

# Point Based Value Iteration with Optimal Belief Compression for Dec-POMDPs

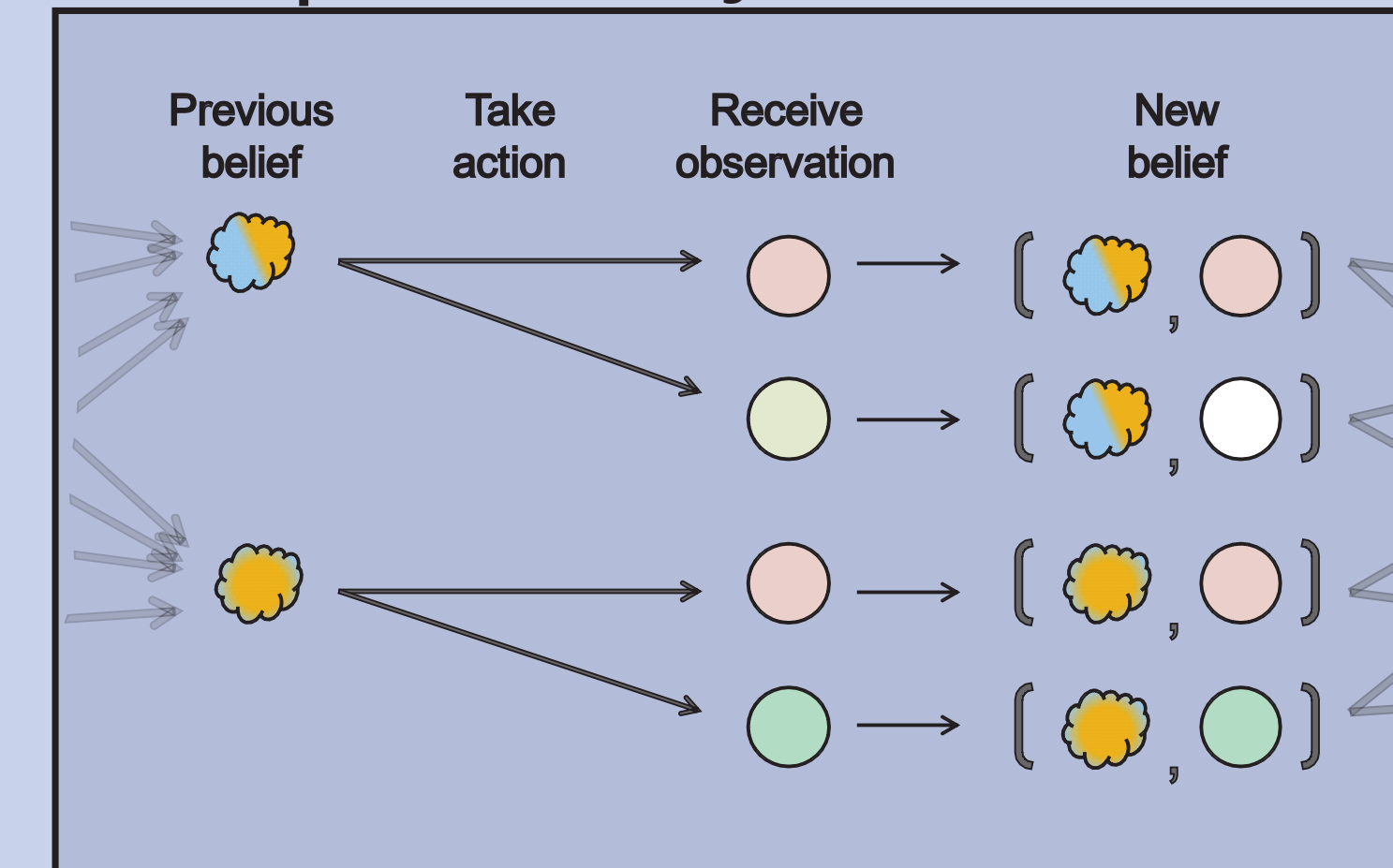
**Problem:** An agents' belief set has doubly exponential growth with the horizon due to nested beliefs.

**Key Idea:** Allow agent's to decide how to best compress their own beliefs and fold this belief choice into the model itself thereby allowing our model solver to compute an optimal compression scheme.

