

Hierarchical Imitation and Reinforcement Learning:

<https://arxiv.org/pdf/1803.00590.pdf>

Paper Summary Notes

1 Main Ideas, High Level Assumptions

Central Question of the Paper: *When experts are available, how can we most effectively leverage their feedback?*

- **Context:**
 - standard imitation learning can solve long-horizon sparse rewards, but may require lots of expert behaviour
 - hierarchy helps long horizon problems by decomposing into sub-problems and acting recursively optimally
 - *Hierarchical Guidance*: framework where high level expert focuses (guides) low-level learner.
 - * High level expert ensures low-level learning only occurs when necessary, and only over relevant parts of state space, thereby speeding up learning.
 - * Omitting feedback on already mastered subtasks reduces expert effort.
- **Goals:**
 - reduce *expert cost* which defines how much effort the expert needs to make the system work
 - decompose problem into hierarchy to reduce cost of exploration
 - faster learning than standard hierarchical RL, and more label-efficient than normal Imitation learning

2 Formulation

2.1 Hierarchical RL Basics

```

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  $a \leftarrow \pi_g(s)$ 

```

- *low level trajectory*: $\langle s_t, a_t \rangle$ tuples of low level actions taken for a given goal g
- *high level trajectory*: $\langle s_t, g_t \rangle$ tuples of states and goals, where s_{t+1} is given by the terminal state from low level policy acting from state s_t under goal g_t
- *full trajectory*: concatenation of all low level actions throughout the entire hierarchical trajectory

2.2 Expert Feedback Supervision Types

- **Hierarchical Demonstration** HierDemo(s): expert executes hierarchical policy starting from s and returns the hierarchical trajectory (passive IL)
- **High Level Labelling** LabelHigh(τ_{high}): expert provides next subgoal at each state s of a high level trajectory giving labelled dataset of tuples $\langle s, g, \rangle$
- **Low Level Labelling** LabelLow(τ_{low}): expert provides next primitive actions towards subgoal g at each state s of a low level trajectory giving labelled dataset of tuples $\langle s, a, \rangle$

- **Full Level Labelling** $\text{LabelFull}(\tau_{full})$: expert labels the agent full trajectory ignoring hierarchical structure, giving labelled dataset of tuples $\langle s, a, \rangle$
- **Low level Inspection** $\text{InspectLow}(\tau_{low}; g)$: expert just verifies whether subgoal g was accomplished in the trajectory
- **Full Inspection** $\text{InspectFull}(\tau_{full})$: expert verifies whether agent’s overall goal was accomplished

Inspect ops are *lazy* version of label op, requiring less effort (easier to say whether or not robot is in elevator than providing an actual path to the elevator)

idea: if hierarchical trajectory agrees with expert on high level, and low level trajectories pass inspection, the the resulting full trajectory must also pass full inspection (hierarchical policy need not always agree with expert execution at low level to succeed in the overall task)

2.3 Dagger Algorithm Background

1

Dagger (Dataset aggregation) is an iterative reinforcement learning algorithm in the imitation learning space, that aims to trains a deterministic policy that achieves good performance in an online fashion (contrary to passive imitation learning which does rather poorly)

”DAGGER proceeds by collecting a dataset at each iteration under the current policy and trains the next policy under the aggregate of all collected datasets. The intuition behind this algorithm is that over the iterations, we are building up the set of inputs that the learned policy is likely to encounter during its execution based on previous experience (training iterations). This algorithm can be interpreted as a Follow-The-Leader algorithm in that at iteration n we pick the best policy π_{n+1} in hindsight”

3 Hierarchical Guidance

Hierarchical Guidance is paradigm where feedback from high-level expert guides low level learner in 2 ways:

1. high level expert ensures that low level expert is only queried when necessary (when subtasks have not yet been mastered)
2. low level learning is limited to relevant parts of the state space

¹Dagger Algorithm is used as a subroutine for the algorithms in the current paper.

3.1 Hierarchical Behaviour Cloning

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  $\mu \leftarrow \text{Train}(\mu, \mathcal{D}_{\text{HI}})$ 

```

- expert provides hierarchical demonstrations, with both low and high level trajectories, getting dataset for supervised learning which we can optimize via sgd
- then we train subpolicies predicting action from state, and meta controller predicting goal from state
- here hierarchical guidance is automatic since subpolicy demonstrations only occur in relevant parts of the state space
- but this is passive learning, we iterate on these ideas in the online learning context next

3.2 Hierarchically Guided Dagger

Algorithm 2 Hierarchically Guided DAGger (**hg-DAGger**)

