
Hierarchical Imitation and Reinforcement Learning

Hoang M. Le¹ Nan Jiang² Alekh Agarwal² Miroslav Dudík² Yisong Yue¹ Hal Daumé III^{3,2}

ICML-2018

Hierarchical RL

Idea: choosing subtasks, finish subtasks.

Ex. Leave a building: go to the elevator → take the elevator down → walk out

simultaneously learn **HI-level** policy $\mu : S \rightarrow G$, called the **metacontroller**, and **LO-level** policy $\pi_g : S \rightarrow A$, called the **subpolicies**.

- 1: **for** $h_{\text{HI}} = 1 \dots H_{\text{HI}}$ **do**
- 2: observe state s and choose subgoal $g \leftarrow \mu(s)$
- 3: **for** $h_{\text{LO}} = 1 \dots H_{\text{LO}}$ **do**
- 4: observe state s
- 5: **if** $\beta_g(s)$ **then break**
- 6: choose action $a \leftarrow \pi_g(s)$

Termination
function

types of supervision

- $\text{HierDemo}(s)$: *hierarchical demonstration*. The expert executes its hierarchical policy starting from s and returns the resulting hierarchical trajectory $\sigma^* = (s_1^*, g_1^*, \tau_1^*, s_2^*, g_2^*, \tau_2^*, \dots)$, where $s_1^* = s$.
- $\text{Label}_{\text{HI}}(\tau_{\text{HI}})$: *HI-level labeling*. The expert provides a good next subgoal at each state of a given HI-level trajectory $\tau_{\text{HI}} = (s_1, g_1, s_2, g_2, \dots)$, yielding a labeled data set $\{(s_1, g_1^*), (s_2, g_2^*), \dots\}$.
- $\text{Label}_{\text{LO}}(\tau; g)$: *LO-level labeling*. The expert provides a good next primitive action towards a given subgoal g at each state of a given LO-level trajectory $\tau = (s_1, a_1, s_2, a_2, \dots)$, yielding a labeled data set $\{(s_1, a_1^*), (s_2, a_2^*), \dots\}$.
- $\text{Inspect}_{\text{LO}}(\tau; g)$: *LO-level inspection*. Instead of annotating every state of a trajectory with a good action, the expert only verifies whether a subgoal g was accomplished, returning either *Pass* or *Fail*.
- $\text{Label}_{\text{FULL}}(\tau_{\text{FULL}})$: *full labeling*. The expert labels the agent's full trajectory $\tau_{\text{FULL}} = (s_1, a_1, s_2, a_2, \dots)$, from start to finish, ignoring hierarchical structure, yielding a labeled data set $\{(s_1, a_1^*), (s_2, a_2^*), \dots\}$.
- $\text{Inspect}_{\text{FULL}}(\tau_{\text{FULL}})$: *full inspection*. The expert verifies whether the agent's overall goal was accomplished, returning either *Pass* or *Fail*.

Algorithm 1

Algorithm 1 Hierarchical Behavioral Cloning (**h-BC**)

- 1: Initialize data buffers $\mathcal{D}_{\text{HI}} \leftarrow \emptyset$ and $\mathcal{D}_g \leftarrow \emptyset, g \in \mathcal{G}$
 - 2: **for** $t = 1, \dots, T$ **do**
 - 3: Get a new environment instance with start state s
 - 4: $\sigma^* \leftarrow \text{HierDemo}(s)$
 - 5: **for all** $(s_h^*, g_h^*, \tau_h^*) \in \sigma^*$ **do**
 - 6: Append $\mathcal{D}_{g_h^*} \leftarrow \mathcal{D}_{g_h^*} \cup \tau_h^*$
 - 7: Append $\mathcal{D}_{\text{HI}} \leftarrow \mathcal{D}_{\text{HI}} \cup \{(s_h^*, g_h^*)\}$
 - 8: Train subpolicies $\pi_g \leftarrow \text{Train}(\pi_g, \mathcal{D}_g)$ for all g
 - 9: Train meta-controller $\mu \leftarrow \text{Train}(\mu, \mathcal{D}_{\text{HI}})$
-

Algorithm 2 Hierarchically Guided Dagger (hg-Dagger)

