



Summary

and outperform baselines.

- not deployment time.
- problems that require information-seeking behavior

Introduction

training time.

Motivating Example (Sphinx Domain)

- \checkmark One of the boxes contains a +1 reward.





Information Seeking Dilemma



Learning Belief Representations for Partially Observable Deep RL

Andrew Wang^{*12} Andrew C. Li^{*12} Toryn Q. Klassen¹²³ Rodrigo Toro Icarte⁴⁵ Sheila A. McIlraith¹²³

¹Department of Computer Science, University of Toronto ²Vector Institute ³Schwartz Reisman Institute ⁴Pontificia Universidad Católica de Chile ⁵Centro Nacional de Inteligencia Artificial * denotes equal contribution

information appears in	discarded
both s_t and o_t	
information appears in	$\psi(o_t)$
both s_t and o_t	
information appears in	$\phi(s_t)$
s_t but not o_t	

Experiments

Domains



Fig. 1: Goal: reach the box containing +1 reward on each episode. $8 \times 8 \times 3$ grid. NoisyTV variation: s_t, o_t flash random colours. Lying sphinx variation: sphinx lies 50% of time.

Main Results

tial VAE.





Fig. 4: Varying the cost of information for the Sphinx domains. Most approaches learn to seek information when cost of information is low, but when cost of information is high, only our method can make any progress.

akeaways		
	Standard PO-RL	Our Approach
Goal	Learn a POMDP policy $\pi(a_t h_t)$	Learn a POMDP policy $\pi(a_t h_t)$
Learns	Agent learns by observing	Agent learns by observing
	observations and rewards	observations, states, and
		rewards
Predicts	Agent directly predicts actions	Agent predicts belief states
	+ values given history.	given history, and actions +
		values given belief states.









cookie. $9 \times 9 \times 3$ grid.



Fig. 3: Goal: Find the exit pointed to by statues across the map. Continuous-state, continuous-action environment. Observations represented as 100×100 image.

Our method improves upon baselines in all domains. We compared vs. Recurrent PPO, Asymmetric Actor-Critic, Unbiased Asymmetric Actor-Critic, Sequen-

Fig. 5: A visualization of various belief states learned in the Sphinx task.