

1 Dear all reviewers, we highly appreciate your valuable comments and will reflect your comments in the revision.

2 » Reviewer 1

3 » “1. I think more could be done on the experimental front. ... in better RL in a number of settings.”

4 In all the environments except for Walker2D-disc, OC3 and SoftRobust successfully acquired options corresponding to
5 decomposed skills required for solving the tasks. Especially, OC3 produces robust options. In HopperIceBlock-disc,
6 for example, OC3 produces options corresponding to walking on slippery grounds and jumping onto a box (Figures
7 1 and 2). In HalfCheetah-disc, OC3 produces an option for running (highlighted in green in Figure 3) and an option
8 for stabilizing the cheetah-bot’s body (highlighted in red in Figure 3), which is used mainly in the rare-case model
9 parameter setups. Acquisition of such decomposed robust skills is useful when transferring to a different domain, and
10 they can be also used for post-hoc human analysis and maintenance, which is an important advantage of option-based
RL over flat-policy RL.

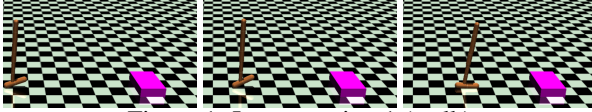


Figure 1: Learnt option 1 (walk)

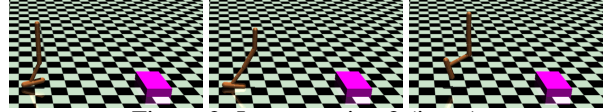


Figure 2: Learnt option 2 (jump)

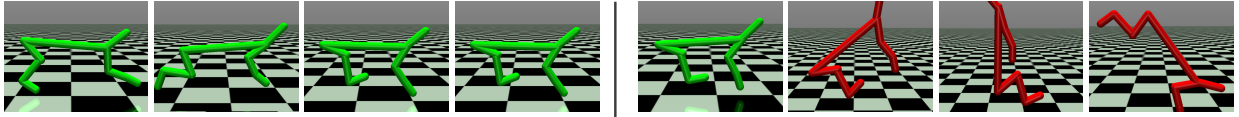


Figure 3: Learnt options (shown in green when option 1 is taking the control and in red otherwise) in Halfcheetah-cont. Left: the case with a normative model parameter setup. Right: the case with an abnormal model parameters setup.

11 Left: the case with a normative model parameter setup. Right: the case with an abnormal model parameters setup.

12 » “2. The soft robust loss baseline is not that bad ... to demonstrate the benefits of the proposed framework.”

13 We have conducted additional experiments with the condition that the model parameter distribution and the parameter
14 value ranges in the test phase are more uncertain (in terms of entropy) than those in the training phase. The results
15 confirmed that OC3 outperforms SoftRobust with respect to multiple performance measures (average return/loss and
16 CVaR) (the results will be reported in the revised version). The sim2real experiment couldn’t be done in this time frame.

17 » Reviewer 2

18 » “However, correct me if I misunderstood, ... it is unclear to me what the originality is of the proposed algorithm.”

19 Let us clarify the difference between our algorithm and Algorithm 2 in [6]¹. In addition to the type of parameters you
20 pointed out, Algorithm 2 in [6] and our algorithm are primarily different in the form of augmented MDPs. In Algorithm
21 2 in [6], the augmented MDPs take the form of $\langle S', A, R', \gamma, T', \mathbb{P}'_0 \rangle$, while, in our algorithm, the augmented MDPs
22 take the form of $\langle S', A, R', \gamma, \mathbb{E}_p [T'_p], \mathbb{P}'_0 \rangle$. Here T' is the transition function without model parameters, and the other
23 elements are the same as those in our paper. In sum, the types of transition function are different.

24 This difference comes from the difference of optimization objectives. In [6], the optimization objective is the expected
25 loss without a model parameter uncertainty ($\mathbb{E}_{\mathcal{R}'} [\mathcal{R}']$), while in our optimization objective, the model parameter uncer-
26 tainty is introduced by taking the expectation of the model parameter distribution ($\mathbb{E}_{\mathcal{R}', p} [\mathcal{R}'] = \sum_p \mathbb{P}(p) \mathbb{E}_{\mathcal{R}'} [\mathcal{R}' | p]$).

27 You may be concerned that modifying the augmented MDPs in [6] to use $\mathbb{E}_p [T'_p]$ for dealing with the parameter
28 uncertainty is straightforward. It turned out, however, that this seemingly simple modification required five pages of
29 proofs for justification, and any theoretical discussion for this has not been provided in previous work (e.g., in [6], [32],
30 and [1]). This is the reason why we described theoretical discussions (i.e., Eq. 11 ~ 14 based on our option policy
31 gradient theorems in Appendix) for this modification, which should help many researchers in the field.

32 » Reviewer 3 : Thank you very much for your comments. We will modify our paper as you pointed out.

33 » Reviewer 4

34 » “1) While I understand the paper is proposing ... is only to extend this CVaR MDP framework to include options.”

35 In the derivation of our CVaR option critic, the optimization of the soft robust loss $\mathbb{E}_{\mathcal{R}', p} [\mathcal{R}']$ in augmented MDPs is
36 necessary (please see our answer to Reviewer 2’s Q). So, as a part of the derivation, we also propose a new option critic
37 architecture to deal with this loss. Note that, due to the model parameter uncertainty (i.e., soft robust loss $\mathbb{E}_{\mathcal{R}', p} [\mathcal{R}']$),
38 the standard option-critic framework (the Bacon’s one) w. Chow’s CVaR PG cannot be applied to our setting. Both of
39 the previous frameworks (and Tamer’s framework) cannot deal with the case with model parameter uncertainty.

40 » “3) While the experimental ... the vanilla CVaR PG algorithm compare with the option-critic counterpart?”

41 OC3 produces some useful portable options even in Halfcheetah and Walker2D-cont (see our answer to Reviewer 1’s
42 Q). Note that, since Chow’s framework [6] (and Tamer’s framework) cannot deal with model parameter uncertainty,
43 these cannot be applied to our setting. Finally, we will modify the notations as you suggested.

¹We assume that you mean Algorithm "2" (Actor-Critic Algorithms for CVaR Optimization), not Algorithm "1", in [6].