

1 We thank the reviewers for their helpful suggestions, and will do our best to incorporate them into our paper. Overall,
2 we want to emphasize that the goal of SPECTRL is to make it easier to apply RL to tasks with complex objectives. In
3 particular, SPECTRL enables the user *program* what the agent needs to do at a high level; then, SPECTRL automatically
4 learns a policy that tries to best satisfy the user intent.

5 **R1, LSTMs:** It is true that LSTMs can be used to solve RL problems with non-Markovian specifications. However,
6 the task monitor introduced in this paper not only eliminates the need to learn the internal state, but also enables us to
7 perform reward shaping. As we show in our experiments, reward shaping is crucial for learning complex tasks. Thus,
8 even if the LSTM learns a perfect encoding of the state, our approach would still substantially outperform it.

9 **R2, Curiosity-driven exploration:** These approaches aren't really applicable to addressing our challenges, which
10 involve complex objectives rather than hard exploration. For example, the task may be for the agent to visit a sequence
11 of goals in a particular order. The agent may know how to reach each of these goals individually, but this is not sufficient
12 for the agent to complete the task.

13 **R2, Sub-goal based rewards.** Many challenges arise when considering the details of how sub-goal based rewards
14 would be implemented. For example, how does achieving a sub-goal count compared to violating a constraint? How do
15 we handle sub-goals that can be achieved in multiple ways? How do we ensure the agent does not repeatedly obtain a
16 reward for a sub-goal that it has already completed? As tasks get more complex/deeply nested, manually determining
17 rewards for sub-goals is very non-trivial. These challenges are exactly what our system is designed to solve.

18 **R2, Baselines:** We used ARS since it has been demonstrated to be state-of-the-art for continuous control. For fair
19 comparison, we use ARS with TLTL as well. Also, like TLTL, our approach can be used with with other algorithms.

20 **R2, More challenging environments:** Our experiments already show that our algorithm can solve problems that the
21 state-of-the-art ARS algorithm cannot. Benchmarks such as MuJoCo may involve more complex dynamics, but they
22 focus largely on short-term control tasks—e.g., running in a straight line as fast as possible. In contrast, the benefit of
23 our approach is on handling complex, long-term control tasks, where the agent must perform a variety of actions to
24 achieve complex objectives—e.g., reaching multiple goals in sequence and then coming back. While there has been
25 some success solving these problems using traditional RL, these approaches typically rely on enormous amounts of
26 computation. In contrast, our algorithm converges quickly without a large amount of computation (all our experiments
27 are run on a single GPU).

28 **R2, Cart-Pole:** Our cart-pole benchmark differs from the standard OpenAI one in two ways. (i) The goal in the
29 standard one is just to keep the pole balanced for as long as possible. In ours, the goal is to move the cart to the right
30 and then back to the starting position without letting the pole fall. We find this task to be substantially more challenging
31 than the standard one. (ii) The standard benchmark has discrete actions, whereas ours has continuous actions, which
32 also makes the task more challenging.

33 **R2, TLTL:** We have two contributions over TLTL. (i) While it is possible to extend the state space to handle non-
34 Markovian TLTL specifications, this could not previously be done automatically. One of our contributions is to provide
35 a way for automatically extending the state space; our algorithm can be modified to work for TLTL formulas as well.
36 (ii) We perform reward shaping, whereas TLTL does not. We find that reward shaping is critical for good performance.

37 **R3, Reward machines:** We differ from the paper on reward machines in two ways. First, in their paper, the specification
38 is given directly as a finite state machine along with reward functions for each state. Their key contribution is in the
39 ability to learn multiple tasks simultaneously by applying the Q-learning updates across different specifications. In
40 contrast, our contribution is to automatically generate a task monitor from a high-level logical specification. Second,
41 our notion of task monitor is more expressive since it has a finite set of registers that can store real values. In contrast,
42 finite state machines cannot store quantitative information. We will add a discussion to our related work.

43 **R3, barriers:** Some of our benchmarks include barriers (i.e., obstacles). Our approach is able to solve these tasks.

44 **Presentation:** To improve readability, we plan to add pictures describing the high level approach in a final version. We
45 will include error-bar plots for Figure 3, which we show below. To exclude outliers, we omitted one best and one worst
46 run out of the 5 runs. While the variance is sometimes high, our insights continue to hold—in particular, our algorithm
47 consistently outperforms the baselines. We will add a conclusion as well.

