



南京航空航天大学
Nanjing University of Aeronautics and Astronautics

Learning Representations in Model-Free Hierarchical Reinforcement Learning

Jacob Rafati David C. Noelle

AAAI 2018

How to Learn with sparse delayed reward ?

(1)



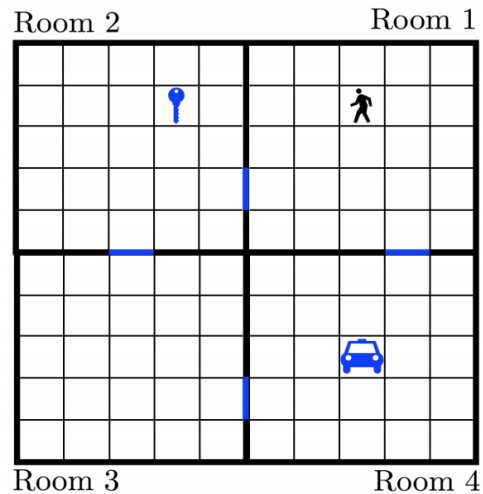
take “play” :

$r_{t+1}=1, r_{t+2}=...=r_{t+99}=0, r_{t+100}=-100;$

take “learn” :

$r_{t+1}=-1, r_{t+2}=...=r_{t+99}=0, r_{t+100}=100;$

(2)



action = {North, South, East, West}

$r = +10$ reaching the key

$r = +100$ moves to the car while carrying the key

$r = -2$ bumping to the wall

no reward or punishment for exploring the space

Main Idea :

Rely on an upper-level policy to decompose the entire task, and then use the lower-level policy to gradually execute。

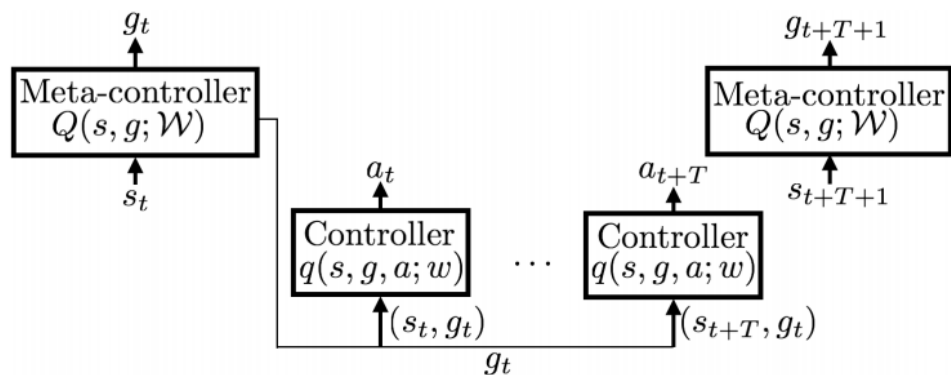
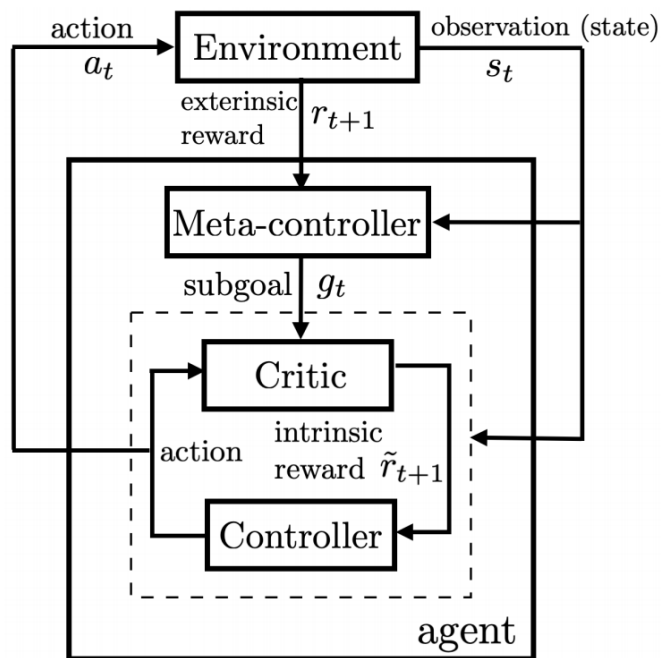
Problems :