**Formal Solution:** Given a DecPOMDP  $\langle N, A, S, O, P, R, s^{(0)} \rangle$  (with states  $\omega \in S$  and observations  $\sigma \in O$ ) we define the BB-DecPOMDP approximation model  $\langle N', A', O', S', P', R', s'^{(0)} \rangle$  (which optimally compresses beliefs) with belief set size parameters  $t_1, \dots, t_n$  as:

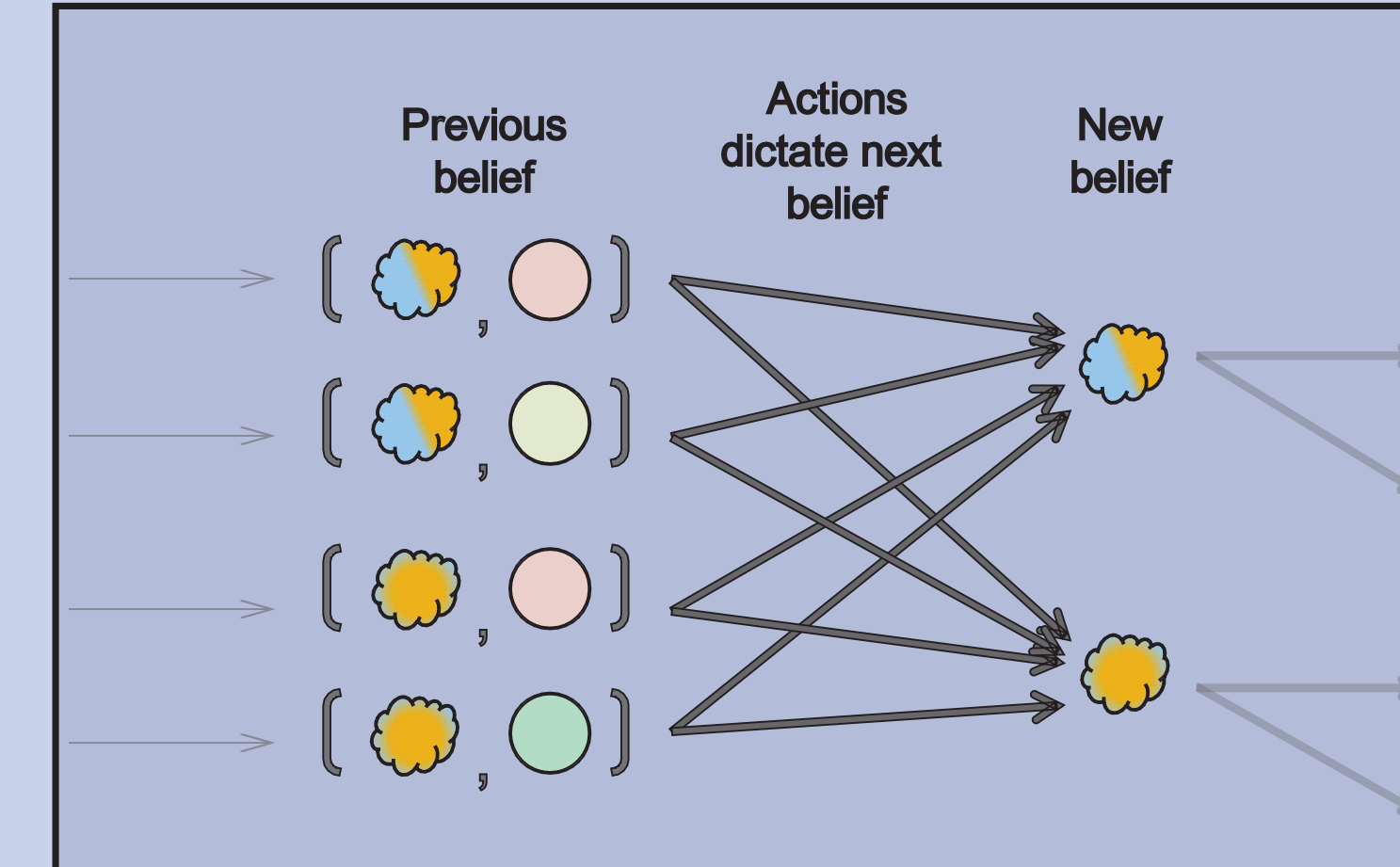
- $N' = N$  and  $A'_i = \{a_{i1}, \dots, a_{i \max(|A_i|, t_i)}\}$
- $O'_i = \{O_i, \emptyset\} \times \{1, 2, \dots, t_i\}$  with factored observation  $o_i = \langle \sigma_i, \theta_i \rangle \in O'$
- $S' = S \times O'$  with factored state  $s = \langle \omega, \sigma_1, \theta_1, \dots, \sigma_n, \theta_n \rangle \in S'$
- $P'(s'|s, a) = \begin{cases} P(\omega', \langle \sigma'_1, \dots, \sigma'_n \rangle | \omega, a) & \text{if } \forall i: \sigma_i = \emptyset, \sigma'_i \neq \emptyset \text{ and } \theta'_i = \theta_i \\ 1 & \text{if } \forall i: \sigma_i \neq \emptyset, \sigma'_i = \emptyset \text{ and } \theta'_i = a_i, \omega' = \omega \\ 0 & \text{otherwise} \end{cases}$
- $R'(s, a) = \begin{cases} R(\omega, a) & \text{if } \forall i: \sigma_i = \emptyset \\ 0 & \text{otherwise} \end{cases}$
- $s'^{(0)} = \langle s^{(0)}, \theta_1, 1, \dots, \theta_n, 1 \rangle$  is the initial state distribution

