

Motivation

MARL agents often converge to arbitrary conventions.



Split by ingredient



Split by side

Can we generate a diverse set of conventions to learn to work with people?

Diversity Definition

Want to maximize nearest-neighbors score for set of conventions, D:

$$S(D) = \mathbb{E}_{\pi \in D_{test}} [\max_{\pi^* \in D} J(\pi^*, \pi)]$$

For ad-hoc coordination, train a single convention-aware agent:

$$L(\hat{\pi}, D) = -J(\hat{\pi}, \hat{\pi}) - \frac{\lambda}{D} \sum_{\pi \in D} \mathbb{E}_{(o,a) \sim \pi} [\log(\hat{\pi}(a \ o))]$$

To evaluate the performance of the convention-aware agent:

$$J(\hat{\pi}, D_{test}) = \mathop{\mathbb{E}}_{\pi \in D_{test}} \left[J(\hat{\pi}, \pi) \right]$$

Statistical Diversity

Some approaches add a term to induce statistical variations in trajectories:

$$C_{ADAP}(D) = \mathop{\mathbb{E}}_{s \in S} \left[\mathop{\mathbb{E}}_{\substack{\pi_1, \pi_2 \in D \\ \pi_1 \neq \pi_2}} \exp(-D_{\mathsf{KL}} \begin{pmatrix} \pi_1(s) & \pi_2(s) \end{pmatrix} \right]$$

These fail to identify trivial variations in navigation:



Diverse Conventions for Human-Al Collaboration

Bidipta Sarkar, Andy Shih, and Dorsa Sadigh

Stanford University

Cross-Play Minimization

Adding a new convention π_n guarantees an increase in the score S:

$$S(D_n) \ge S(D_{n-1}) + p(\pi_n)(J(\pi_n, \pi_n) - \max_{\pi^* \in D_{n-1}} J(\pi_n, \pi^*))$$

Maximizing this lower bound gives a cross-play minimization algorithm.

– Handshakes and Mixed-Play -

Pure cross-play minimization incentivizes sabotages under cross-play.



Mixed-play ensures that conventions always act in good faith.



CoMeDi: Cross-Play and Mixed-Play

Loss function to generate convention π_n

$$L(\pi_n) = -J(\pi_n, \pi_n) + \alpha J(\pi_n, \pi^*) - \beta J_{M}(\pi_n, \pi^*)$$

Where π^* is the most-compatible previously discovered convention.





Simulation Results

In Blind Bandits, CoMeDi finds both the S and G conventions:





In Balance Beam, tuning the mixed-play weight eliminates handshakes.



β	SP↑	$\mathrm{XP}\downarrow$	$\mathrm{HS}\downarrow$	$PX\downarrow$	LS ↑	$RS\downarrow$
0.00	1.616	-0.392	0.16	0.00	1.128	-0.656
0.25	1.808	0.096	0.00	0.00	1.272	0.096
0.50	2.000	0.200	0.00	0.00	1.384	0.128
1.00	1.904	0.632	0.00	0.24	1.232	0.160

In Overcooked, CoMeDi's convention-aware agent has strong cross-play performance with held-out conventions.

layout	CoMeDi	XP	ADAP	SP	PPO_BC	layout	CoMeDi	XP	ADAP	SP	PBT
simple	4.55	3.49	2.79	4.36	3.75	simple	4.91	4.46	2.77	4.51	3.75
random1	3.42	1.58	0.72	2.70	2.29	random1	3.43	1.81	1.10	2.38	2.29
random3	1.19	0.49	0.03	0.33	1.17	random3	1.91	1.1	0.03	1.13	1.17
unident_s	5.27	3.62	1.84	2.51	2.58	unident_s	2.94	3.16	2.59	4.52	2.58
random0	0.37	0.49	0.00	0.21	1.44	random0	1.56	0.88	0	0.25	1.44

Scores with PBT

Scores with PPO_BC

Overcooked User Study —

CoMeDi significantly outperforms baselines in terms of pure score and user opinion, even surpassing human-level performance.



Live Demo, Code, and Videos



Visit our website: https://iliad.stanford.edu/Diverse-Conventions/