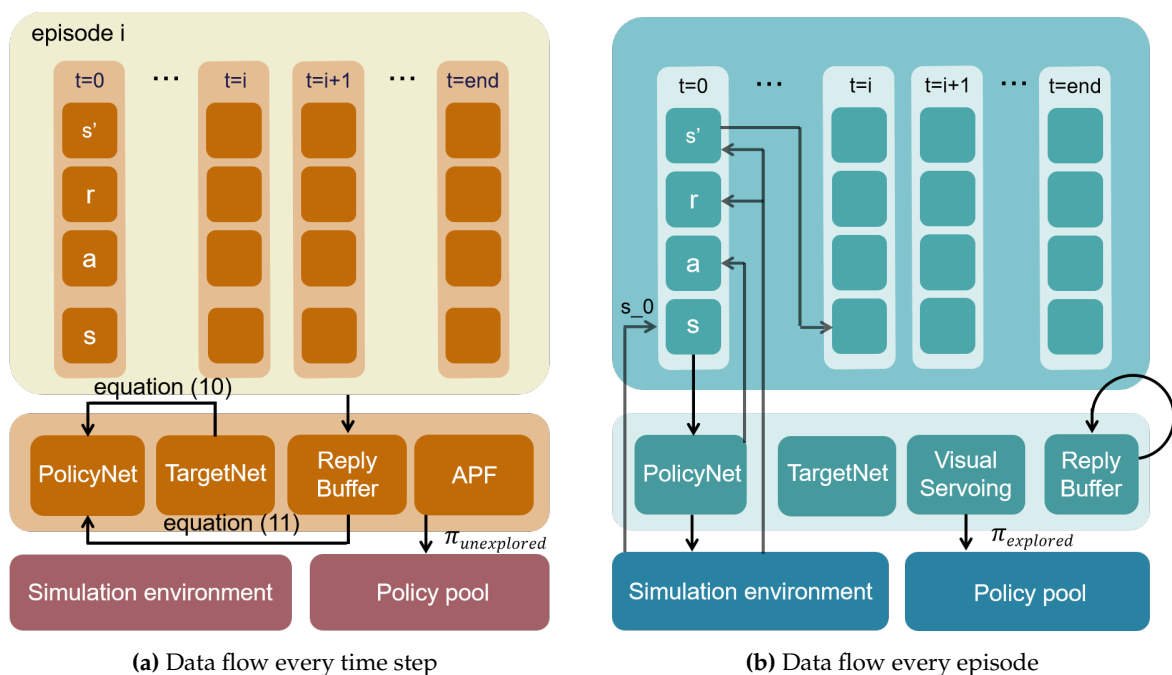


Figure 4. Data flow of the hierarchical policy.

Low-level grid map: Within each selected high-level block, the low-level strategy navigates through the individual grid cells. The policy network produces action instructions for robotics when visual and state information is obtained in grid cells.

$$s_{low} \in \{\text{cell}_{1,1}, \text{cell}_{1,2}, \dots, \text{cell}_{m,m}\} \quad (7)$$

(2) RL structure

We model after the design of DQN and incorporate this structure into the proposed framework: a policy network and a target network. As shown in Figure 5, the policy network is updated continuously based on the Bellman equation and the target network and generates policies for explored areas $\pi_{explored}$. APF is applied to generate policies for unexplored areas $\pi_{unexplored}$. $\pi_{explored}$ and $\pi_{unexplored}$ are combined to form policy π_{global} for the target net. The target network's parameters are updated less frequently, specifically, every 100 episodes, by comparing the global strategy with the most recent policy network.

$$Q(s, a; \theta) \rightarrow \text{policy network} \quad (8)$$

$$Q'(s, a; \theta^-) \rightarrow \text{target network} \quad (9)$$

The action selection is based on the ϵ -greedy policy, balancing exploration and exploitation. The Q-value updates follow the standard Bellman equation:

$$Q(s_t, a_t) = r_t + \gamma \max_{a'} Q'(s_{t+1}, a'; \theta^-) \quad (10)$$

where r_t is the reward at time t , γ is the discount factor, and θ^- are the parameters of the target network.

Additionally, the policy network parameters are updated based on the target network to ensure stability:

$$\theta \leftarrow \theta + \alpha (y - Q(s, a; \theta)) \nabla_{\theta} Q(s, a; \theta) \quad (11)$$

where $y = r_t + \gamma \max_{a'} Q'(s_{t+1}, a'; \theta^-)$.

