1 We thank the reviewers for insightful comments and suggestions. We hope the rebuttal will clarify all the doubts.

2 **Motivation and real-world application** In the LfL [1,2] setting, we do not ignore the near-optimal trajectories,
3 but we can use all the training data to recover the reward function. In some cases, such as in multi-agent IRL,
4 the agent is interested in learning the rewards of other agents before the other agents become experts, to adopt
5 strategic behavior. Also in autonomous driving, when there is a new circuit, instead of waiting for the driver to
6 become an expert, we can learn his intentions in the initial learning phase. In the appendix, we reported an experi-
7 ment that uses the trajectories collected from an autonomous driving simulator (used in real-world problems [3]) to
8 recover the reward function. We think this experiment may show a more **realistic application** of the LfL problem.

9 **Gridworld** In the gridworld experiment, we show LOGEL's performance when the *gradient-*
10 *based learner assumption* is *violated*. We experiment with three no-gradient learners (Value
11 Iteration, Q-Learning, and Soft Policy Iteration) and one gradient learner (GPOMDP). The
12 results show that LOGEL works well even when the hypothesis is violated. In Figure 2) (10
13 runs) we show the performance of LOGEL using only the **online samples** of Q-Learning.

14 **MuJoCo** In these experiments, the learner uses the original reward function, but we construct
15 the reward features for LOGEL from the state and action features. In the Hopper environment,
16 we add the alive feature. The LfL algorithm uses, as in the original paper, a reward function
17 which is very informative as they model the reward with a policy that has the same structure
18 of the learner policy. In Figure 3) (6 runs) we show the behavior of LOGEL in the Hopper
19 environment without the alive feature, using only the l2 distance between the current state
20 and the next state and the norm, the squared norm, and the cubed norm of the action. The
21 good behavior of the learned policies with the recovered reward weights is because the optimal
22 reward function to learn a task can be different from the real one [5].



1)



2)



3)

23 **R1** We answered questions 1 and 2 above. In the MuJoCo experiments, we use a different
24 reward space than the learner one, so we cannot show convergence to the original reward
25 weights. As suggested, we report another experiment with different reward features, Fig.3.
26 Other questions: L138 the noise is due to the estimation of the gradient from trajectories; the
27 intrinsic bias refers to the bias introduced by the estimated gradient that cannot be reduced
28 with more learning steps; $\alpha_r > \epsilon$ because the learning rate is by assumption greater than 0;
29 in gridworld experiment we used four learners: SVI, SPI, Q-learning and G(PO)MDP, we
30 collected their trajectories and we estimate the reward function with LOGEL; L254 at the end of
31 (paper) Fig.1 the weight difference is zero; L286 the Reacher and the Hopper experiments use
32 the original reward functions for the learner, but different reward features for LOGEL created
33 from the state and action information. In the appendix we compared LOGEL with TREX.

34 **R2** We can think that in the MuJoCo environment the linearity assumption is violated as we use
35 a different reward space for the recovered reward function. Behavioral cloning, as in LfL, is a
36 requirement for running our algorithm. However, we think that by using the behavioral cloning
37 of the previous iteration as a starting point, we can learn a good approximation of the current
38 policy after a few learning steps. In the MuJoCo experiment, and the Gridworld experiment with the GPOMDP learner,
39 we only use a subset of the demos that the learner uses to update her policy and we can assume that in a real-world
40 scenario we have access to such demos. Thank you for pointing out another related work, we will add the citation.

41 **R3** In the gridworld experiment, we use a value-based and a Q-Learning learner that violate the gradient-based
42 assumption. In Figure 2) we show the behavior of LOGEL using only the online updates of a QLearning learner. In
43 Figure 1) (10 runs), we compare the performance of LOGEL with GIRL [4], which is an IRL batch model-free algorithm.
44 The experiment shows that the algorithm cannot recover the correct reward weights from suboptimal trajectories. We
45 answer to the Hopper question in the MuJoCo section above. $\Psi$ is bounded in l2-norm, $\lambda_{min}$ is the smallest eigenvalue,
46 the confidence interval is computed with t-distribution.

47 **R4** In the Hopper domain, as in the original paper [1], the LfL baseline does not recover a correct reward function; this
48 behavior is due to the starting behavior of the simulated robot which at first often falls to the ground; the LfL algorithm
49 tends to consider these absorbing states as good ones. The domain tasks in which LOGEL outperforms LfL are tasks
50 in which monotonous improvement is violated, in which there are absorbing states and where the learner is not a soft
51 policy improvement algorithm (as the gridworld experiment shows). In the behavioral cloning, we *do not assume the*
52 *knowledge of the policy* model of the learner: for example, in the gridworld environment, the Q-Learning and the SVI
53 learner are $\epsilon$-greedy actors, but we model their policies as Boltzman policies; also, in the appendix, in the autonomous
54 driving experiment, the learner has a rule-based policy, while we use a linear layer to approximate the policy. In Figure
55 1) we show how a batch-model free IRL algorithm cannot learn using only learning trajectories. We think that the
56 merits of LOGEL go beyond removing the monotonically improving assumption: we provide theoretical insights into
57 the correctness of the proposed method and we empirically show that the proposed method works well also with learner
58 using value-based methods (as value-iteration) and soft policy improvement algorithms, outperforming in most of the
59 experiments the LfL baseline.

60 [1] Jacq, A., Geist, M., Paiva, A., Pietquin, O. (2019). Learning from a Learner. ICML
61 [2] Kubala, V., Konidaris, G., Greenwald, A. (2019). Inverse Reinforcement Learning from a Learning Agent, RLDM
62 [3] Likmeta, A., Metelli, A. M., Tirinzoni, A., Giol, R., Restelli, M., Romano D., (2020) Combining reinforcement learning with rule-based controllers for transparent and general decision-making
63 in autonomous driving, Robotics and Autonomous Systems
64 [4] Pirotta, M. and Restelli, M. (2016). Inverse Reinforcement Learning through policy gradient minimization. AAAI
65 [5] Ng et al. "Policy invariance under reward transformations: Theory and application to reward shaping." ICML 1999