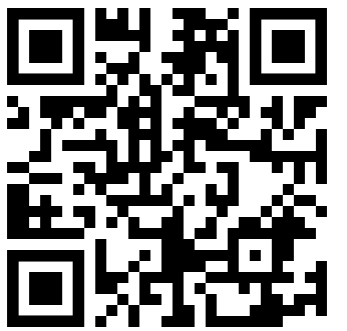




REMEMBERING THE MARKOV PROPERTY IN COOPERATIVE MARL

K. TESSERA, L. HINCKELDEY, R. ZAMBONI, D. ABEL, A. STORKEY
k.tessera@ed.ac.uk



Do current MARL environments test genuine cooperative reasoning, the kind that requires **behaviours grounded in observations and memory**?

1) PROBLEM FORMULATION

THE PROBLEM: RECOVERING A MARKOV SIGNAL

Setting: Dec-POMDP

$\mathcal{M} = (\mathcal{N}, \mathcal{S}, T, O, \mu, \{\mathcal{A}^i\}, \{\mathcal{O}^i\}, R, \gamma)$; agents receive partial o_t^i and act with $\pi^i(a_t^i | h_t^i)$.

Goal: From histories h_t^i , recover a *Markovian signal*: a belief/approx. over environment state and teammates' behaviour that is sufficient for control.

Question: Do current MARL environments require reconstructing this signal? Do MARL agents do this?

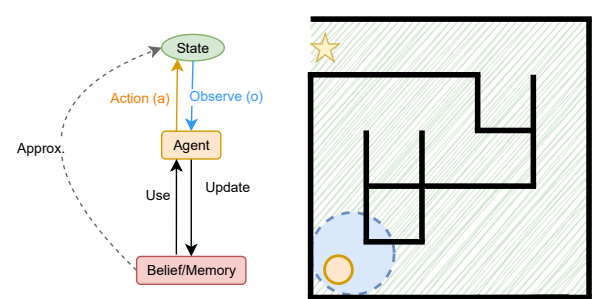
MUTUAL INFORMATION

Observation-grounding:

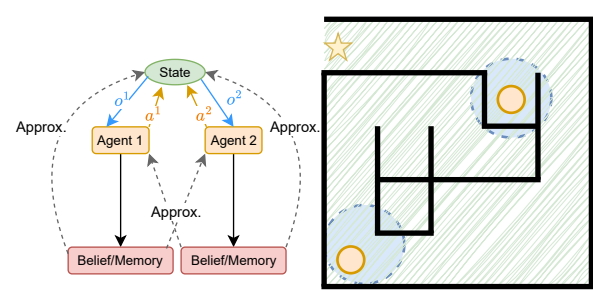
- $\mathbb{I}(O; A)$: dependence between actions and current observations
- $\mathbb{I}(H; A)$: dependence between actions and agent history (e.g., RNN hidden state)

Memory:

- Empirical test:** Performance with and without memory
- Compare $\mathbb{I}(H; A)$ and $\mathbb{I}(O; A)$**
If $\mathbb{I}(H; A) \gg \mathbb{I}(O; A)$, agents rely more on memory than immediate observations.

