

1 We thank all reviewers for providing insightful comments and helpful suggestions. We have revised some parts of our
2 paper accordingly to improve the clarity. The following are our responses to some specific topics.

3 **Regarding adding a PAC-style analysis for VIEC using learned models (Reviewer 1).** Including such an analysis
4 is indeed good for the sake of completeness, but we feel that it also misses the main points of this paper: (1) an
5 exploration strategy with PAC guarantees (R-MAX, MBIE, etc.) can still be far from optimal in terms of exploration
6 cost; (2) explicit planning has remarkable potentials in improving it. Such potentials can only be shown by comparing
7 the exploration costs of the best exploration scheme (the optimal one for the true MDP) and the ones actually taken by
8 the existing algorithms. PAC-style analysis stresses the worst cases, and thus cannot be used to show such potentials.
9 This is a little like guaranteeing that one can run 100m in at most 1 minute, but it does not say how (or even whether)
10 one can achieve 9.6s (which is the main concern of this paper).

11 It is also very difficult (if not impossible) to use a PAC-style analysis to expose the weaknesses of algorithms such as
12 distance / reward traps, because such an analysis usually occurs at an abstract level and does not show exactly when and
13 how wrong decisions are made during exploration.

14 Therefore, while the PAC-style analysis is an interesting suggestion and can surely bring some more insight, we do not
15 consider it essential for the purpose of this paper.

16 **Regarding the reward traps, exploration vs exploitation, and the Six Arms example (Reviewer 1).** Although
17 algorithms like R-MAX and MBIE are conceptually designed for trading-off exploration and exploitation, in reality
18 they achieve it by doing “pure” exploration first and then converging to some policy at certain stage. As Reviewer 1
19 mentioned in the Six Arms example, MBIE first explores and tries out different arms, then stops exploration and starts
20 exploiting the best arm after it has sufficient evidence for confirming the optimality of that arm with a high confidence.
21 This is a clear example of exploration-then-exploitation behaviour with exactly one phase change in the process.

22 As we have discussed in Section 3.1, the minimum sample size for confirming the optimality of a policy / action is given
23 by the Hoeffding’s or Chernoff’s inequalities, which implicitly result in a fixed demand matrix for a given MDP. Given
24 the same demand matrix, the optimal exploration scheme never chooses actions more than necessary (by definition),
25 while MBIE can be distracted by immediate rewards and eventually choose suboptimal actions more often than needed
26 (which we call reward traps). Therefore, no matter how the domain is designed, MBIE always needs more steps to
27 fulfil a demand matrix (i.e. finish its exploration phase) than the optimal exploration scheme (unless the demand matrix
28 happens to match the actual behaviour of MBIE exactly, in which case the costs are the same).

29 Even in the case where the total reward is of primary concern (and thus “balancing” exploration and exploitation is
30 important), being lured to the reward traps only prolongs the exploration phase of MBIE, which leads to strictly less
31 total reward in the long run. Therefore, what we called reward traps are still traps in such cases.

32 We understand that this part can be somewhat tricky (especially to readers accustomed to traditional concepts of
33 exploration). More discussion is added to the final version to help clarify our ideas.

34 **Regarding the computational cost (Reviewer 1).** We agree that the computational cost for VIEC is too high to be
35 of practical use. However, an efficient approximation is practical, since a large portion of the demand space is not
36 crucial in deciding the (near-)optimal scheme. We have a simple implementation of the Monte-Carlo approximation
37 that can achieve $O(nmd)$ in tower MDPs like the optimal scheme (albeit with a worse constant factor) which computes
38 much faster than VIEC. This, of course, is still not as fast as heuristic based algorithms and could benefit from further
39 improvement. This will be our future work.

40 **Regarding the lack of numerical experiments (Reviewer 2 & 3).** We agree that including a section of empirical
41 results might make the paper more complete. In fact, it is actually more difficult to visualise the exploration behaviours
42 through empirical results than through the analysis done in Section 4. If we only provide the exploration cost curves, it
43 will not provide more useful information than the theoretical results shown in Table 1 of Section 4. We consider that
44 providing more details about how existing algorithms behave during exploration and why some decisions are bad is
45 more helpful for readers to understand the weaknesses of those algorithms as well as the potential of explicit planning,
46 hence the focus on the theoretical analysis in this paper.

47 **Regarding the relationship with Bayesian RL (Reviewer 3).** The resemblance between the two is a little superficial.
48 Although both do planning in augmented MDPs, the objectives of planning, the augmented state spaces, and the
49 definitions of optimality are all different. In the context of this paper, Bayesian RL is essentially yet another family of
50 heuristic-based methods. It chooses actions that are (approximately) Bayesian optimal, which are not guaranteed to be
51 optimal with respect to exploration costs. In fact, since Bayesian methods take immediate rewards into consideration,
52 they are prone to reward traps just like the IE family elaborated in Section 4.5. The ones that use discounts in planning
53 are also prone to distance traps.