(3) Base policy and reward settings

Visual servoing: For areas already explored, we employ the visual-act model to map input images to action directives. The network processes the visual input and outputs an action for the agent. During every step, these action instructions are added to a policy pool, which serves as the optimal policy to determine the parameters of the target net.

$$a_{(x,y)} = \pi(s_{low} | I_{(x,y)}) \quad (12)$$

$$\pi_{explored} \leftarrow \bigcup_{(x,y) \in \mathbf{E}}^n a_{(x,y)} \quad (13)$$

where $a_{(x,y)}$ is the action at the coordination (x, y) , $I_{(x,y)}$ is the visual input, and π is the policy function parameterized by the pre-trained visual-act decision model. $\pi_{explored}$ is the policy developed by the pre-trained model for the explored area.

Artificial potential field for unexplored regions: For high-level blocks that have not been explored, we utilize the artificial potential field method to generate navigation strategies. The APF method uses an attractive potential to pull the agent towards the goal and a repulsive potential to avoid obstacles. The direction generated by the APF is then decomposed to produce specific actions. The potential field is represented as U , which can be decomposed into two parts, U_{att} and U_{rep} . The formula is represented as below:

$$U_{att} = \frac{1}{2} k_{att} d^2 \quad (14)$$

$$U_{rep} = \begin{cases} \frac{1}{2} k_{rep} \left(\frac{1}{d_o} - \frac{1}{d_{o0}} \right)^2 & \text{if } d_o \leq d_{o0} \\ 0 & \text{if } d_o > d_{o0} \end{cases} \quad (15)$$

where k_{att} and k_{rep} are constants, d is the distance to the goal, d_o is the distance to the obstacle, and d_{o0} is the influence distance of the obstacle.

The combined potential field U determines the direction \vec{F} for navigation:

$$\pi_{unexplored} \leftarrow \vec{F} = -\nabla U \quad (16)$$

This direction is then translated into specific actions within the unexplored high-level blocks as the $\pi_{unexplored}$. The whole process is conducted after every episode is done.

Global strategy: It is synthesized by combining the policies derived from both the explored and unexplored regions. The synthesized strategy is then compared with the historical strategy obtained from past episodes. The target network parameters are updated based on the comparison to ensure the new strategy optimizes over the past strategies.

$$\pi_{global} = \pi_{explored} \cup \pi_{unexplored} \quad (17)$$

Rewards settings: The reward functions are listed below. The s_{next} represents the next state of the agent. The s_{target} represents the target state. The $R_{collision}$, R_{safe} represent the reward value the agent receives in collision condition and safe condition respectively.

1) target distance

$$R_1(s, a) = -\|s_{next} - s_{target}\| \quad (18)$$

2) obstacle avoidance

$$R_2(s, a) = \begin{cases} R_{collision} & \text{if collision} \\ R_{safe} & \text{otherwise} \end{cases} \quad (19)$$

3) whether the agent reach the target

$$R_3(s, a) = \begin{cases} R_{target} & \text{if } s_{next} = s_{target} \\ R_{step} & \text{otherwise} \end{cases} \quad (20)$$

(4) Adaptive experience replay management

The management of the experience replay buffer [28] is adaptively adjusted based on the cumulative average reward, which serves as an indicator of the learning progress and the efficiency of the current policy. The cumulative average reward, R_{avg} , is computed as follows:

$$R_{avg} = \frac{1}{N} \sum_{i=1}^N R_i \quad (21)$$

where R_i is the immediate reward received after the i -th action, R_g represents rewards of achieving the goal, R_o represents rewards of encountering obstacles, and N is the total number of actions taken up to the current point in time. If the agent achieves the goal or encounters obstacles, R_g or R_o is added to R_{avg} and suspend this episode instantly.