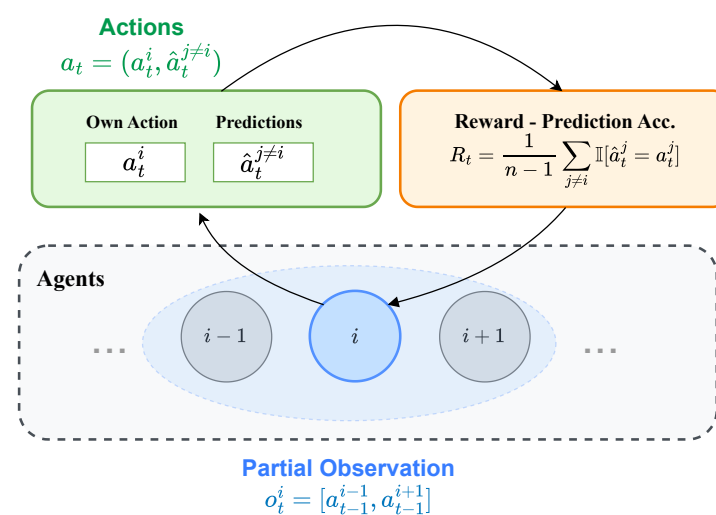


(a) **POMDP:** If you can't see, you must remember

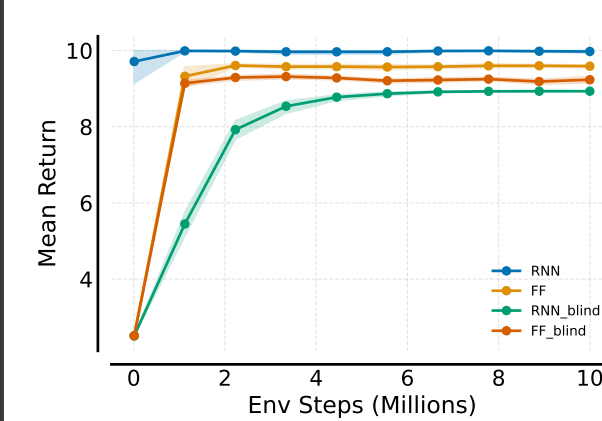


(b) **Dec-POMDP:** If you can't see, you must predict

2) CASE STUDY: BRITTLE CONVENTIONS VS. ROBUST COORDINATION

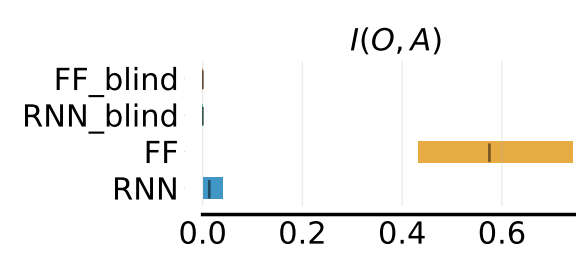


(a) **Prediction Game**



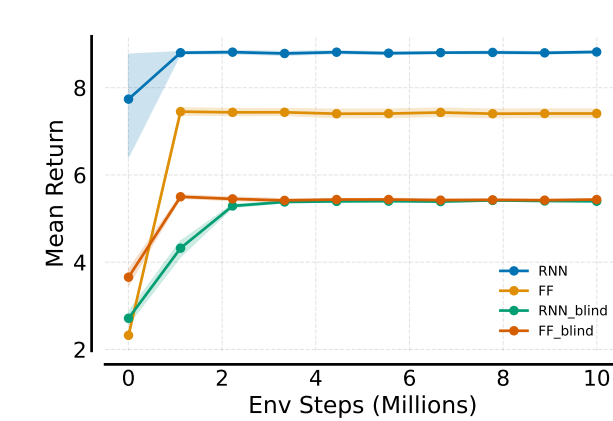
(a) Co-adapt: High Perf.

(b) **Scenario A: Co-adapting**

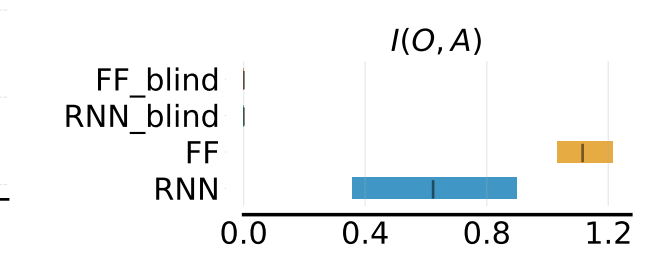


(b) **Low $I(O, A)$**

(c) **Scenario B: Mixed**



(c) Mixed: High Perf.