```

1: Initialize data buffers  $\mathcal{D}_{\text{HI}} \leftarrow \emptyset$  and  $\mathcal{D}_g \leftarrow \emptyset, g \in \mathcal{G}$ 
2: Run Hierarchical Behavioral Cloning (Algorithm 1)
   up to  $t = T_{\text{warm-start}}$ 
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  $\mu \leftarrow \text{Train}(\mu, \mathcal{D}_{\text{HI}})$ 

```

- start by pretraining using cloning alg 1
- learner executes hierarchical policy, including choosing a subgoal, executing low level trajectory, and terminating according to end signal
- expert only gives feedback if the agent fails entire task
- if failed, expert first gets the subgoals that should have been taken by using **LabelHigh**
- only perform labelling on low level if meta controller chooses correct subgoal, but subpolicy fails
- since all preceding subgoals were chosen and executed correctly, and current subgoal is also correct, low level learning is in the relevant part of state space, and since subgoal failed, this is where we need to learn

3.3 Hierarchically Guided Imitation Learning

Algorithm 3 Hierarchically Guided DAgger / Q -learning
(hg-DAgger/ Q)

input Function $\text{pseudo}(s; g)$ providing the pseudo-reward
input Predicate $\text{terminal}(s; g)$ indicating the termination of g
input Annealed exploration probabilities $\epsilon_g > 0, g \in \mathcal{G}$
 1: Initialize data buffers $\mathcal{D}_{\text{HI}} \leftarrow \emptyset$ and $\mathcal{D}_g \leftarrow \emptyset, g \in \mathcal{G}$
 2: Initialize subgoal Q -functions $Q_g, g \in \mathcal{G}$
 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: $a \leftarrow \epsilon_g\text{-greedy}(Q_g, s)$
 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 τ and update $s \leftarrow \tilde{s}$
 13: **until** $\text{terminal}(s; g)$
 14: Append (s_{HI}, g, τ) to σ
 15: **until** end of episode
 16: Extract τ_{FULL} and τ_{HI} from σ
 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 $\mu \leftarrow \text{Train}(\mu, \mathcal{D}_{\text{HI}})$

- high level expert provides decomposition and pseudo-reward for each subgoal and picks a new subgoal at each step
- *hierarchical guidance impacts how the low level learner accumulates experience:*
 - as long as meta controller subgoal g agrees with expert, agent's experience is added to buffer
 - if meta controller selects wrong subgoal, no experience is there, ensuring buffer contains only relevant data