Based on R_{avg} , we adjust the buffer size, B_{size} , and the batch size, b_{size} , of the experience replay buffer to enhance the learning efficiency. The adjustments are made according to the following rules:

$$B_{size} = B_{min} + (B_{max} - B_{min}) \cdot \min\left(\frac{R_{avg}}{R_{target}}, 1\right) \quad (22)$$

$$b_{size} = b_{min} + (b_{max} - b_{min}) \cdot \min\left(\frac{R_{avg}}{R_{target}}, 1\right) \quad (23)$$

where B_{min} and B_{max} represent the minimum and maximum buffer sizes, respectively; similarly, b_{min} and b_{max} denote the minimum and maximum batch sizes. R_{target} is a target average reward that indicates an optimal learning performance. These formulas ensure that as the average reward approaches the target, both the buffer size and the batch size increase, allowing the model to learn from a larger set of experiences. Conversely, if the performance drops, the model focuses on a smaller, potentially more relevant set of experiences to adjust its policy more rapidly.

4. Simulation

In this section, we build up joint simulation environments based on the V-rep platform and Pycharm. We conduct simulations on the environments and validate the following conclusions:

- 1) The leveraging visual cues for actions selection model does focus on certain objects in images to form correspondent advice for action selection.
- 2) To assess the modal's capability to generate across environments, we conduct path planning simulation in unseen environments, which proves that the proposed method can reduce training time and increase exploring efficiency.
- 3) The proposed method performs well when combined with reinforcement learning.

4.1. Simulation Setting

The simulation settings including equipment and environment settings are listed as below.

Equipment: The simulations are conducted on Windows 11 operating system powered by 13th Gen Intel(R) Core(TM) i9-13900HX, with training executed on an NVIDIA GeForce 4090 GPU. The simulation experiment is conducted on V-rep.

Simulation environment: Table 1 presents the parameters of the model car and obstacles in the virtual simulation platform.

Table 1. Parameters of the robot.

Parameter	Value
Wheel radius(m)	0.5
Tread(m)	1
On-board Camera Resolution	512×512
Robot cabinet width(m)	0.9
Robot cabinet length(m)	2
Radius of obstacle(m)	1

Network Parameters of the proposed method: Table 2 presents the parameters of the pre-trained visual information alignment model. The high-level grid size is calculated by the viewing range of the vision sensor.

Table 2. Hyper-parameter of HRL model.

Parameter	Value
Learning rate (High-level) α_1	0.002
Learning rate (Low-level) α_2	0.002
ϵ_{decay}	0.995
ϵ_{min}	0.1
Γ	0.99
High-level grid size	4×4
Low-level grid size	32×32
Goal reward	1000
Max step number	500
Obstacle reward	-200
Buffer size ascending rate	0.9
Buffer size ascending threshold	1.2
Buffer size descending rate	0.9
Buffer size descending threshold	1.0
Initialized buffer size	10000
Initialized batch size	128
Buffer size (min-max)	1000 – 20000
Batch size (min-max)	16 – 256

4.2. Simulation Process

The comparative and validation simulation experiments follow a routine as below:

Step 1: Establish communication between the Python terminal and V-rep, set the simulation step size to 0.05s, and calibrate the initial state of the robot. The maximum range of the vision sensor is 5 meters, and each grid covers 0.625×0.625 square meters, so we would include 8×8 grids in a high-level grid. This grid is then used in the high-level grid map. Initialize parameters of grid map for path planning.

Step 2: The designed control algorithm runs on the Python side, generating the trajectory. This is then converted into control rates for the four wheels at each moment. The control laws are sent to the model car's rotational joints to execute the movement according to the control instructions.

Step 3: In the environment, record the state of key robot nodes and send this information to the Python terminal through the application programming interface.

Step 4: Repeat steps 2 and 3, iterating through the set simulation time and frequency until the entire simulation is completed.

4.3. Assessment of the Visual-Act Model

In this part, we trained our model on several diverse environments and tested its performance on a separate set of unseen environments that follow different distributions with the train set. The training environments utilize cityscapes [29]. Labels are modified to fulfill the model input requirements. It includes various layouts, obstacle types, and goal locations in real road scenes.

As shown in Figure 5, observations of the loss curve revealed a contracting trend during the training process, indicating a consistent enhancement in model accuracy. The training process converges around 200 epochs because the zero-shot image encoder is perfectly pre-trained, which can well reflect the semantic and spatial information of objects in the image in the coding and accelerate the convergence speed in the proposed method training process.

