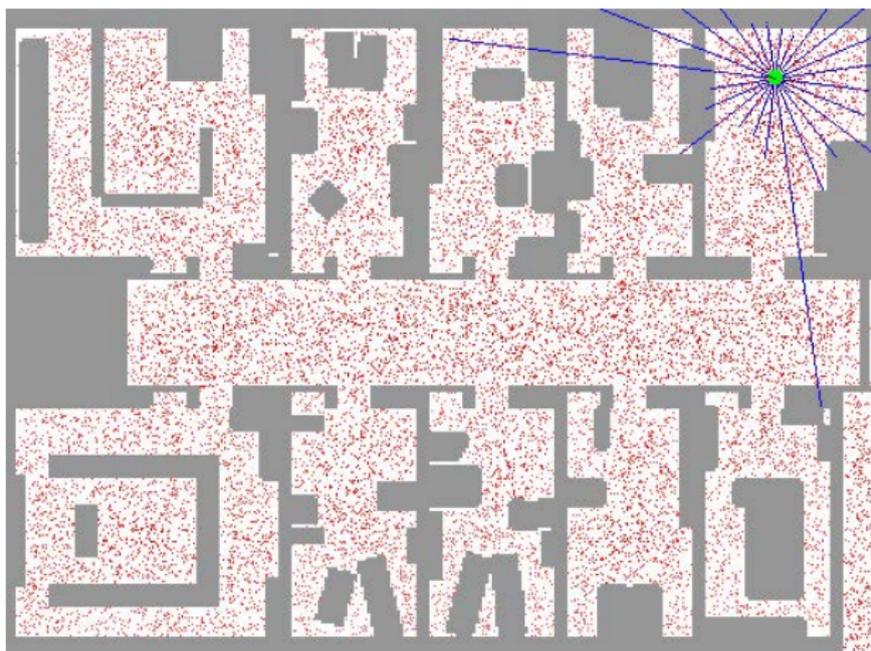


Information Gathering and Reward Exploitation of Subgoals for POMDPs

Hang Ma and Joelle Pineau
McGill University

AAAI January 27, 2015



http://www.cs.washington.edu/ai/Mobile_Robotics/mcl/animations/global-floor.gif

POMDPs

A partially observable Markov decision process (POMDP) is a tuple $\langle \mathcal{S}, \mathcal{A}, \Omega, T, O, R, b_0, \gamma \rangle$

- \mathcal{S} is a finite set of states, \mathcal{A} is a finite set of actions, Ω is a finite set of observations;
- $T(s, a, s') = p(s'|s, a)$ is the transition function that maps each state and action to a probability distribution over states;
- $O(s', a, o) = p(o|s', a)$ is the observation function that maps a state and an action to a probability distribution over possible observations;
- $R(s, a)$ is the reward function, $b_0(s)$ is the initial belief state, and γ is the discount factor.

Beliefs

A belief state $b \in \mathcal{B}$ is a sufficient statistic for the history, and is updated after taking action a and receiving observation o as follows:

$$b^{a,o}(s') = \frac{O(s', a, o) \sum_{s \in \mathcal{S}} T(s, a, s') b(s)}{p(o|a, b)}, \quad (1)$$

where $p(o|a, b) = \sum_{s \in \mathcal{S}} b(s) \sum_{s' \in \mathcal{S}} T(s, a, s') O(s', a, o)$ is a normalizing factor.

Policies

A policy is a mapping from the current belief state to an action.
A value function $\mathcal{V}_\pi(b)$ specifies the expected reward gained starting from b followed by policy π :

$$\mathcal{V}_\pi(b) = \sum_{s \in \mathcal{S}} b(s) R(s, \pi(b)) + \gamma \sum_{o \in \Omega} p(o|b, \pi(b)) \mathcal{V}_\pi(b^{\pi(b), o}).$$

POMDP Planning

Find an optimal policy that maximizes its value function.

Challenges

Computational complexity:

- We must reason in a $(n - 1)$ -dimensional continuous belief space (the curse of dimensionality).
- Complexity also grows fast with the length of planning horizon (the curse of history).