 - requires defining a $\text{success}(s, g)$ signal determining if subgoal g is complete and pseudo reward indicator

4 Results

4.1 Cost Savings

Theorem 1. *Given finite classes \mathcal{M} and Π_{LO} and realizable expert policies, the total cost incurred by the expert in hg-Dagger by round T is bounded by*

$$TC_{\text{FULL}}^I + (\log_2 |\mathcal{M}| + |\mathcal{G}_{\text{opt}}| \log_2 |\Pi_{\text{LO}}|)(C_{\text{HI}}^L + H_{\text{HI}} C_{\text{LO}}^I) + (|\mathcal{G}_{\text{opt}}| \log_2 |\Pi_{\text{LO}}|) C_{\text{LO}}^L, \quad (1)$$

Theorem 2. *Given the full policy class $\Pi_{\text{FULL}} = \{(\mu, \{\pi_g\}_{g \in \mathcal{G}}) : \mu \in \mathcal{M}, \pi_g \in \Pi_{\text{LO}}\}$ and a realizable expert policy, the total cost incurred by the expert in flat DAgger by round T is bounded by*

$$TC_{\text{FULL}}^I + (\log_2 |\mathcal{M}| + |\mathcal{G}| \log_2 |\Pi_{\text{LO}}|) C_{\text{FULL}}^L. \quad (2)$$

$$\propto \frac{H_{\text{HI}} + H_{\text{LO}}}{H_{\text{HI}} H_{\text{LO}}}$$

Therefore the cost ratio is

4.2 Maze

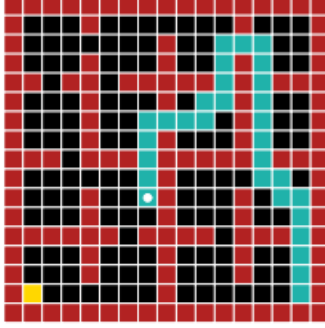
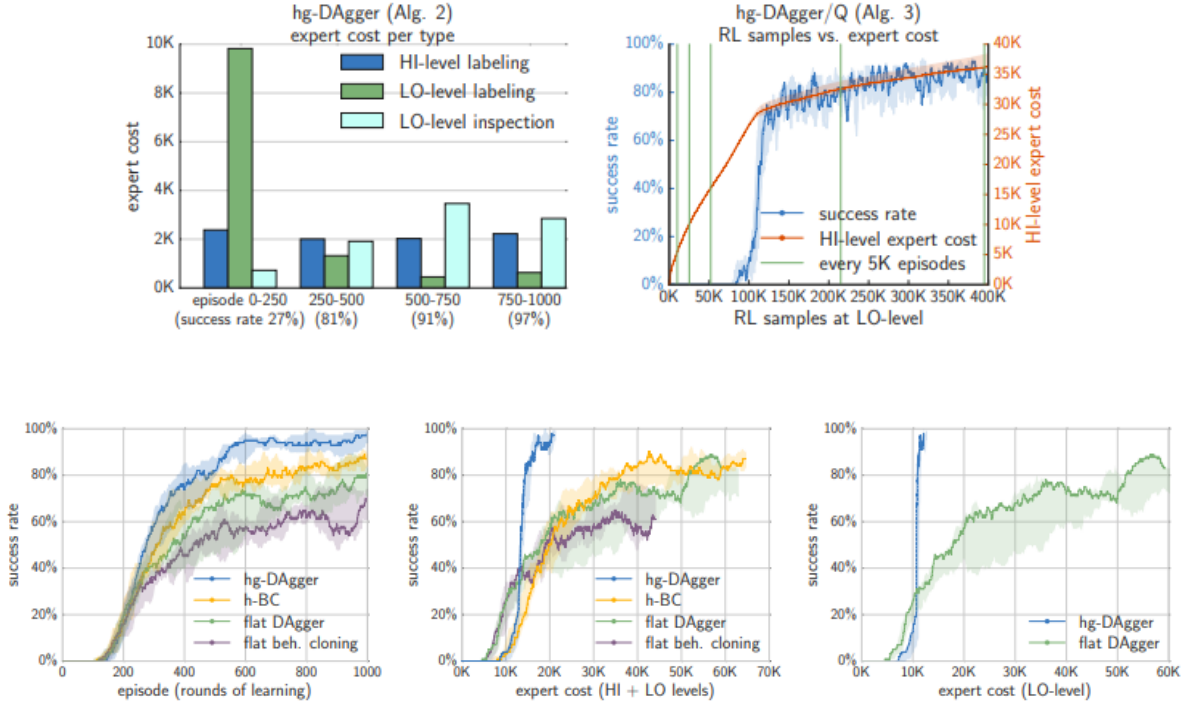


Table 1. Network Architecture—Maze Domain

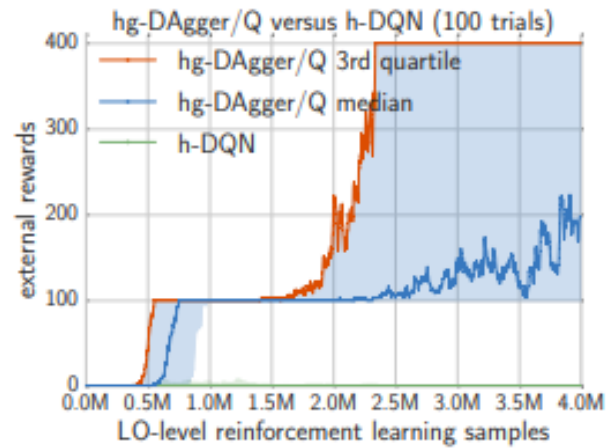
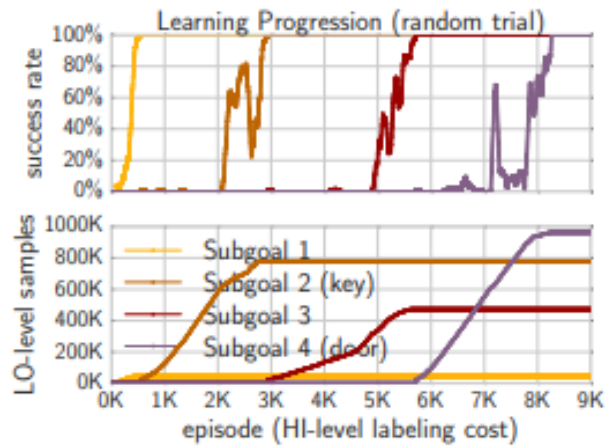
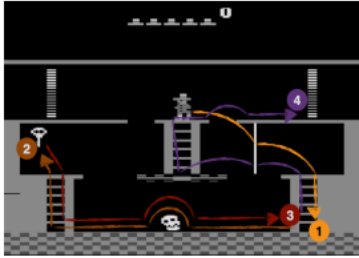
1: Convolutional Layer	32 filters, kernel size 3, stride 1
2: Convolutional Layer	32 filters, kernel size 3, stride 1
3: Max Pooling Layer	pool size 2
4: Convolutional Layer	64 filters, kernel size 3, stride 1
5: Convolutional Layer	64 filters, kernel size 3, stride 1
6: Max Pooling Layer	pool size 2
7: Fully Connected Layer	256 nodes, relu activation
8: Output Layer	softmax activation (dimension 4 for subpolicy, dimension 5 for meta-controller)

- Primitive Actinos: left, right, up, down
- Meta Actions: go to north, south, west, east room, go to target



- Hierarchical Version Outperforms Flat
- Using guidance saves alot of agent cost due to hierarchy
- number of High level labels rapidly icnreases initially and then flattens out after the learner becomes more successful

4.3 Montezuma Revenge



- The algo is able to get non zero rewards on this challenging first room task
- *h-sqn* may be picking the wrong subgoals, and hence getting noisy experience in the replay buffer, making it difficult to convergence