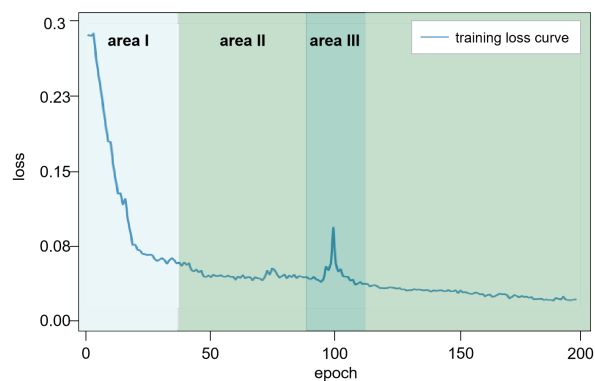


Figure 5. Loss curve of the pre-trained model.

The loss of the model training drops rapidly in Area I, but it picks up in Area III. We speculate that when setting the loss function, we treat the similarity using the sigmoid function to guarantee the artificial number in logarithm loss calculation to be positive, and the gradient is amplified during the loss backward process; thus, the fluctuation occurs.

As shown in Figure 6, the test environments are designed to differ significantly from the training ones regarding layout complexity and obstacle placements, ensuring a rigorous assessment of the model's generalization ability. The virtual environment is primarily designed to imitate a wilderness setting, and its significant differences from the urban road image data in the training set and the inexperience of the pre-trained model in such a scenario further contribute to the rigorosity of the assessment.

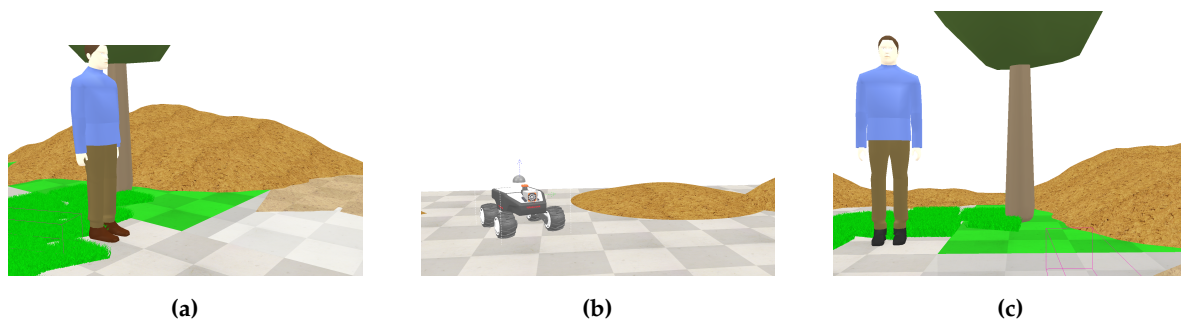


Figure 6. Test set samples.

We conduct a comparative simulation experiment to confirm the focused area of the visual-act model. I_0 represents the embeddings produced by the original image with CLIP. I_1 , I_2 , and I_3 represent the image embeddings processed by the visual-act model. The $\text{sim}(I_0)$ value represents the similarity between the image embeddings and semantic vector of certain objects.

In Table 3, data labeled in bold have the highest similarity to the corresponding object after processing with the attention mechanism in all directions. For example, as shown in Figure 7, the tree is in front of the robot, and the similarity comparison represents that the correspondent processed embeddings raise focus on this obstacle from 0.2607 to 0.4121 while decreasing focus in irrelevant directions. It shows that the pre-trained model focuses on certain obstacles related to correspondent action while disregarding the unrelated ones.

Table 3. Similarities comparison.

Object	$\text{Sim}(I_0)$	$\text{Sim}(I_1)$	$\text{Sim}(I_2)$	$\text{Sim}(I_3)$
Tree	0.2607	0.4121	0.1298	0.0736
People	0.2473	0.5223	0.2954	0.1371
Floor	0.1849	0.0718	0.1035	0.1302
Sand	0.3001	0.1727	0.0931	0.1127
Grass	0.2131	0.1981	0.1703	0.0689

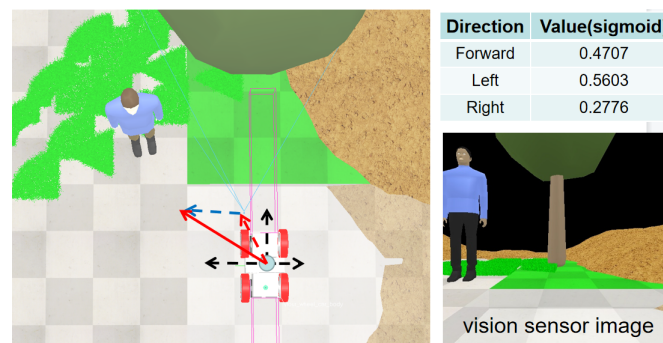


Figure 7. high-level policy combined with action instruction.

