

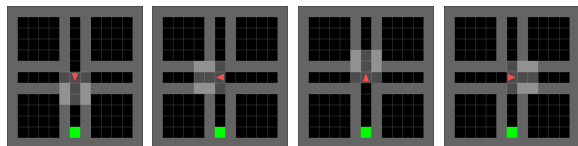


Privileged Experts and Mismatch With Students

Experts have access to full state
Student sees partial (egocentric) observations

Intuitive Example

Expert sees:



Expert recommends:

Move forward Turn left Turn left/right Turn right

But student's partial view makes these recommendations confusing...



Hence, student learns an averaged-policy:

Prop. 1: A student's policy is the average of the teacher's policy.

$$\text{If } \pi^{\text{IL}} = \underset{\text{student policies } \pi}{\text{argmin}} \mathbb{E}[\text{CrossEntropy}(\pi, \pi^{\text{exp}})]$$

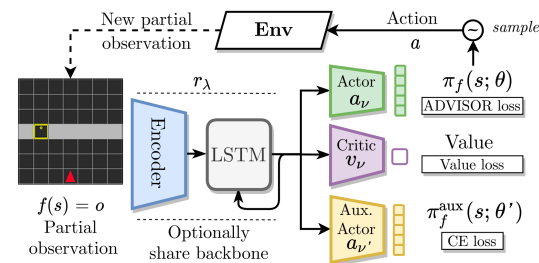
$$\text{Then } \pi^{\text{IL}}(o) = \mathbb{E}[\pi^{\text{exp}}(S) \mid f(S) = o]$$

Adaptive Insubordination (ADVISOR)

Idea:

1. Student estimates if it can imitate the expert
2. Accordingly, weigh IL & RL losses at each step

Schematic:

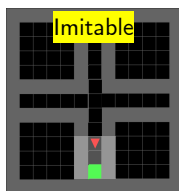


$$\mathcal{L}^{\text{ADV}}(\theta) = \mathbb{E}_{\mu}[w(S) \cdot \text{CE}(\pi^{\text{exp}}(S), \pi_f(S; \theta)) + (1 - w(S)) \cdot L(\theta, S)]$$

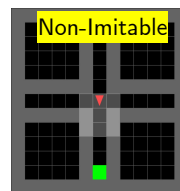
Auxiliary actor:

- Trained only by IL
- Estimates imitability of the current state

Outcome:



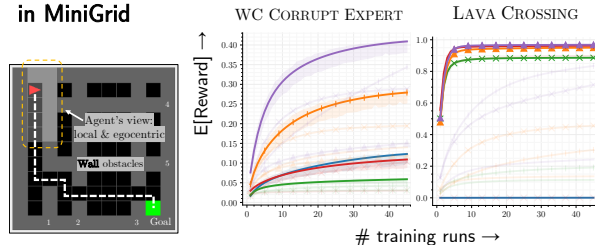
Small Aux. Actor Loss
High $w(S)$ → More IL



High Aux. Actor Loss
Low $w(S)$ → More RL

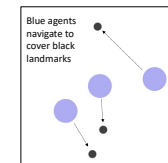
Results

Exhaustive evaluation in MiniGrid

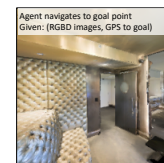


Previous methods				ADVISOR (ours)	
RL only:	IL only:	IL + RL & IL → RL:	Demonstration-based	12. ADV	13. \uparrow → ADV
1. PPO	2. BC	5. BC + PPO (static)	9. BC ^{obs}	14. BC ^{sil-1} → ADV	15. AD ^{obs} + PPO
	3. DAGger	6. BC → PPO	10. BC ^{obs} + PPO		
	4. BC ^{sil-1}	7. \uparrow → PPO	11. GAIL		
		8. BC ^{sil-1} → PPO			

Success across diverse tasks



COOPNAV
(in MPE)



POINTNAV
(in AIHABITAT)



OBJECTNAV
(in A12-THOR)

Tasks → Training routines ↓	PointGoal Navigation SPL		ObjectGoal Navigation SPL		Cooperative Navigation Reward	
	@10%	@100%	@10%	@100%	@10%	@100%
RL only	30.9	54.7	6.7	13.1	-561.8	-456.0
IL only	30.1	68.7	3.8	9	-460.3	-416.7
IL + RL static	48.9	71.5	6.5	11.3	-475.5	-424.6
ADVISOR (ours)	57.7	77.1	11.9	14.1	-419.9	-405.6