6-11: execute agent's policy

13-19: label

- High-level
- Low-level

```
3: for  $t = T_{\text{warm-start}} + 1, \dots, T$  do
4:   Get a new environment instance with start state  $s$ 
5:   Initialize  $\sigma \leftarrow \emptyset$ 
6:   repeat
7:      $g \leftarrow \mu(s)$ 
8:     Execute  $\pi_g$ , obtain LO-level trajectory  $\tau$ 
9:     Append  $(s, g, \tau)$  to  $\sigma$ 
10:     $s \leftarrow$  the last state in  $\tau$ 
11:   until end of episode
12:   Extract  $\tau_{\text{FULL}}$  and  $\tau_{\text{HI}}$  from  $\sigma$ 
13:   if  $\text{Inspect}_{\text{FULL}}(\tau_{\text{FULL}}) = \text{Fail}$  then
14:      $\mathcal{D}^* \leftarrow \text{Label}_{\text{HI}}(\tau_{\text{HI}})$ 
15:     Process  $(s_h, g_h, \tau_h) \in \sigma$  in sequence as long as
        $g_h$  agrees with the expert's choice  $g_h^*$  in  $\mathcal{D}^*$ :
16:       if  $\text{Inspect}(\tau_h; g_h) = \text{Fail}$  then
17:         Append  $\mathcal{D}_{g_h} \leftarrow \mathcal{D}_{g_h} \cup \text{Label}_{\text{LO}}(\tau_h; g_h)$ 
18:         break
19:       Append  $\mathcal{D}_{\text{HI}} \leftarrow \mathcal{D}_{\text{HI}} \cup \mathcal{D}^*$ 
20:   Update subpolicies  $\pi_g \leftarrow \text{Train}(\pi_g, \mathcal{D}_g)$  for all  $g$ 
21:   Update meta-controller  $\mu \leftarrow \text{Train}(\mu, \mathcal{D}_{\text{HI}})$ 
```

Algorithm 3

Hierarchically Guided Dagger / Q-learning (hg-Dagger/Q)

