1  We thank the reviewers for their valuable suggestions. Please find our answers for each reviewer below.

2  **Reviewer 1**

3  We thank the reviewer for the positive assessment of our work. Below, we provide a concrete plan of incorporating
4  reviewer's feedback in the updated version of the paper.

5  **Extended experimental analysis.** As suggested by the reviewer, we will add a more detailed analysis about the
6  experimental results in the paper. In particular, we will add the following experiments/details/results: *(i)* evaluation
7  of learner-aware teaching under unknown constraints for L3-L5 (the findings are similar as for the already presented
8  experiments); *(ii)* experiments illustrating the effect of $C_r$ and $C_c$ in soft preference constraints; *(iii)* additional details
9  and discussion of parameter choices in our experiments; *(iv)* reporting the run time of our algorithms, and illustrating
10  scalability w.r.t. the problem size; and *(v)* reporting standard errors in Figure 3 (b) (the currently reported results in the
11  paper are significant at significance level $0.1$).

12  **Ideas for outlook and future work.** To the best of our knowledge this paper is the first to consider IRL with preference
13  constraints. Hence, we primarily focused on developing the theoretical framework and algorithms (for both the known
14  and unknown constraint settings). Nevertheless, we agree that the directions suggested by the reviewer (more complex
15  domains and human subject experiments; suboptimal demonstrations and implications on performance; addressing the
16  problem from a learner's perspective) are important. We will add a discussion on these directions in the revised paper.

17  **Technical clarifications.** Below we answer the technical questions raised by the reviewer.

18  • The values $C_r$ and $C_c$ describe a learner's relative importance to mimic the teacher's demonstrations and following
19    its own preferences, respectively, and are thus properties of a learner and not the parameters of a teacher.
20  • The performance of AWARE-BIL decreases for increasing learner's constraints because the learner's preferences to
21    avoid certain cells conflicts with the goal to go to certain cells to accumulate rewards. Note that this decrease is due
22    to the experimental setup and not due to limitations of AWARE-BIL.
23  • $\delta_r^{\text{hard}}$ and $\delta_r^{\text{soft}}$ are used to characterize a learner's reward feature matching behaviour as part of the learner's opti-
24    mization objective: While a mismatch of up to $\delta_r^{\text{hard}}$ between the learner's and teacher's reward feature expectations
25    incurs no cost regarding the optimization objective, a mismatch larger than $\delta_r^{\text{hard}}$ incurs a cost of $C_r \cdot \|\delta_r^{\text{soft}}\|_p$. Please
26    also note that $\delta_r^{\text{hard}}$ is a fixed parameter, while $\delta_r^{\text{soft}}$ is an optimization variable. In equation 1, $m$ is the number of
27    preference constraints of the learner. In general, $m \neq d_c$. Note that we have a typo in the paper in line 108 which
28    might have caused some confusion: we incorrectly wrote $\delta_c^{\text{soft}} \in \mathbb{R}^{d_c}$ but we wanted to say $\delta_c^{\text{soft}} \in \mathbb{R}^m$. We will
29    correct this typo and elaborate on the notation in the revised paper.
30  • $\delta_r^{\text{soft,low}}$ and $\delta_r^{\text{soft,up}}$ are auxiliary variables used to rewrite the constraints on the absolute value of the mismatch in a
31    form more convenient for optimization. We will add clarification and more details to the revised paper.

32  **Reviewer 2**

33  We thank the reviewer for providing useful suggestions and high-level comments on the paper structure.

34  As suggested, we will remove the auto-pilot example from the introduction and elaborate more on the other two
35  examples. We will also emphasize that the learner-aware teacher with full-knowledge of the learner allows us to
36  formalize the problem and introduce a theoretical/algorithmic framework to study the limitations of learner-agnostic
37  teaching. The real use-case of learner-aware teaching is for incomplete knowledge of the learner. We believe that in this
38  paper we consider an important new direction for inverse reinforcement learning which we would like to make available
39  to the community in a timely manner by a conference publication. However, we will revise the paper to include more
40  details on the algorithms in Section 5.

41  **Reviewer 3**

42  We thank the reviewer for appreciating the novelty of the problem setting and providing suggestions for improvements.

43  **Regarding linearity of the reward function.** It is true that our results are currently for the linear setting. However,
44  we believe that it is worthwhile to first thoroughly understand this setting. Moreover, as we don't constrain the feature
45  maps $\phi_r$ and $\phi_c$, the features we consider can be nonlinear functions of a set of "basic" features, which in principle
46  makes it possible to accommodate quite general situations in our setting. Nevertheless, we agree that a natural next step
47  is to investigate to what extent our ideas can be extended to nonlinear reward settings.

48  **Regarding experimental evaluation on more realistic tasks.** Generally, we agree with the reviewer's suggestions
49  and believe that evaluating our algorithms on more realistic tasks is a natural direction for future work. We would
50  like to reemphasize that the paper's primary focus is on introducing an important problem setting for IRL, developing
51  algorithms for the problem, and empirically understanding the performance of these algorithms. We will further extend
52  the experimental analysis in the paper as outlined in our response to Reviewer 1.