Subprolem1 : Learning a meta-policy to choose a subgoal

Subprolem2 : Exploring the state space while learning subtask through intrinsic motivation

Subprolem3 : Subgoal discovery

Meta-controller/Controller Framework



meta-controller :

$$G_t = \sum_{t'=t}^{\tau} \gamma^{t'-t} r_{t'}$$

$$Q(s, g) = \mathbf{E}_{\pi_g} [G_t | s_t = s, g_t = g]$$

$$\mathcal{L}(\mathcal{W}) \triangleq \mathbb{E}_{(s, g, G, s_{t'}) \sim \mathcal{D}_2} [(G + \gamma \max_{g'} Q(s_{t'}, g'; \mathcal{W}) - Q(s, g; \mathcal{W}))^2]$$

controller :

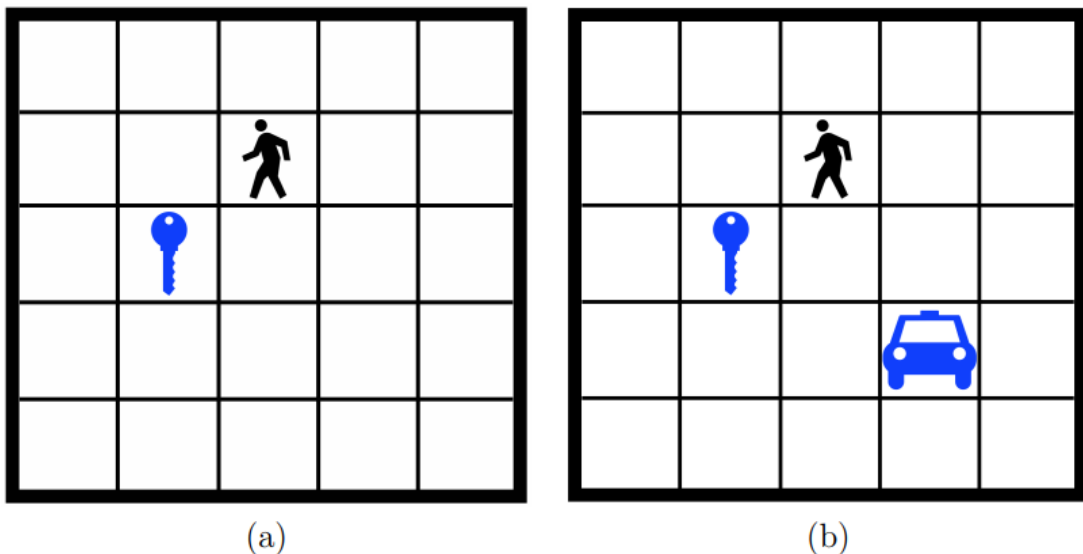
$$\tilde{G}_t = \sum_{t'=t}^{t+T} \gamma^{t'-t} \tilde{r}_{t'}(g)$$

$$q(s, g, a) = \mathbf{E}_{\pi_{ag}} [\tilde{G}_t | s_t = s, g_t = g, a_t = a]$$

$$L(w) \triangleq \mathbb{E}_{(s, g, a, \tilde{r}, s') \sim \mathcal{D}_1} [(\tilde{r} + \gamma \max_{a'} q(s', g, a'; w) - q(s, g, a; w))^2]$$

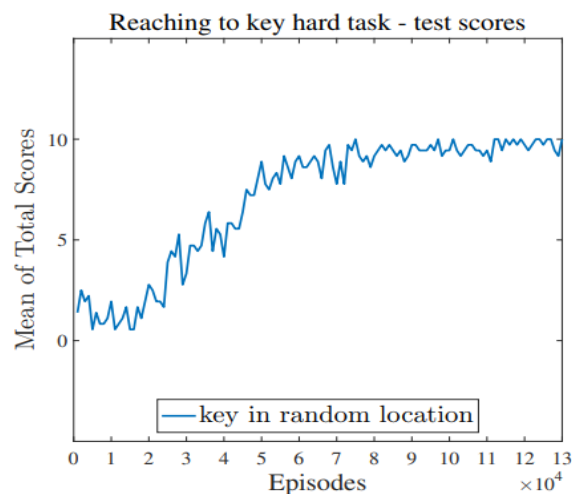
Assuming that there is an oracle to give almost good subgoals , at least two benefits can get:

- (1) exploration of large scale state spaces
- (2) enabling the reuse of skills in varied environments

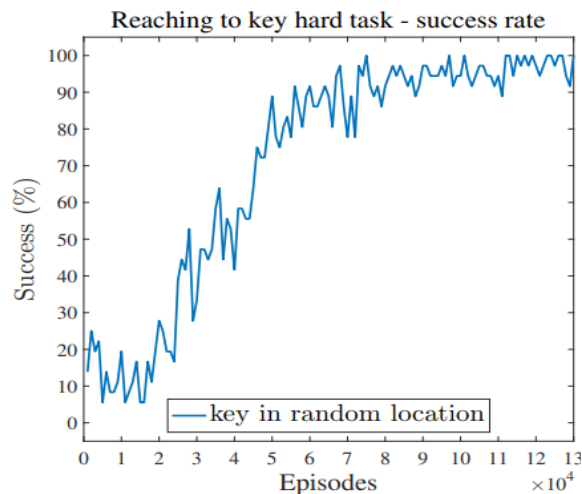


$$\tilde{r}_{t+1} = \begin{cases} \min(r_{t+1}, -1) & \text{if } s_{t+1} \text{ is not terminal} \\ +1 & \text{if } s_{t+1} \text{ achieves the goal, } g_t \end{cases}$$

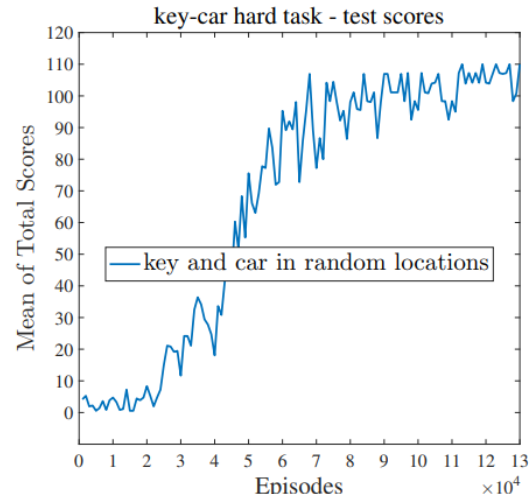
- **Key Task, Hard Placement.** In this simplified version of the task, the agent was trained to move to the key, producing a policy, π_{ag} , for reaching a randomly located goal g (key). This is illustrated in Figure 7(a). For each starting $s \in S$, a random goal, g , was assigned and the cumulative reward was obtained. We report the average reward scores and the average success percentage in Figure 8 (a) and (b), respectively.
- **Key Task, Easy Placement.** This version of the task is the same as the last, except that the goal, g , was always randomly placed in a location adjacent to the starting state, s . (See Figure 7 (a).) We report the average reward scores and the average success percentage in Figure 8 (c) and (d), respectively.
- **Key-Car Task, Hard Placement.** In this version of the task, both the key, g_{key} , and the car, g_{car} , were randomly placed. The agent received positive reward when the agent moved to the key (+10) and subsequently moved to the car (+100). (See Figure 7 (b).) We report the average scores and the average success percentage in Figure 9 (a) and (b), respectively.
- **Key-Car Task, Easy Placement.** This version of the task is the same as the last, except that the key was always located at (0,0), and the car was always located at (1,1). We report the average reward scores and the average success percentage in Figure 9 (c) and (d), respectively.



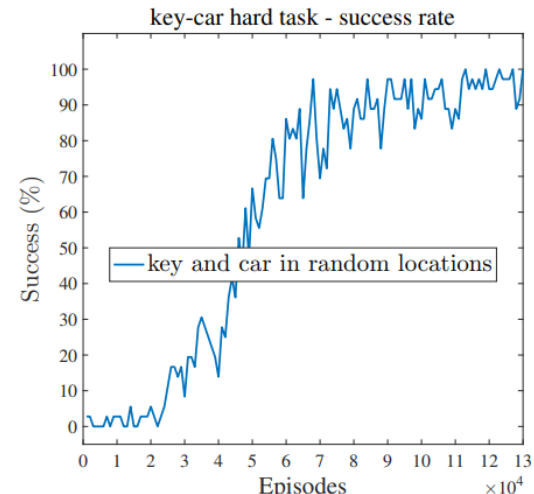
(a)