Train subpolicy with RL
based on pseudo reward

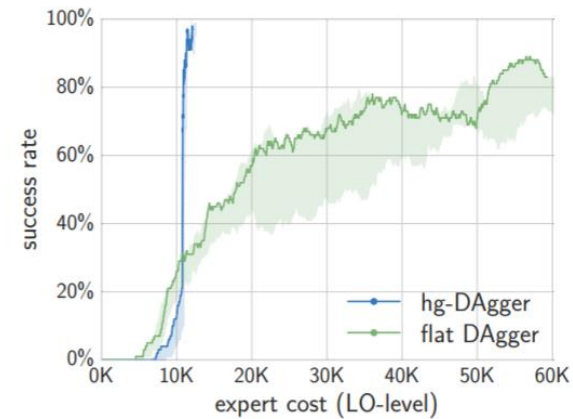
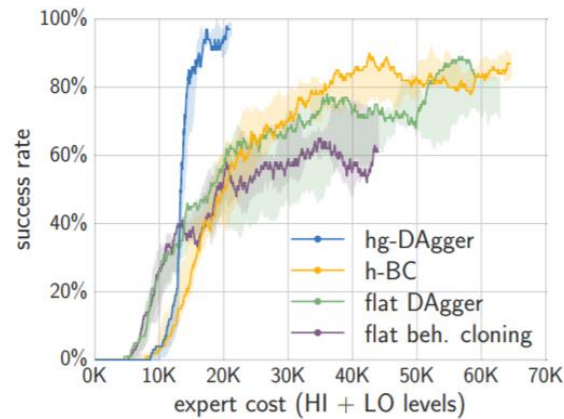
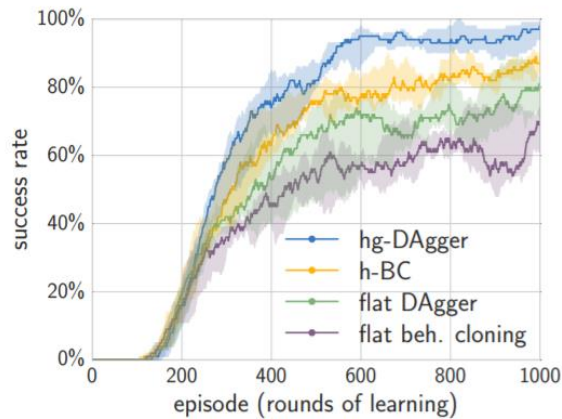
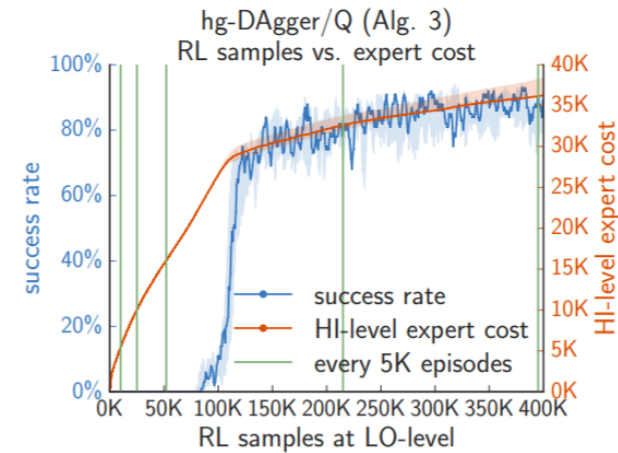
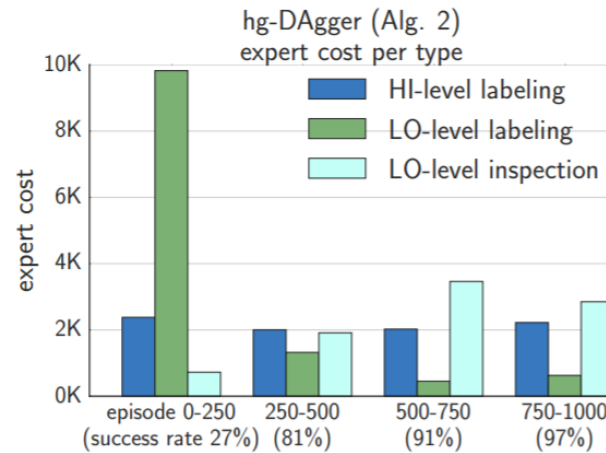
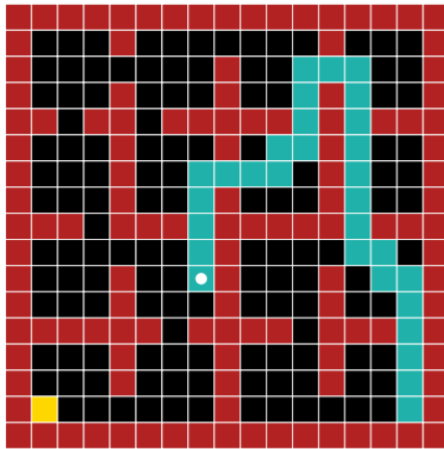
$$\begin{cases} 1 & \text{if success}(s; g) \\ -1 & \text{if } \neg \text{success}(s; g) \text{ and } \text{terminal}(s; g) \\ -\kappa & \text{otherwise,} \end{cases}$$

```

3: for  $t = 1, \dots, T$  do
4:   Get a new environment instance with start state  $s$ 
5:   Initialize  $\sigma \leftarrow \emptyset$ 
6:   repeat
7:      $s_{\text{HI}} \leftarrow s$ ,  $g \leftarrow \mu(s)$  and initialize  $\tau \leftarrow \emptyset$ 
8:     repeat
9:        $a \leftarrow \epsilon_g$ -greedy( $Q_g, s$ )
10:      Execute  $a$ , next state  $\tilde{s}$ ,  $\tilde{r} \leftarrow \text{pseudo}(\tilde{s}; g)$ 
11:      Update  $Q_g$ : a (stochastic) gradient descent step
                           on a minibatch from  $\mathcal{D}_g$ 
12:      Append  $(s, a, \tilde{r}, \tilde{s})$  to  $\tau$  and update  $s \leftarrow \tilde{s}$ 
13:     until  $\text{terminal}(s; g)$ 
14:     Append  $(s_{\text{HI}}, g, \tau)$  to  $\sigma$ 
15:   until end of episode
16:   Extract  $\tau_{\text{FULL}}$  and  $\tau_{\text{HI}}$  from  $\sigma$ 
17:   if  $\text{Inspect}_{\text{FULL}}(\tau_{\text{FULL}}) = \text{Fail}$  then
18:      $\mathcal{D}^* \leftarrow \text{Label}_{\text{HI}}(\tau_{\text{HI}})$ 
19:     Process  $(s_h, g_h, \tau_h) \in \sigma$  in sequence as long as
                            $g_h$  agrees with the expert's choice  $g_h^*$  in  $\mathcal{D}^*$ :
20:     Append  $\mathcal{D}_{g_h} \leftarrow \mathcal{D}_{g_h} \cup \tau_h$ 
                           Append  $\mathcal{D}_{\text{HI}} \leftarrow \mathcal{D}_{\text{HI}} \cup \mathcal{D}^*$ 
21:   else
22:     Append  $\mathcal{D}_{g_h} \leftarrow \mathcal{D}_{g_h} \cup \tau_h$  for all  $(s_h, g_h, \tau_h) \in \sigma$ 
23:   Update meta-controller  $\mu \leftarrow \text{Train}(\mu, \mathcal{D}_{\text{HI}})$ 

```

Maze Navigation Domain



Montezuma's Revenge

