

1 We want to thank all reviewers for their time and feedback.

2 **To Reviewer #1:**

3 \* “Do you believe that your method could nonetheless be applied to more realistic settings?” **A1:** Yes. Our method is  
4 inspired by the influential empirical paper by [Tang et al. 2017] as we mentioned in Line 188-190.

5 \* “Is it fair/accurate to say that you are efficiently exploring?...” **A2:** Our method does guide exploration. Our  
6 framework relies on a provably efficient tabular RL algorithm which has to have an efficient exploration component,  
7 such as UCB or Thompson sampling. We want to emphasize that a tabular algorithm that uses random policy, e.g.,  
8  $\epsilon$ -greedy, *does not* have the desired properties we need because it requires an exponential number of samples. See the  
9 lower bound in Section A of [Jin et al. 2018]. We will add more clarifications on this point.

10 \* “Why the reward function depends on the time step?” **A3:** This setting is more general than the case where the  
11 reward function is independent of the time step. This setting has also been studied in many previous works. We will add  
12 more clarifications.

13 \* “Does your method support learning if  $H$  is unknown?” **A4:** Yes, our algorithm still works as long as we know an  
14 upper bound of  $H$ . We will add discussions about this. Thanks for pointing it out.

15 \* “Are levels the same as time steps within the episode?” **A5:** Yes.

16 \* “...one embedding for each level?” **A6:** In our description of the framework, we need one decoding function per  
17 level. But this is mainly for the convenience of analysis and is not necessary. Observation data can be mixed up among  
18 levels and used to train ONE decoding function for all levels. This is what we implemented in the numerical test. We  
19 will add more clarification about this part in the next version. Thank you for asking.

20 \* “Why run  $N$  times?” **A7:** If we directly require every output policy (Line 11, Algorithm 1) to be near-optimal with  
21 probability  $1 - \delta$ , the sample complexity will be  $\mathcal{O}(1/\delta)$  instead of  $\mathcal{O}(\log(1/\delta))$ . To avoid this, we only require each  
22 output policy to be near-optimal with a parameter-free probability e.g., 0.99 and then we run  $N = \mathcal{O}(\log(1/\delta))$  times  
23 to select the best policy out of  $N$  and achieve the final goal.

24 \* “...continuous low-dimensional representations...” **A8:** We are using clustering methods on the low-dimensional  
25 embedding space to reduce it to a finite-state space. This is also the implementation in [Tang et al. 2017].

26 **To Reviewer #3:** Thank you for the positive review.

27 **To Reviewer #4:**

28 \* “...the framework doesn’t consider the probability that the ULO encodes inappropriately the observations.” **A1:** In  
29 our definition of ULO (Definition 1), we do consider the probability that ULO encodes wrongly, which is described by  
30 the term  $g(n, \delta)$ . Besides, in the definition, the correctness of a function  $\hat{f}$  is on-average, i.e., if for some state  $s$ ,  $\mu(s)$  is  
31 close to 0,  $\hat{f}$  can be very inaccurate on  $s$ .

32 \* “...the training of ULO is not explained in detail.” **A2:** The training of ULO is described in Line 3-12 Algorithm  
33 2. We also explained in Line 218-232. Due to the page limit, we defer some parts in the appendix. We use ULO as a  
34 black-box oracle for abstraction and generality. We discuss some specific examples in Section B. We will add some  
35 details to the main text in the final version. Thanks for asking.

36 \* “...examine and/or discuss the applicability of the framework in more challenging environments.” **A3:** We have more  
37 discussion in the appendix about, for example, the choice of ULO one can consider solving real-world problems. We  
38 will consider this in the next version.

39 \* “...the label matching process and parts of Alg.2 should be explained in a more clear way. For instance, what is  
40 the purpose of the training data?” **A4:** Thanks for the advice. Due to the page limit, we defer some descriptions to the  
41 appendix. We will improve this part in the next version. For your question, the training data is collected to feed into the  
42 ULO and generate a more refined decoding function for each iteration.

43 **To Reviewer #5:**

44 \* “...some curves lack error bars and clarification in the caption.” **A1:** Thanks for mentioning this. We double-checked  
45 the graphs and we are sure that no curve lacks an error bar in the current version. They might be too tiny to observe.  
46 Yes, it is 1 standard deviation. We will add more clarification in the caption in the next version.

47 \* “There are some small typos.” **A2:** Thank you. We will double-check and revise them in the next version.

48 \* “More can be added such as the difference between their setting and real-world settings.” **A3:** Very good advice. We  
49 will consider this in the next version.