4.4. Performance Assessment under Singleton Environment

In this part of the study, we conduct simulation experiments in a singleton environment, where the training and testing datasets are identical. This experimental setup is designed to assess the robustness and effectiveness of our proposed reinforcement learning method for path planning, by directly comparing it to other conventional RL methods under controlled conditions.

(1) Settings

The Singleton environment simplifies the problem space by using the same scenarios in both the training and testing phases, allowing us to focus on the learning and optimization capabilities of the algorithms without the variability introduced by unseen test conditions. The experimental area measures 25×25 meters, and high-level grids are planned based on camera parameters to provide better guidance for macro-action. The representation of one train environment sample is shown in Figure 8.

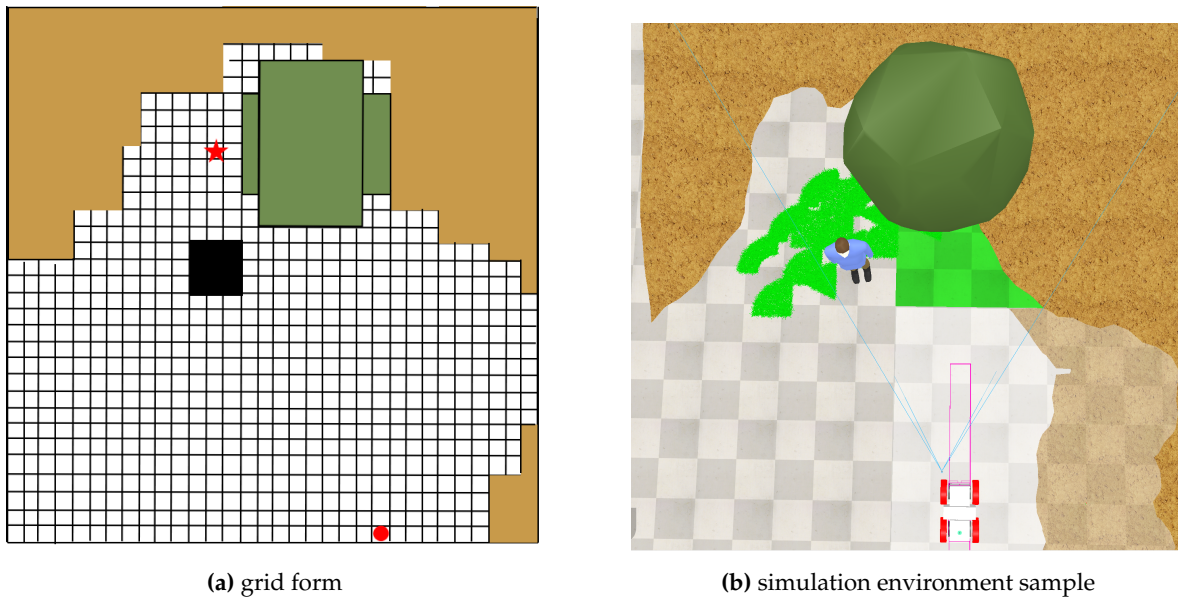


Figure 8. Environment representation, the simulation environments can be represented by grid forms. The red star represents the target, the green circle represents the starting point, the yellow blocks represent sand piles, the green blocks represent trees, and the black blocks represent people.

(2) Metrics

The metrics are listed as below:

- 1) efficiency - the step costs compared to the shortest possible paths.
- 2) final convergence value - the value of cumulative average reward when the reward curve converges.
- 3) convergence degree - the fluctuation trend during the training process, the number of epochs each algorithm took to converge, and the final convergence value.

(3) Outcomes

Figure 9 illustrates the cumulative average reward curve when incorporating the visual-act network into hierarchical reinforcement learning. We can corroborate that the proposed method has a good convergence trend and a low degree of fluctuation. Compared to the traditional HRL method, the convergence episodes decrease from 780 to 480, the final convergence value increases from 24 to 297, which shows great progress in performance.