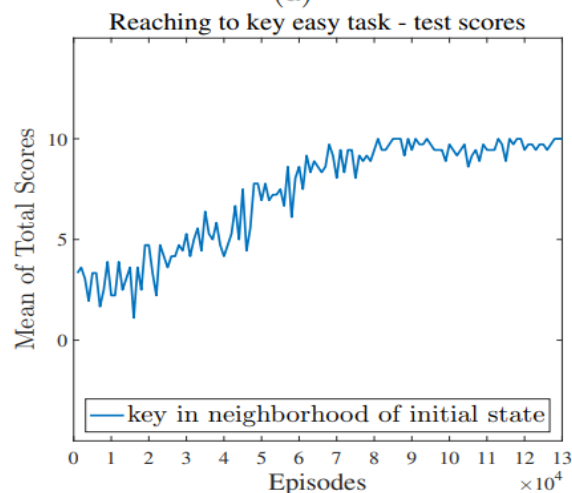
(b)



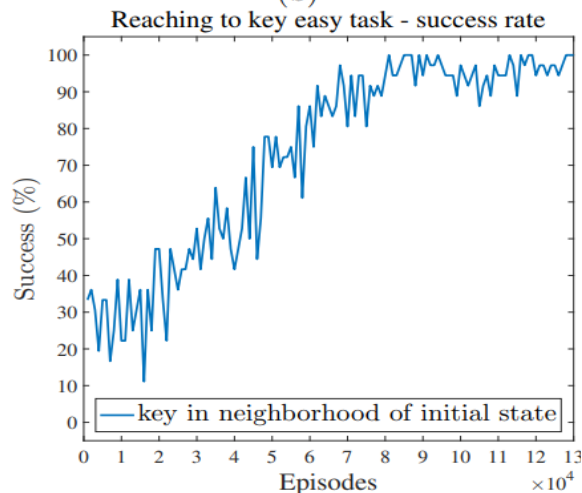
(a)



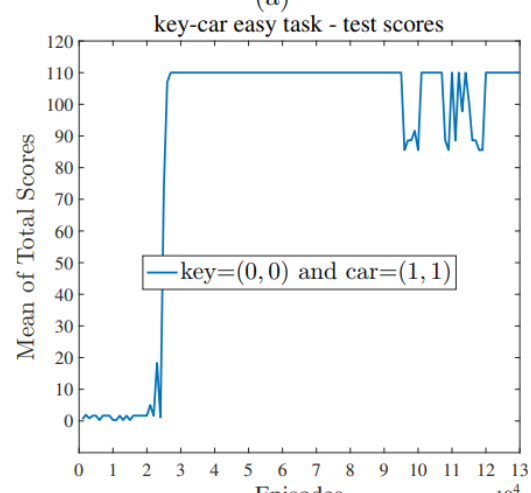
(b)



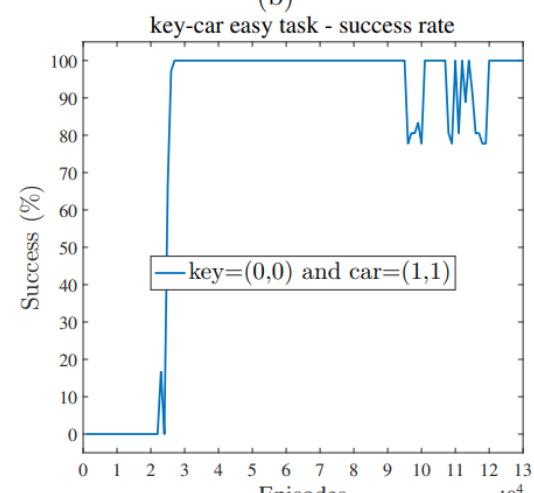
(c)



(d)