Challenges

Computational complexity:

- We must reason in a $(n - 1)$ -dimensional continuous belief space (the curse of dimensionality).
- Complexity also grows fast with the length of planning horizon (the curse of history).

Practical challenges:

- How to carry out intelligent information gathering in a large high dimensional belief space.
- How to scale up planning with long sequences of actions and delayed rewards.

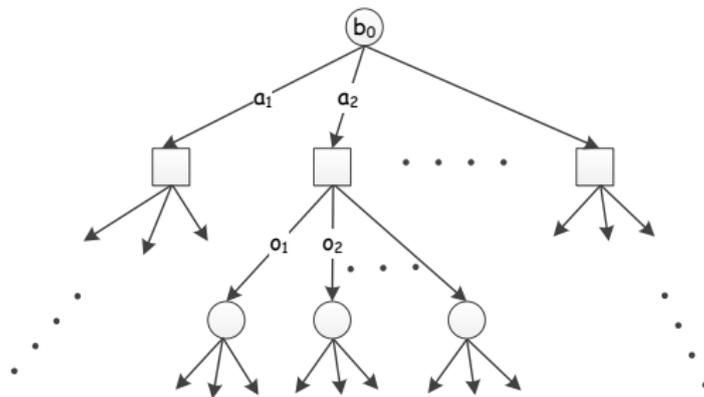


Figure: A belief tree rooted at b_0 .

Existing Solvers

- PBVI (Pineau et al., 2003)
- HSVI/HSVI2 (Smith and Simmons, 2004, 2005)
- FSVI (Shani et al., 2007)
- SARSOP (Kurniawati et al., 2008)
- RTDP-Bel (Bonet and Geffner, 2009)
- MiGS (Kurniawati et al., 2011)
- Some online solvers, e.g. PUMA (He et al., 2010), POMCP (Silver and Veness, 2010)

Our Approach: Information Gathering and Reward Exploitation of Subgoals (IGRES)

Main ideas:

Our Approach: Information Gathering and Reward Exploitation of Subgoals (IGRES)

Main ideas:

- 1 Sample potentially important states as subgoals.

Our Approach: Information Gathering and Reward Exploitation of Subgoals (IGRES)

Main ideas:

- 1 Sample potentially important states as subgoals.
- 2 Generate macro-actions (sequences of actions) for transitions to subgoals.

Our Approach: Information Gathering and Reward Exploitation of Subgoals (IGRES)

Main ideas:

- 1 Sample potentially important states as subgoals.
- 2 Generate macro-actions (sequences of actions) for transitions to subgoals.
- 3 Generate macro-actions for information gathering and reward exploitation in the neighborhood of subgoals.

Capturing Important States

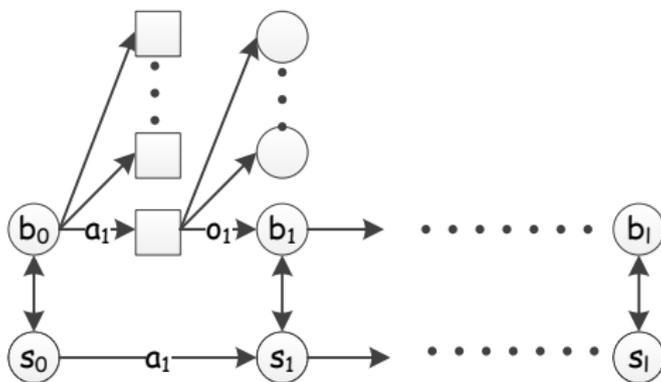
- 1 Identify importance heuristic functions with respect to:
 - Immediate reward;
 - Information gain.
- 2 Then a state is sampled as subgoal with probability of its importance.