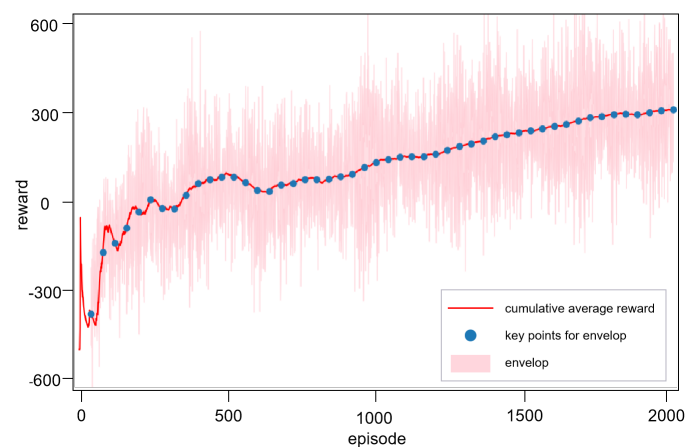


Figure 9. cumulative average rewards curve.

Figure 10 presents gradients formed in the high-level grid and path in inexperienced environments. We can see that the proposed method adequately provides efficient high-level gradient information for path-planning tasks, resulting in an efficient and collision-free final path-planning outcome.

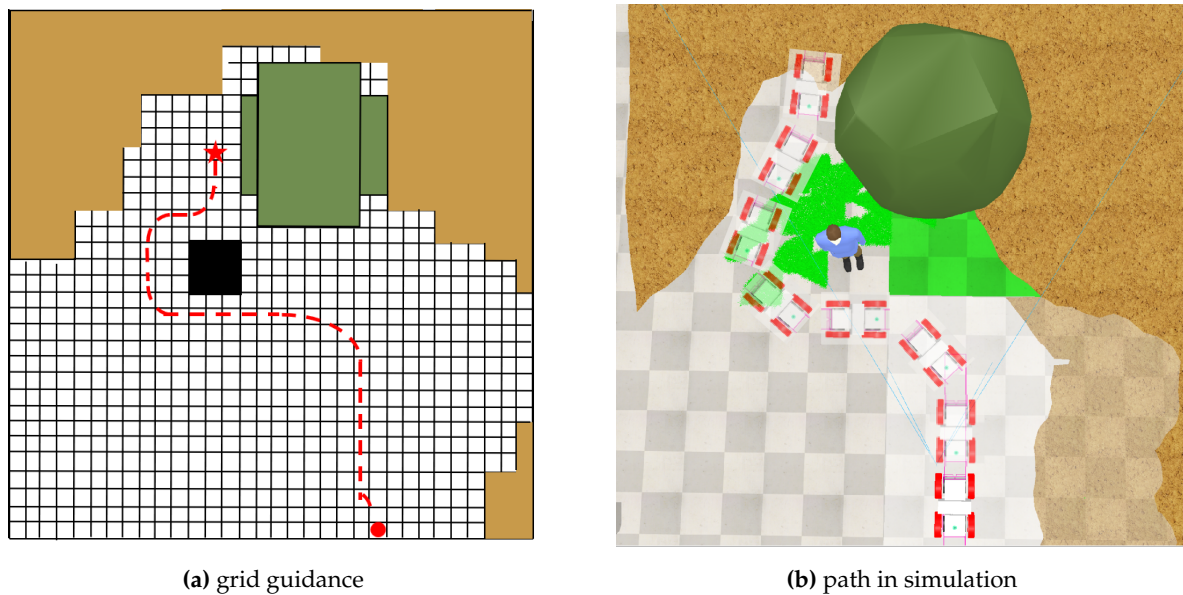


Figure 10. path planning with train environment.

4.5. Performance Assessment under Test Environments

In this part, we conduct the generalization simulation under several environments modeled after the trained ones.