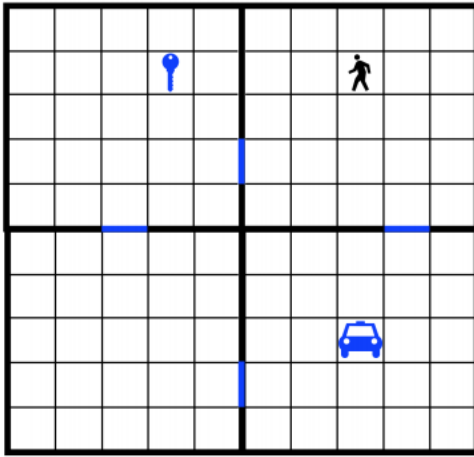


(c)

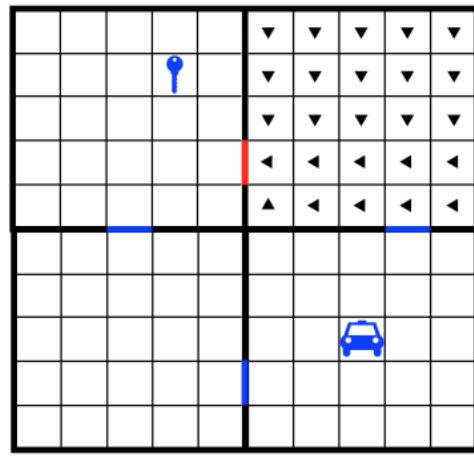


(d)

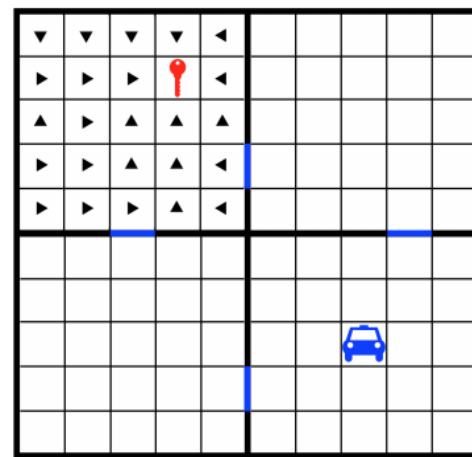
Reusing Learned Skills



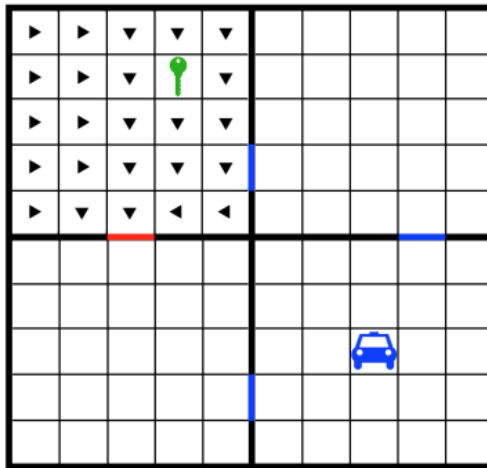
(a)



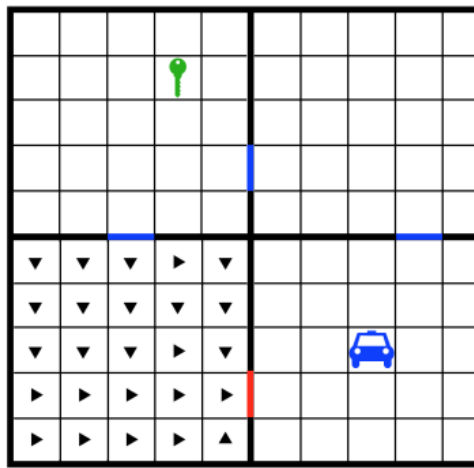
(b)



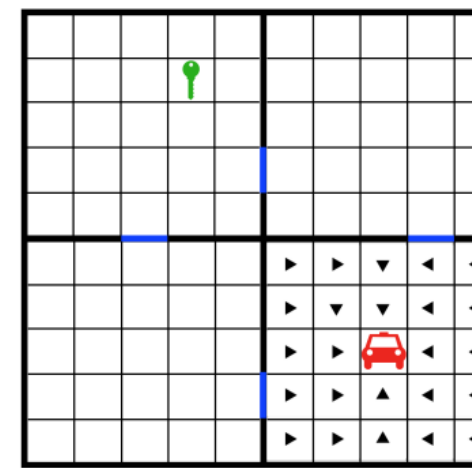
(c)



(d)



(e)



(f)

Good Subgoal Assumptions

- (1) attending to the states associated with anomalous transition experiences.
(large rewards、 large changes in state features)
- (2) clustering experiences based on a similarity measure and collecting the set of associated states into a potential subgoal.