Leverage the Structure

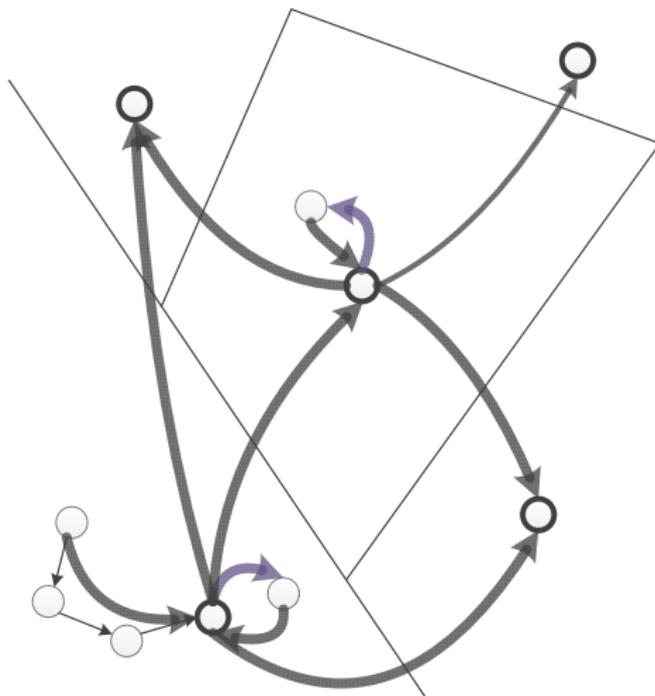
- 1 Specify distance between states with respect to the approximate similarity of their value.
- 2 Group all the states by the distance to the subgoals.

Sampling Belief States with Macro-actions

- 1 Associate each belief with an estimated current state.
- 2 Generate a macro-action towards the corresponding subgoal.
- 3 Gather information and exploit rewards around the subgoal.



Overview of IGRES



Conclusion

- Reduce planning complexity by using macro-actions.
- Sample potentially most useful beliefs (based on subgoal states).
- Capability of planning in large state space for a long horizon.

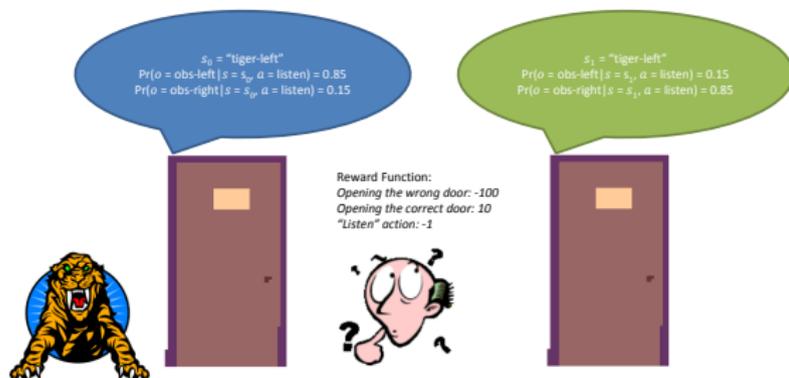


Figure: Tiger Domain

- Deterministic motions
- Noisy sensor for *Rock-goodness*
- +10 for sampling good
- -10 for sampling bad
- +10 for exiting
- No other cost/reward

