

1 We thank the reviewers for their comments.

2 We feel that the contributions of this work and the motivation for it are well understood by the reviewers: We investigate
3 the convergence properties of an online planning algorithm, given access to a lookahead policy oracle. Indeed, we
4 study the performance of such algorithms, both in the exact form and in several approximate settings, and compare
5 them to their approximate dynamic programming (ADP) counterparts. To the best of our knowledge, there is no
6 theoretical analysis regarding guarantees of lookahead policies in online planning algorithms. Furthermore, we believe
7 the generality of the presented techniques may be found useful in the analysis of and development of online planning
8 algorithms.

9 There is no dispute on the importance of empirical work. However, we believe that the theoretical results provided
10 in this work are important on their own. Our analysis spans not only the exact, but also three approximate settings,
11 and provide detailed comparison to the performance of ADP. Furthermore, unlike the scarcity of theoretical results in
12 online planning with lookahead policies, there are many works that study the empirical performance of different online
13 planning algorithms that are based on lookahead policy. Yet, the existing empirical works are heuristic; this stresses the
14 importance of theoretical results on online planning algorithms with lookahead policies. We hope the rigorous approach
15 pursued in our work will stimulate further theoretical research on the interplay between the lookahead horizon and the
16 performance of the online planning algorithm. There are, as always, interesting and important theoretical questions to
17 be solved.

18 As mentioned in response to Reviewer 3 and also clearly highlighted by Reviewer 4, a thorough empirical comparison
19 of RTDP, MCTS and ADP is very important and useful for the community, and may shed light on several unanswered
20 questions about the MCTS algorithm. However, due to the extent of our theoretical results and the length of current
21 paper, we feel a thorough empirical study of these algorithms is outside the scope of this work.

22 **R1.** For the comment on the needed empirical work, please see the above paragraphs that have been written for all
23 the reviewers. The second question that you raised is definitely of interest. In fact, in our opinion, answering such a
24 question deserves a work on its own, as there are probably several answers for it, such as which assumptions should be
25 made? What are the relevant structural properties? Can we choose it in an online manner? In a somewhat different
26 context, this question is equivalent to asking how to choose a hyperparameter of an algorithm? Indeed, it is an important
27 question, however it is outside the scope of the current work.

28 In the ICML-2020 paper "Multi-step Greedy Reinforcement Learning Algorithms" by Tomar, Efroni, and Ghavamzadeh,
29 the authors consider the framework of *model-free RL*, whereas we focus on online planning. For this reason, the
30 works are not very much related (we shall clarify this in the text, thanks!). Furthermore, this empirical paper is a follow
31 up to two theoretical papers ('Beyond the one step greedy approach in RL' ICML 18' and 'Multiple-Step Greedy
32 Policies in Approximate and Online RL' NeurIPS 19'). In our opinion, this clearly shows the importance of theoretical
33 results that can guide the design of practical algorithms.

34 **R2.** We would like to thank the reviewer for the supportive review. It definitely encourages us to keep investigating
35 online planning algorithms and expand our current understanding. We will revise the sentence on line 207.

36 Previous results which analyzed the performance of the UCT algorithm showed its worst-case sample complexity
37 depends exponentially (or even worst) on to the horizon of the problem. Furthermore, the sample complexity to find an
38 ϵ optimal action using the UCT usually scales as $O(1/\epsilon^2)$. Unlike the UCT, the sample complexity of RTDP depends
39 on the size of the state space (or the abstract state space as we showed in this work) and does not depend exponentially
40 on the horizon, only polynomially. The dependence on ϵ scales as $O(1/\epsilon)$. These results definitely implies on a possible
41 superiority of RTDP over MCTS from a theoretical perspective. We will emphasize it in the discussion part.

42 **R3.** We agree with the reviewer about the need for a thorough empirical work that compares RTDP, MCTS and
43 ADP. Such work is very important to the community in our opinion and might resolve the hardness of tuning the
44 hyper-parameters of the MCTS algorithm. Due to the extent of the theoretical results supplied in the paper, that
45 comprehensively study h-RTDP in its exact form and in three approximate settings, we feel a thorough empirical study
46 of these algorithms is outside the scope of the current paper.

47 We now address the additional feedback the reviewer gave. -) It is the same motivation as in RTDP. We will discuss this
48 in the final version of the paper. -) We tried to stress this as much as we can. We will add it to the abstract to avoid
49 confusions. -) Due to lack of space, we will add a table to the appendix to fully specify this computational complexity.
50 -) Thanks! we will fix this.

51 **R4.** We would like to thank the reviewer for the positive feedback. We also thank you for the minor comments which
52 helps us to improve the work. We will fix them all in the final version of the paper.