Methods:

merges anomaly (outlier) detection with the K-means clustering of experiences.

Algorithm 4 Unsupervised Subgoal Discovery Algorithm

```
for each  $e = (s, a, r, s')$  stored in  $\mathcal{D}$  do  
    if experience  $e$  is an outlier (anomaly) then  
        Store  $s'$  to the subgoals set  $\mathcal{G}$   
        Remove  $e$  from  $\mathcal{D}$   
    end if  
end for
```

Fit a K -means Clustering Algorithm on \mathcal{D} using previous centroids as initial points
Store the updated centroids to the subgoals set \mathcal{G}

Algorithm 5 Unified Model-Free HRL Algorithm

```

Pretrain controller using Algorithm 2 on a set of random subgoals  $\mathcal{G}'$ 
Initialize experience memories  $\mathcal{D}$ ,  $\mathcal{D}_1$  and  $\mathcal{D}_2$ 
Walk controller for  $M'$  episodes on random subgoals  $\mathcal{G}'$ , and store  $(s, a, s', r)$  to  $\mathcal{D}$ 
Run Unsupervised Subgoal Discovery on  $\mathcal{D}$  to initialize  $\mathcal{G}$ 
for episode = 1, ...,  $M$  do
    Initialize state  $s_0 \in \mathcal{S}$ ,  $s \leftarrow s_0$ 
     $G \leftarrow 0$ 
     $g \leftarrow \text{EPSILON-GREEDY}(Q(s, \mathcal{G}; \mathcal{W}), \epsilon_2)$ 
    repeat for each step  $t = 1, \dots, T$ 
        compute  $q(s, g, a; w)$ 
         $a \leftarrow \text{EPSILON-GREEDY}(q(s, g, \mathcal{A}; w), \epsilon_1)$ 
        Take action  $a$ , observe  $s'$  and external reward  $r$ 
        Compute intrinsic reward  $\tilde{r}$  from internal critic
        Store controller's intrinsic experience,  $(s, g, a, \tilde{r}, s')$  to  $\mathcal{D}_1$ 
        Store agent's transition experience,  $(s, a, r, s')$  to  $\mathcal{D}$ 
        Sample  $J_1 \subset \mathcal{D}_1$  and compute  $\nabla L$ 
        Update controller's parameters,  $w \leftarrow w - \alpha_1 \nabla L$ 
        Sample  $J_2 \subset \mathcal{D}_2$  and compute  $\nabla \mathcal{L}$ 
        Update meta-controller's parameters,  $\mathcal{W} \leftarrow \mathcal{W} - \alpha_2 \nabla \mathcal{L}$ 
         $s \leftarrow s'$ ,  $G \leftarrow G + r$ 
        Decay exploration rate of controller  $\epsilon_1$ 
        if experience  $e$  is an outlier (anomaly) then
            Store  $s'$  to the subgoals set  $\mathcal{G}$ 
            Remove  $e$  from  $\mathcal{D}$ 
        end if
    until  $s$  is terminal or subgoal  $g$  is attained
    Decay exploration rate of meta-controller  $\epsilon_2$ 
    Store meta-controller's experience,  $(s_0, g, G, s')$  to  $\mathcal{D}_2$ 
    Fit a  $K$ -means clustering on  $\mathcal{D}$  every  $N$  step to update centroids of  $\mathcal{G}$ 
end for

