

1 We thank the reviewers for their comments. We will address all style suggestions and minor points.

2 **Explaining limitations of our approach:** The main comment in all 3 reviews is the suggestion to include more
3 discussion or experiments with small problems to illustrate the limitations of the approach, which we agree is a good
4 idea. First note that the main assumptions of SNAP are the independence assumption in rollouts and the fact that the
5 rollout plans do not depend on observations beyond the first step. This last restriction can be lifted at the cost of being
6 exponential in depth, as in other algorithms, but as experiments show speedup is crucial in large problems.

7 Deterministic transitions are not necessarily bad for factored representations because a belief focused on one state is
8 both deterministic and factored and this can be preserved by the transition function. Both the T-maze and the rocksample
9 domains that were proposed in the reviews are actually suitable for SNAP. The reason is that one step of observation is suf-
10 ficient and the reward does not depend in a sensitive manner on correlation among variables. The [Pajarinen] paper shows
11 that factoring works well for rocksample. We encoded the T-maze domain in RDDL and show the results in the table.

12 We use 3 seconds per step for each algorithm, and the T-maze length
13 is shown in the first row. YES means that the algorithm finds optimal
14 actions for all steps, NO means the algorithm does not, and X means
15 the result is not available as DESPOT only runs with pomdpX and the
16 RDDL to pomdpX translator failed for this problem size. The result shows that SNAP can solve the T-maze problems.

		5	10	12	15	30	50
SNAP	YES	YES	YES	YES	YES	YES	YES
DESPOT	YES	YES	YES	X	X	X	X
POMCP	YES	NO	NO	NO	NO	NO	NO

17 However, we can illustrate the tradeoff with two other simple domains. The first has 2 state variables x_1, x_2 , 3
18 action variables a_1, a_2, a_3 and one observation variable o_1 . The initial belief state is uniform over all 4 assignments
19 which when factored is $b_0 = (0.5, 0.5)$, i.e., $p(x_1 = 1) = 0.5$ and $p(x_2 = 1) = 0.5$. The reward is if $(x_1 ==$
20 $x_2)$ then 1 else -1 . The actions a_1, a_2 are deterministic where a_1 deterministically flips the value of x_1 , that is:
21 $x'_1 =$ if $(a_1 \wedge x_1)$ then 0 else if $(a_1 \wedge \bar{x}_1)$ then 1 else x_1 . Similarly, a_2 deterministically flips the value of x_2 . The
22 action a_3 gives a noisy observation testing if $x_1 == x_2$ as follows: $p(o = 1) =$ if $(a_3 \wedge x'_1 \wedge x'_2) \vee (a_3 \wedge \bar{x}'_1 \wedge$
23 $\bar{x}'_2)$ then 0.9 else if a_3 then 0.1 else 0. In this case, starting with $b_0 = (0.5, 0.5)$ it is obvious that the belief is not
24 changed with a_1, a_2 and calculating for a_3 we see that $p(x'_1 = 1 | o = 1) = \frac{0.5 \cdot 0.9 + 0.5 \cdot 0.1}{(0.5 \cdot 0.9 + 0.5 \cdot 0.1) + (0.5 \cdot 0.9 + 0.5 \cdot 0.1)} = 0.5$ so
25 the belief does not change. In other words we always have the same belief and same expected reward (which is zero).
26 Therefore, for this problem factoring implies that the search is blind. On the other hand, a particle based representation
27 of the belief state will be able to concentrate on the correct two particles (00,11 or 01,10) using the observations.

28 The second problem has the same state and action variables, same reward, and a_1, a_2 have the same dynamics. We have
29 two sensing actions a_3 and a_4 and two observation variables. Action a_3 gives a noisy observation of the value of x_1 as
30 follows: $p(o_1 = 1) =$ if $(a_3 \wedge x'_1)$ then 0.9 else if $(a_3 \wedge \bar{x}'_1)$ then 0.1 else 0. Action a_4 does the same w.r.t. x_2 . In this
31 case the observation from a_3 does change the belief, for example: $p(x'_1 = 1 | o_1 = 1) = \frac{0.5 \cdot 0.9}{0.5 \cdot 0.9 + 0.5 \cdot 0.1} = 0.9$. That is, if
32 we observe $o_1 = 1$ then the belief is $(0.9, 0.5)$. But the expected reward is still: $0.9 \cdot 0.5 + 0.1 \cdot 0.5 - 0.9 \cdot 0.5 - 0.1 \cdot 0.5 = 0$
33 so the new belief state is not distinguishable from the original one, *unless one uses additional sensing action a_4 to*
34 *identify the value of x_2* . In other words for this problem we must develop a search tree because one level of observations
35 does not suffice. If we were to develop such a tree we can reach belief states like $(0.9, 0.9)$ that identifies the correct
36 action and we can succeed despite factoring, but SNAP will fail because the search is limited to one level of observations.
37 Here too a particle based representation will succeed because it retains the correlation between x_1, x_2 .

38 **Rev1 - Heuristic domain knowledge:** we agree that the performance of MCTS will improve with domain specific
39 knowledge. However, our focus was on domain independent performance of the planners. In addition, since different
40 algorithms might use domain knowledge in a different manner the comparison of algorithms would be less clear.

41 **Rev1 - Prior work with factored belief:** Thanks! We will add a discussion of this and other work to the paper.

42 **Rev2 - Evaluating contributions from SOGBOFA:** We agree that it is important to understand the contribution of
43 different components. Contributions of components that are part of SOGBOFA were reported in the cited papers, albeit
44 in the context of MDPs. The experimental evaluation in this paper is centered on evaluating the the new ideas in this
45 paper, showing the importance of sampling for large observation spaces, and performance sensitivity w.r.t. run time per
46 step and number of samples.

47 **Rev2 and Rev3: line 96:** Yes $p(x = 1)$ should be $p(x = T)$. We will plan to improve the descrip-
48 tion. In the current notation, for a generic variable x assume $p(x = T)$ is given by the RDDL expression
49 if $(cond1)$ then p_1 else if $(cond2)$ then p_2 else p_3 where the conditions form a set of mutually exclusive and exhaustive
50 set of conjunctions. Then the probability of x given that $cond_i$ holds is p_i . In general, assuming discrete variables, the
51 CPT of the variable x can be rewritten in this form and this is facilitated by the RDDL representation.

52 **Rev3 - plan vs. policy:** We will plan to improve the description. By a plan we meant a pre-determined sequence of
53 actions not conditioned on states or observations. The same notion is also known as an open loop policy and as a straight
54 line plan. A policy will condition future action choices on the future states (in MDP) or belief states (in POMDP).