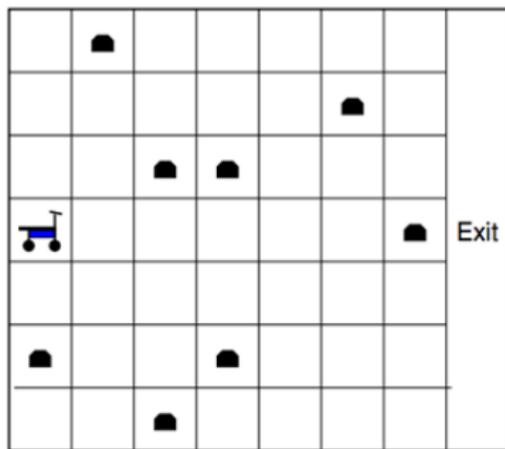


Figure: RockSample Domain

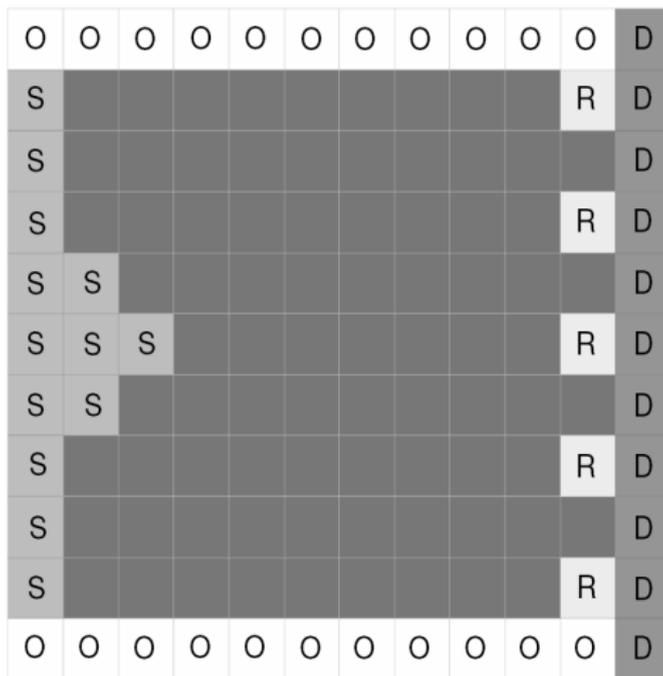


Figure: Underwater Domain

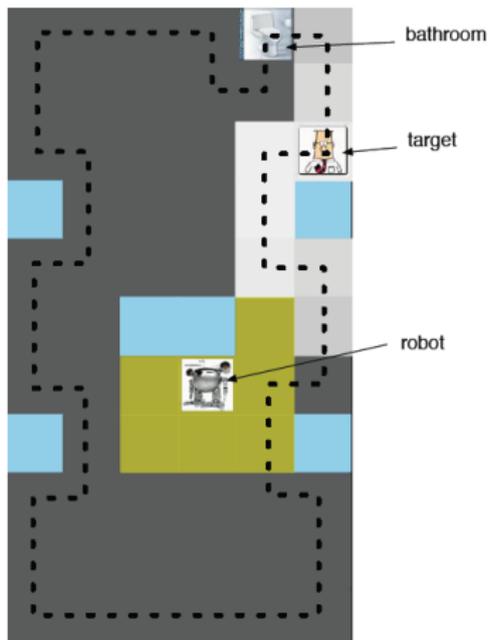


Figure: Homecare Domain

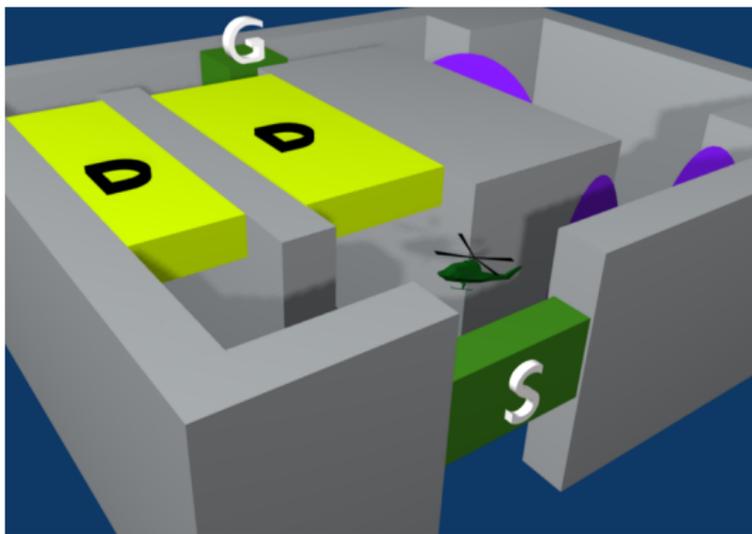


Figure: 3D-Navigation Domain

Results of the benchmark domains.