```

Algorithm 2 Intrinsic Motivation Learning

```

Specify Subgoals space  $\mathcal{G}$ 
Initialize  $w$  in  $q(s, g, a; w)$ 
Initialize controller's experience memory,  $\mathcal{D}_1$ 
Initialize agent's experience memory,  $\mathcal{D}$ 
for episode = 1, ...,  $M$  do
    Initialize state  $s_0 \in \mathcal{S}$ ,  $s \leftarrow s_0$ 
    Select a random subgoal  $g$  from  $\mathcal{G}$ 
    repeat for each step  $t = 1, \dots, T$ 
        compute  $q(s, g, a; w)$ 
         $a \leftarrow \text{EPSILON-GREEDY}(q(s, g, \mathcal{A}; w), \epsilon_1)$ 
        Take action  $a$ , observe  $s'$  and external reward  $r$ 
        Compute intrinsic reward  $\tilde{r}$  from internal critic
        Store controller's intrinsic experience,  $(s, g, a, \tilde{r}, s')$  to  $\mathcal{D}_1$ 
        Store agent's experience,  $(s, a, s', r)$  to  $\mathcal{D}$ 
        Sample  $J_1 \subset \mathcal{D}_1$  and compute  $\nabla L$ 
        Update controller's parameters,  $w \leftarrow w - \alpha_1 \nabla L$ 
         $s \leftarrow s'$ 
        Decay exploration rate of controller  $\epsilon_1$ 
    until  $s$  is terminal or subgoal  $g$  is attained
end for

```

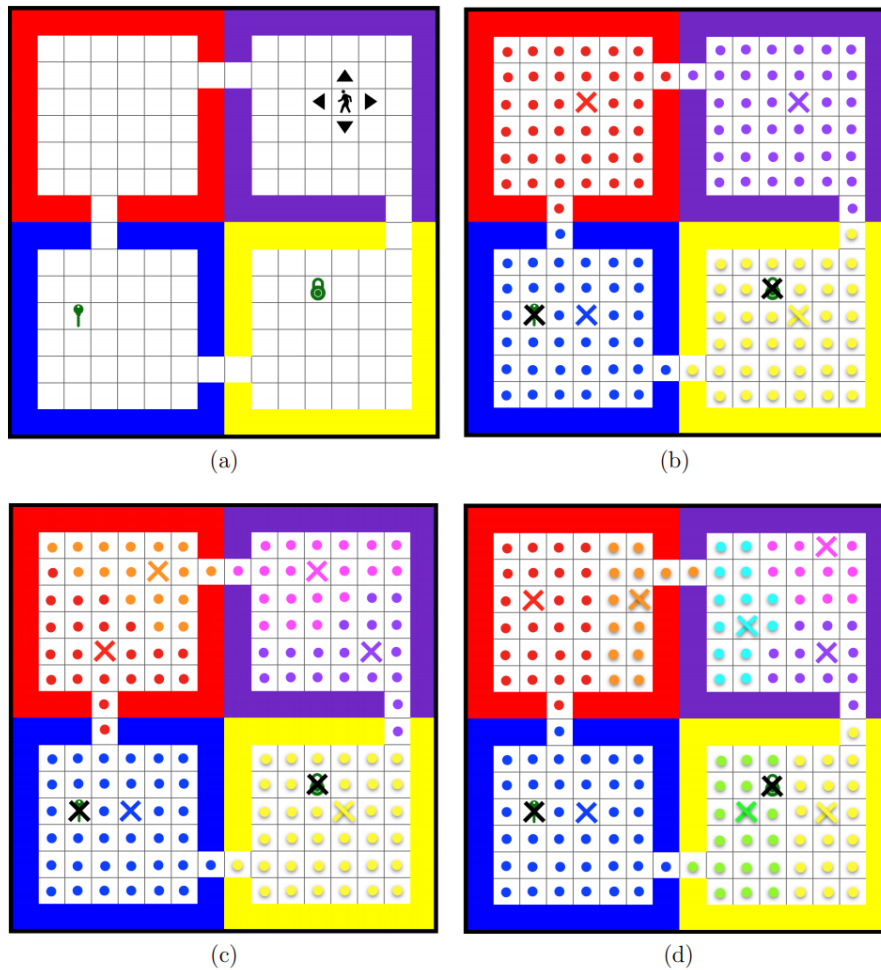
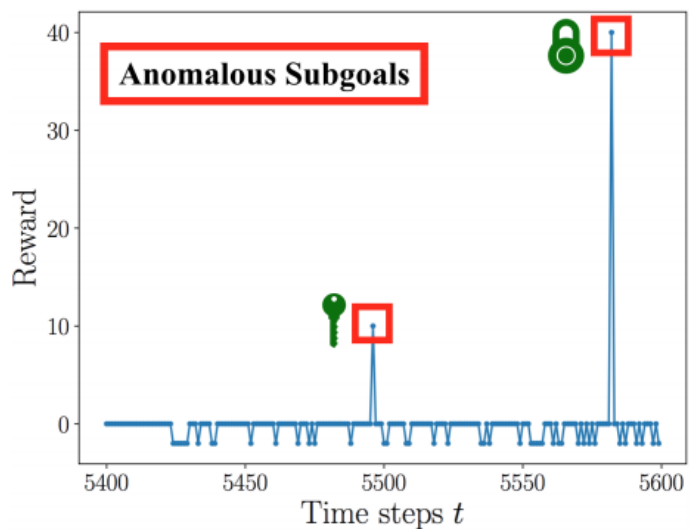
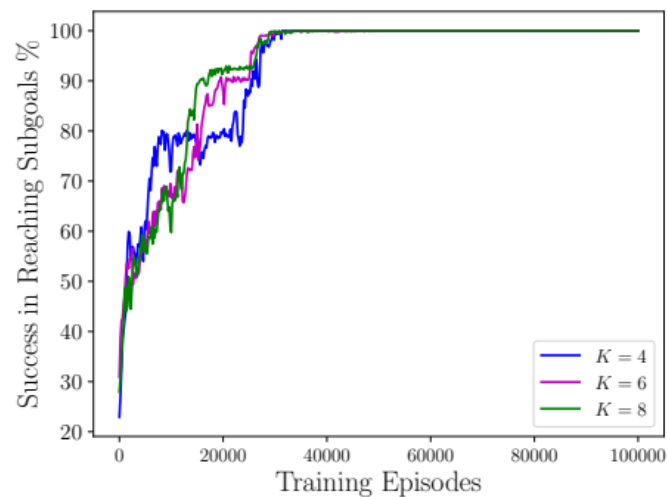