**Belief Expansion (from original model)**



Model alternates between these two phases

**Belief Contraction (new phase)**



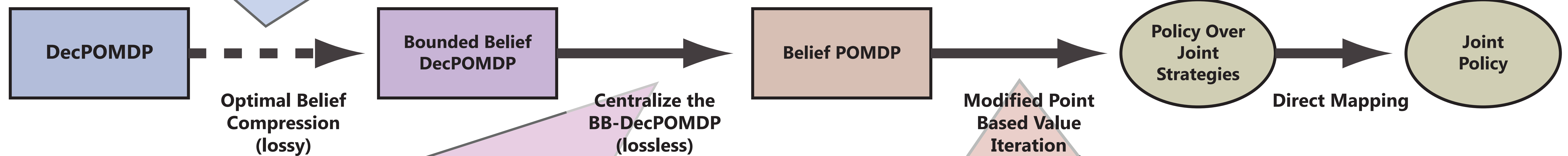
## Results on Benchmark Problems

	S	A <sub>i</sub>	O <sub>i</sub>	Previous Best	1-Belief		2-Beliefs	
				Utility	Utility	Γ	Utility	Γ
Dec-Tiger	2	2	2	13.4486	-20.000	2	4.6161	187
Broadcast	4	2	2	9.1	9.2710	36	9.2710	44
Recycling	4	3	2	31.92865	26.3158	8	31.9291	13
Grid small	16	5	2	6.89	5.2716	168	6.8423	206
Box Pushing	100	4	5	149.854	127.1572	258	223.8674	357
Wireless	64	2	6	-175.40	-208.0437	99	-167.1025	374

	3-Beliefs		4-Beliefs		5-Beliefs	
	Utility	Γ	Utility	Γ	Utility	Γ
Dec-Tiger	13.4486	231	13.4486	801	13.4486	809
Broadcast	9.2710	75	9.2710	33	9.2710	123
Recycling	31.9291	37	31.9291	498	31.9291	850
Grid small	6.9826	276	6.9896	358	6.9958	693
Box Pushing	224.1387	305	-	-	-	-

Utility achieved by our PBVI-BB-DecPOMDP algorithm compared to the previously best known policies on a series of standard benchmarks. Higher is better. Our algorithm beats all previous results except on Dec-Tiger where we believe an optimal policy has already been found.

**Approach:** Solve a DecPOMDP by transforming it into a POMDP



**Problem:** How can we centralize the model in order to make the problem tractable?

**Key Ideas:** Factor the state so that it contains each agent's current belief. Beliefs are only labels and can be decoupled from history. We can convert a BB-DecPOMDP into a POMDP with exponential actions corresponding to strategies of a CBG by utilizing the fact that the common knowledge distribution over joint-beliefs is a sufficient statistic for planning (previously proved)

**Formal Solution:** Belief-POMDP  $\langle A', S', O', P', R', s'^{(0)} \rangle$  converted from BB-DecPOMDP  $\langle N, A, S, O, P, R, s^{(0)} \rangle$  (with belief labels  $\Theta_i$  for each agent).