	Errors	Times
Tiger		
[S] = 2, [A] = 3, [O] = 2		
RTDP-Bul	19.42 ± 0.59	0.50
HV12	19.41 ± 0.69	<1
FV3V*	N/A	
SARSOP	18.59 ± 0.61	0.00
MGL	-19.58 ± 0	100
EGRES (9 subgoals: 8)	19.41 ± 0.59	1
Noise-Open		
[S] = 2, [A] = 3, [O] = 2		
RTDP-Bul	-13.67 ± 0.28	1.22
HV12	-13.69 ± 0.64	<1
FV3V*	N/A	
SARSOP	-13.66 ± 0.18	0.18
MGL	-19.88 ± 0	100
EGRES (9 subgoals: 8)	-13.67 ± 0.18	1
RockSample(4,6)		
[S] = 207, [A] = 9, [O] = 2		
RTDP-Bul	17.94 ± 0.12	10.7
HV12	17.92 ± 0.61	<1
FV3V	17.85 ± 0.18	1
SARSOP	17.75 ± 0.12	0.7
MGL	8.57 ± 0	100
EGRES (9 subgoals: 4)	17.30 ± 0.12	10
RockSample(7,8)		
[S] = 1245, [A] = 11, [O] = 2		
RTDP-Bul	29.55 ± 0.13	103
HV12	31.69 ± 0.80	100
FV3V	29.98 ± 0.20	102
SARSOP	21.85 ± 0.13	100
MGL	7.35 ± 0	100
EGRES (9 subgoals: 8)	19.54 ± 0.12	100
Halfworld		
[S] = 92, [A] = 5, [O] = 17		
RTDP-Bul	0.237 ± 0.006	1004
HV12	0.507 ± 0.001	250
FV3V	0.891 ± 0.007	280
SARSOP	0.530 ± 0.008	300
MGL	0.522 ± 0.008	300
EGRES (9 subgoals: 26)	0.538 ± 0.008	300
Tag		
[S] = 870, [A] = 5, [O] = 30		
RTDP-Bul	-6.32 ± 0.12	372
HV12	-6.36 ± 0.10	400
FV3V	-6.11 ± 0.11	35
SARSOP	-6.06 ± 0.12	30
MGL	-6.00 ± 0.12	30
EGRES (9 subgoals: 26)	-6.12 ± 0.12	30
Underwater Navigation		
[S] = 2613, [A] = 6, [O] = 103		
RTDP-Bul	750.07 ± 0.26	338
HV12	718.17 ± 0.60	400
FV3V	725.88 ± 5.91	414
SARSOP	731.33 ± 1.14	150
MGL	715.59 ± 1.37	400
EGRES (9 subgoals: 26)	749.94 ± 0.30	50
Hammers		
[S] = 5409, [A] = 9, [O] = 928		
RTDP-Bul**	N/A	
HV12	15.07 ± 0.37	2000
FV3V**	N/A	
SARSOP	16.64 ± 0.42	1000
MGL	16.70 ± 0.46	1000
EGRES (9 subgoals: 36)	17.32 ± 0.45	1000
3D-Navigation		
[S] = 10000, [A] = 5, [O] = 14		
RTDP-Bul	-93.83 ± 0.01	2115
HV12	-91.98 ± 0	2000
FV3V**	N/A	
SARSOP	-99.97 ± 0	800
MGL	(2.977 ± 0.512) × 10 ⁴	150
EGRES (9 subgoals: 143)	(3.272 ± 0.160) × 10 ⁴	150

* ArrayIndexOutOfBoundsException in theory.

** Solution is not able to compute a solution given large amount of computation time.

*** OutOfMemoryError in theory.

Results of the adaptive management of migratory birds problem.

	Return	Time(s)
Lesser sand plover		
$ \mathcal{S} = 108, \mathcal{A} = 3, \Omega = 36$		
symbolic Perseus*	4675	10
IGRES (# subgoals: 18)	5037.72 ± 8.82	10
Bar-tailed godwit b.		
$ \mathcal{S} = 972, \mathcal{A} = 5, \Omega = 324$		
symbolic Perseus*	18217	48
IGRES (# subgoals: 36)	19572.41 ± 39.35	60
Terek sandpiper		
$ \mathcal{S} = 2916, \mathcal{A} = 6, \Omega = 972$		
symbolic Perseus*	7263	48
IGRES (# subgoals: 72)	7867.95 ± 2.44	60
Bar-tailed godwit m.		
$ \mathcal{S} = 2916, \mathcal{A} = 6, \Omega = 972$		
symbolic Perseus*	24583	58
IGRES (# subgoals: 72)	26654.06 ± 38.60	60
Grey-tailed tattler		
$ \mathcal{S} = 2916, \mathcal{A} = 6, \Omega = 972$		
symbolic Perseus*	4520	378
IGRES (# subgoals: 72)	4860.91 ± 38.47	60
IGRES (# subgoals: 72)	4927.17 ± 38.14	300

* Results from Nicol et al. (2013).