(1) Settings

The test environments differ from the train ones. As shown in Figure 11, they follow a different distribution. The tests are conducted without extra training within unseen test environments. The distinction between the train environments and test environments is listed below:

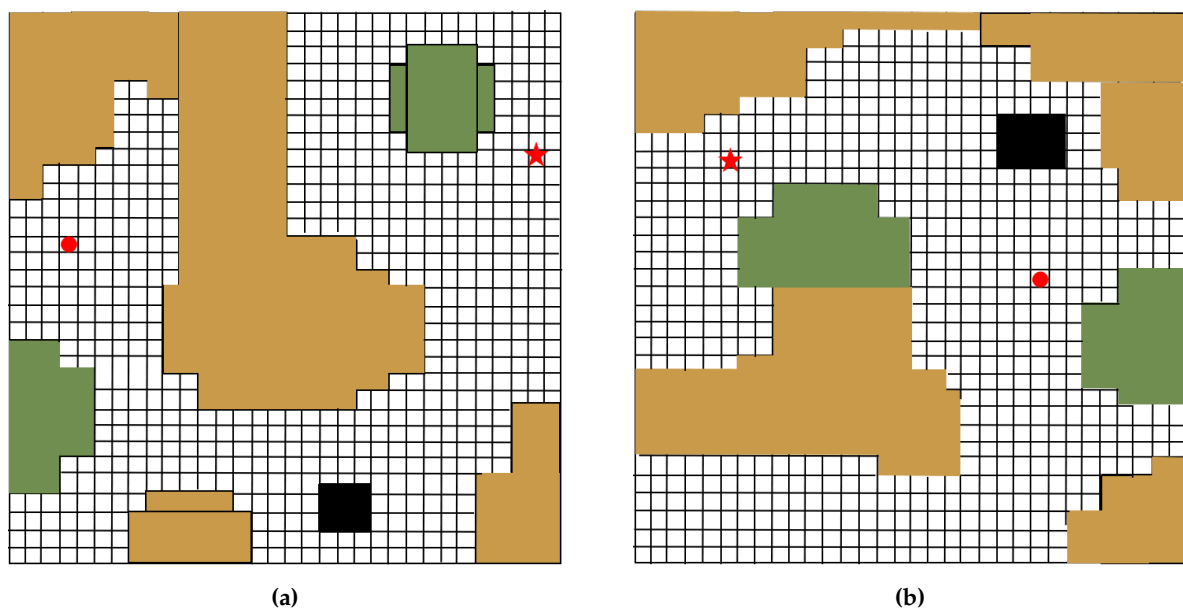


Figure 11. Environment representation samples in grid forms.

1) The starting and target points follow different distributions with the train environments. Specifically, as shown in Figure 11a-c, they all lie in different directions.

detail how the HRL architecture is implemented and propose an adaptive experience pool strategy to balance the exploration capability and convergence dynamically.

We conduct several simulation experiments to validate the proposed method. In cross-environment generalization assessment, we test the proposed method with unseen environments, demonstrating that the proposed method still has good convergence, thus exhibiting strong generalizing performance. The comparative analysis proves that the visual-act model focuses on areas related to decision-making, highlighting the model's capability to leverage visual cues for action selection. Finally, we evaluate the model's performance in path planning when encountering unseen environments. The proposed method improves the robustness and efficiency of robots' path planning and helps to decrease collision rate during the training process.