**Factored states:**  $S' = S \times \prod_{i=1}^n \Theta_i$  with  $\langle \omega, \theta_1, \dots, \theta_n \rangle \in S'$  where  $\omega \in S$  is the underlying state and  $\theta_i \in \Theta_i$  is agent  $i$ 's belief.

**No observations:**  $O' = \{\}$

**Actions are strategies:**  $A' = \prod_{i=1}^n \prod_{j=1}^{|\Theta_i|} A_{ij}$  (one action for each agent for each belief)

**Transition function:**  $P'(s'|s, a) = \sum_{\theta_i} P(\omega', o | \omega, \langle a_{\theta_1}, \dots, a_{\theta_n} \rangle)$  where  $a_{\theta_i}$  is the action agent  $i$  would take in belief  $\theta_i$ .

**Reward function:**  $R'(s, a) = R(\omega, \langle a_{\theta_1}, \dots, a_{\theta_n} \rangle)$

**Initial state:**  $s'^{(0)} = s^{(0)}$

**Problem:** The Belief POMDP has an exponential action space corresponding to joint-strategies.

**Key Idea:** Modify PBVI so that instead of finding the best action (e.g. strategy) exhaustively, solve the corresponding CBG.

**Formal Solution:** Inputs: DecPOMDP  $M$ , discount  $\gamma$ , belief bounds  $|\Theta_i|$ , stopping criterion  $\epsilon_T$ . Output: value function  $\Gamma$

- 1:  $\langle N, A, O, S, P, R, s^{(0)} \rangle \leftarrow$  BB-DecPOMDP approximation of  $M$
- 2:  $B^y \leftarrow$  sampling of states using a random walk from  $s^{(0)}$
- 3:  $\Gamma' \leftarrow \{ \{R_{\min}/\gamma, \dots, R_{\max}/\gamma\} \}$
- 4: **repeat**
- 5:  $B \leftarrow B^y$ ;  $\Gamma \leftarrow \Gamma'$ ;  $\Gamma' \leftarrow \emptyset$
- 6: **while**  $B \neq \emptyset$  **do**
- 7:  $b \leftarrow \text{Rand}(b \in B)$
- 8:  $\alpha \leftarrow \Gamma(b)$
- 9:  $\alpha' \leftarrow$  optimal point of integer program solving CBG at  $b$
- 10: **if**  $\alpha'(b) > \alpha(b)$  **then**
- 11:  $\alpha \leftarrow \alpha'$
- 12:  $\Gamma' \leftarrow \Gamma' \cup \alpha$
- 13: **for all**  $b \in B$  **do**
- 14: **if**  $\alpha(b) > \Gamma(b)$  **then**
- 15:  $B \leftarrow B/b$
- 16: **until**  $\Gamma' - \Gamma < \epsilon_T$
- 17: **return**  $\Gamma$

**Problem:** Problem: How do we solve a Cooperative Bayesian Game (CBG) efficiently given that it's NP complete?

**Key Idea:** Solve for a deterministic agent normal form correlated equilibrium. This is an ILP. The linear relaxation is very often integral.

**Formal Solution:**

**Maximize:**  $V_{\bar{a}, \alpha}(b) = \sum_{s \in S} b(s) \sum_{\bar{a} \in A} \sum_{s' \in S} P(s'|s, \bar{a}) \alpha(s') x_{\bar{a}, \theta \in s}$

**Over:**  $x \in \{0, 1\}^{|\Sigma|}$ ,  $\alpha \in \Gamma$

**Subject to:** For each agent  $i$ , joint-type  $\theta$ , and partial joint-actions of other agents  $\bar{a}_{-i}$

$$\sum_{a_i \in A_i} x_{\bar{a}, \theta} = x_{(\bar{a}_{-i}, \theta_{-i})}$$

$$\text{for each } \theta \in \Theta: \sum_{\bar{a} \in A} x_{\bar{a}, \theta} = 1$$

$$\text{for each } \theta \in \Theta, \bar{a} \in A: x_{\bar{a}, \theta} \geq 0$$