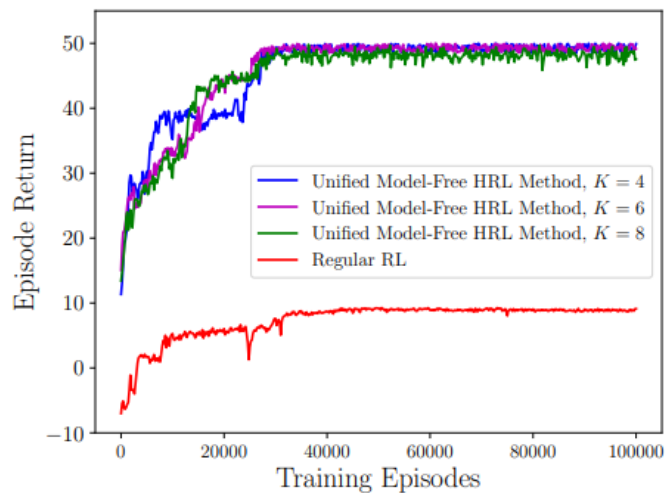
Figure 12: (a) The 4-room task with a key and a lock. (b) The results of the unsupervised subgoal discovery algorithm with *anomalies* marked with black Xs and *centroids* with colored ones. The number of clusters in K-means algorithm was set to $K = 4$. (c) The result of the unsupervised subgoal discovery for $K = 6$. (d) The results of the unsupervised subgoal discovery for $K = 8$.



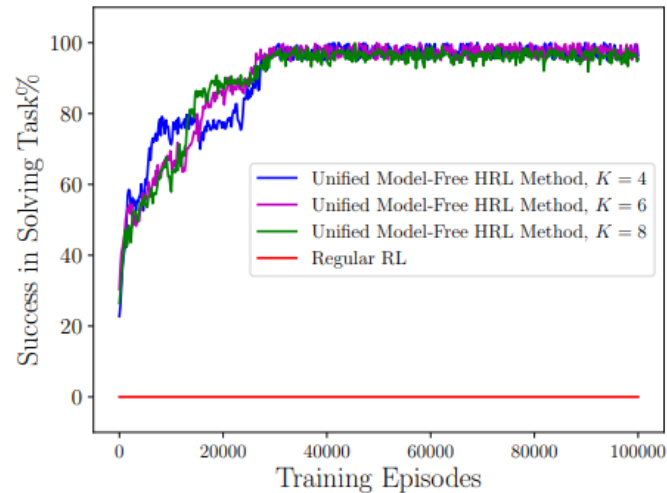
(a)



(b)



(c)



(d)

Montezuma's Revenge