References

1. F. Duchoň. Path planning with modified a star algorithm for a mobile robot. *Procedia engineering* 96: 59-69, 2014.
2. M. Luo, X. Hou and J. Yang. Surface optimal path planning using an extended Dijkstra algorithm. *IEEE Access* 8: 147827-147838, 2020.
3. J. Kennedy and R. Eberhart. Particle swarm optimization. In: *Proceedings of international conference on neural networks*, pp. 1942-1948, Nov 27-Dec 1, 1995, Perth, Australia.
4. Y. Shi and R. Eberhart. A modified particle swarm optimizer. In: *Proceedings of IEEE international conference on evolutionary computation*, pp. 69-73, May 4-9, 1998, Anchorage, USA.
5. G. Grisetti, R. Kümmerle and W. Burgard. A tutorial on graph-based SLAM. *IEEE Intelligent transportation systems magazine*, 2(4): 31-43, 2010.
6. S. Zang, M. Ding and D. Smith. The impact of adverse weather conditions on autonomous vehicles: How rain, snow, fog, and hail affect the performance of a self-driving car. *IEEE vehicular technology magazine*, 14(2): 103-111, 2019.
7. Y. Wang and S. Zou. Policy gradient method for robust reinforcement learning. In: *Proceedings of the international conference on machine learning*, pp. 23484-23526, July 17-23, 2022, Baltimore, USA.
8. J. Zhang, A. Koppel and A. Bedi. Variational policy gradient method for reinforcement learning with general utilities. *Advances in neural information processing systems*, 33: 4572-4583, 2020.
9. C. Lampert, H. Nickisch and S. Harmeling. Learning to detect unseen object classes by between-class attribute transfer. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 951-958, June 20-25, 2009, Miami, USA.
10. W. Wang, W. Zheng and Y. Han. A survey of zero-shot learning: Settings, methods, and applications. *ACM transactions on intelligent systems and technology*, 10(2): 1-37, 2019.
11. L. Zhang, T. Xiang and S. Gong. Learning a Deep embedding model for zero-Shot learning. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3010-3019, July 21-26, 2017, Hawaii, USA.
12. X. Zhu, R. Zhang and B. He. Pointclip v2: Prompting clip and gpt for powerful 3d open-world learning. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2639-2650, June 18-22, 2023, Vancouver, Canada.
13. Y. Song, N. Liang and Q. Guo. MeshCLIP: Efficient cross-modal information processing for 3D mesh data in zero/few-shot learning. *Information processing and management*, 60(6): 103497, 2023.
14. R. Kirk. A survey of zero-shot generalisation in deep reinforcement learning. *Journal of artificial intelligence research*, 76: 201-264, 2023.
15. L. Ball, and H. Parker. Augmented world models facilitate zero-shot dynamics generalization from a single offline environment. In: *Proceedings of international conference on machine learning*, pp. 619-629, July 18-24, 2021, Graz, Austria.
16. Y. Chen. UAV path planning using artificial potential field method updated by optimal control theory. *International journal of systems science*, 47(6): 1407-1420, 2016.
17. Botvinick and M. Michael. Hierarchical reinforcement learning and decision making. *Current opinion in neurobiology*, 22(6): 956-962, 2012.

18. M. Eppe and C. Kerzel. Intelligent problem-solving as integrated hierarchical reinforcement learning. *Nature machine intelligence*, 4(1): 11-20, 2022.
19. J. Wöhlke, F. Schmitt, and Van H.. Hierarchies of planning and reinforcement learning for robot navigation. In: *Proceedings of IEEE international conference on robotics and automation*, pp. 10682-10688, May 30-June 6, 2021, Xi'an, China.
20. J. Duan, S. Eben and Y. Guan. Hierarchical reinforcement learning for self-driving decision-making without reliance on labeled driving data. *Intelligent transport systems*, 14(5): 297-305, 2020.
21. S. Christen, L. Jendele and E. Aksan. Learning functionally decomposed hierarchies for continuous control tasks with path planning. *IEEE robotics and automation letters*, 6(2): 3623-3630, 2021.
22. X. Ye and Y. Yang. Efficient robotic object search via hiem: Hierarchical policy learning with intrinsic-extrinsic modeling. *IEEE robotics and automation letters*, 6(3): 4425-4432, 2021.
23. B. Lin, Y. Zhu and Z. Chen. Adapt: Vision-language navigation with modality-aligned action prompts. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 15396-15406, June 18-24, 2022, New Orleans, USA.
24. E. Kotei and R. Thirunavukarasu. A systematic review of transformer-based pre-trained language models through self-supervised learning. *Information*, 14(3): 187, 2023.
25. A. Jaiswal, A. Babu, and M. Zadeh. A survey on contrastive self-supervised learning. *Technologies*, 9(1): 2, 2020.
26. X. Huo, H. Karimi and X. Zhao. Adaptive-critic design for decentralized event-triggered control of constrained nonlinear interconnected systems within an identifier-critic framework. *IEEE transactions on cybernetics*, 52(8): 7478-7491, 2021.
27. A. Vaswani, N. Shazeer and N. Parmar. Attention is all you need. In: *Proceedings of the conference on neural information processing systems*, pp. 6000-6010, Dec 4-9, 2017, Long Beach, USA.
28. R. Tiwari, K. Killamsetty and R. Iyer. Gcr: Gradient coreset based replay buffer selection for continual learning. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 99-108, June 19-24, 2022, New Orleans, USA.
29. C. Marius, O. Mohamed and R. Sebastian. The cityscapes dataset for semantic urban scene understanding. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3213-3223, June 27-30, 2016, Las Vegas, USA.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.