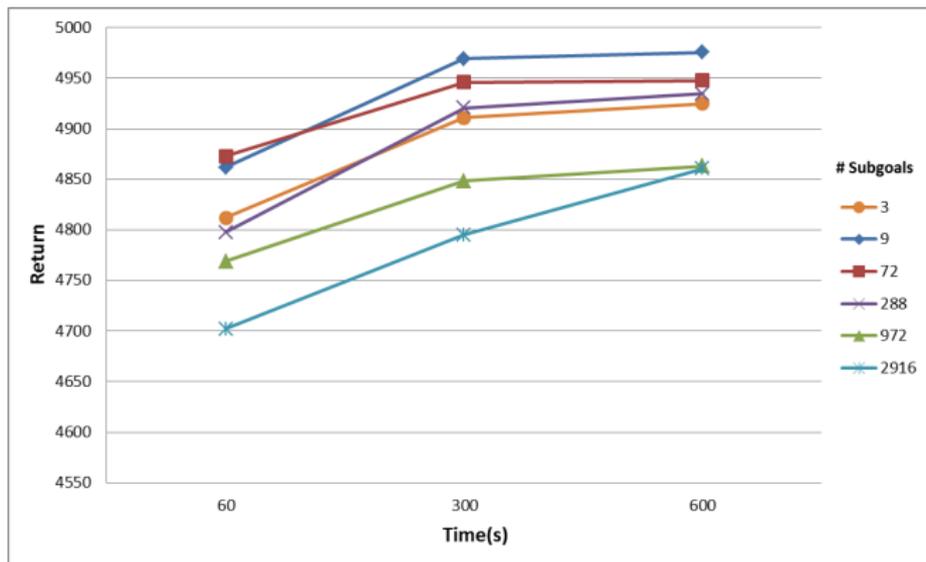


Figure: Performances on Grey-tailed tattler Domain

Thank You!

Reference I

- B. Bonet and H. Geffner. Solving POMDPs: RTDP-bel vs. point-based algorithms. In *International Joint Conference on Artificial Intelligence*, 2009.
- R. He, E. Brunskill, and N. Roy. Puma: Planning under uncertainty with macro-actions. In *National Conference on Artificial Intelligence*, 2010.
- H. Kurniawati, D. Hsu, and W. S. Lee. SARSOP: Efficient point-based POMDP planning by approximating optimally reachable belief spaces. In *Robotics: Science and Systems*, 2008.
- H. Kurniawati, Y. Du, D. Hsu, and W. S. Lee. Motion planning under uncertainty for robotic tasks with long time horizons. *International Journal of Robotics Research*, 30(3):308–323, 2011.

Reference II

- S. Nicol, O. Buffet, T. Iwamura, and I. Chadès. Adaptive management of migratory birds under sea level rise. In *International Joint Conference on Artificial Intelligence*, 2013.
- J. Pineau, G. Gordon, and S. Thrun. Point-based value iteration: An anytime algorithm for POMDPs. In *International Joint Conference on Artificial Intelligence*, 2003.
- G. Shani, R. I. Brafman, and S. E. Shimony. Forward search value iteration for POMDPs. In *International Joint Conference on Artificial Intelligence*, 2007.
- D. Silver and J. Veness. Monte-carlo planning in large POMDPs. In *Advances in Neural Information Processing Systems*, 2010.
- T. Smith and R. G. Simmons. Heuristic search value iteration for POMDPs. In *Uncertainty in Artificial Intelligence*, 2004.

Reference III

T. Smith and R. G. Simmons. Point-based POMDP algorithms: Improved analysis and implementation. In *Uncertainty in Artificial Intelligence*, 2005.