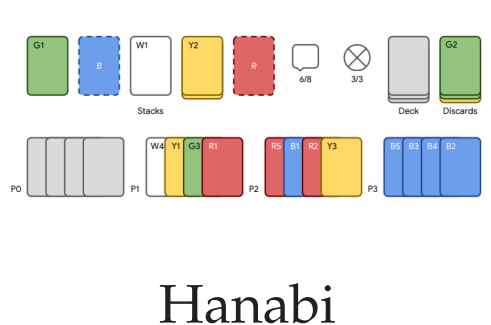
(d) **Higher $I(O, A)$**

BRITTLE CONVENTIONS VS. ROBUST COORDINATION

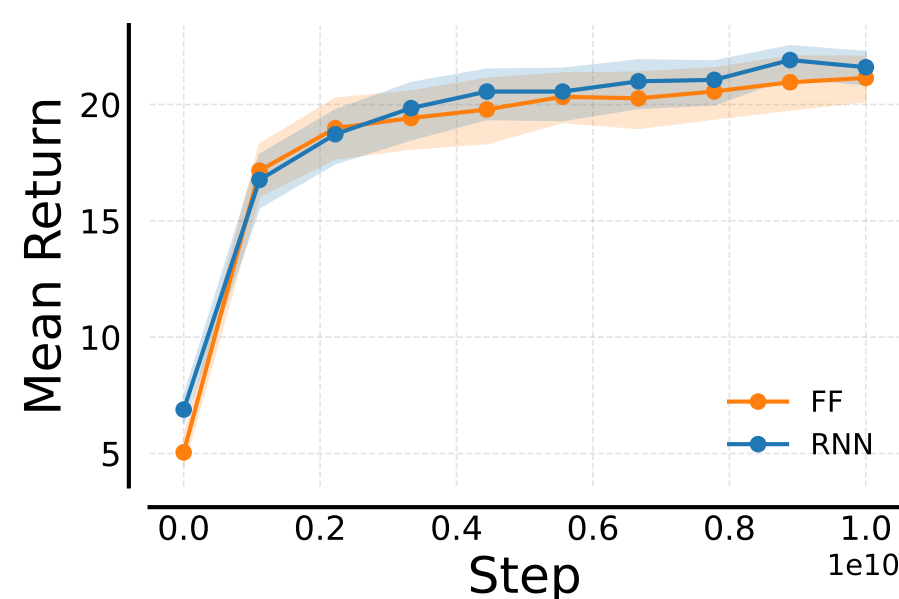
Same method, but the **mechanism** for success changes with environment modifications (partner composition).

Implication: Current MARL environments may **enable** fragile co-adaptation rather than robust cooperation.

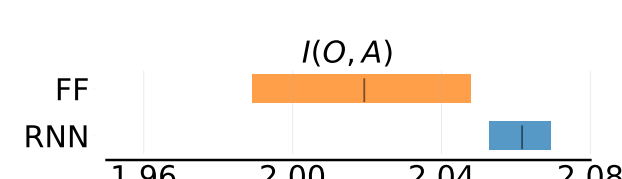
3) EXPERIMENTS: MODERN ENVIRONMENTS



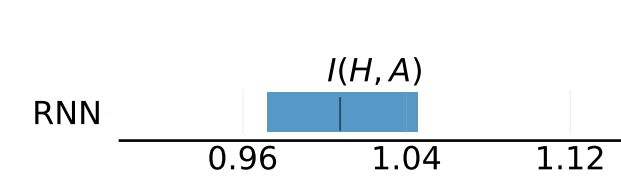
Hanabi



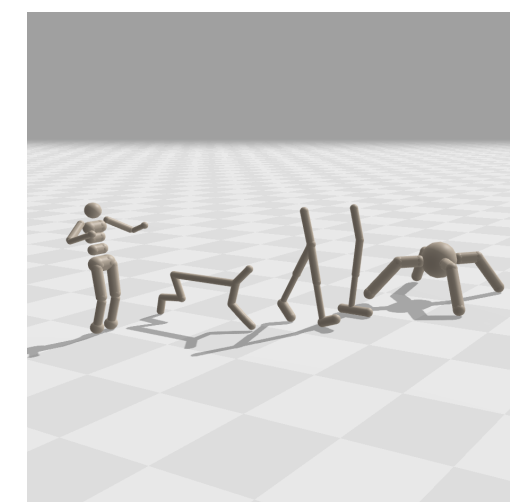
Hanabi (Two-Players).



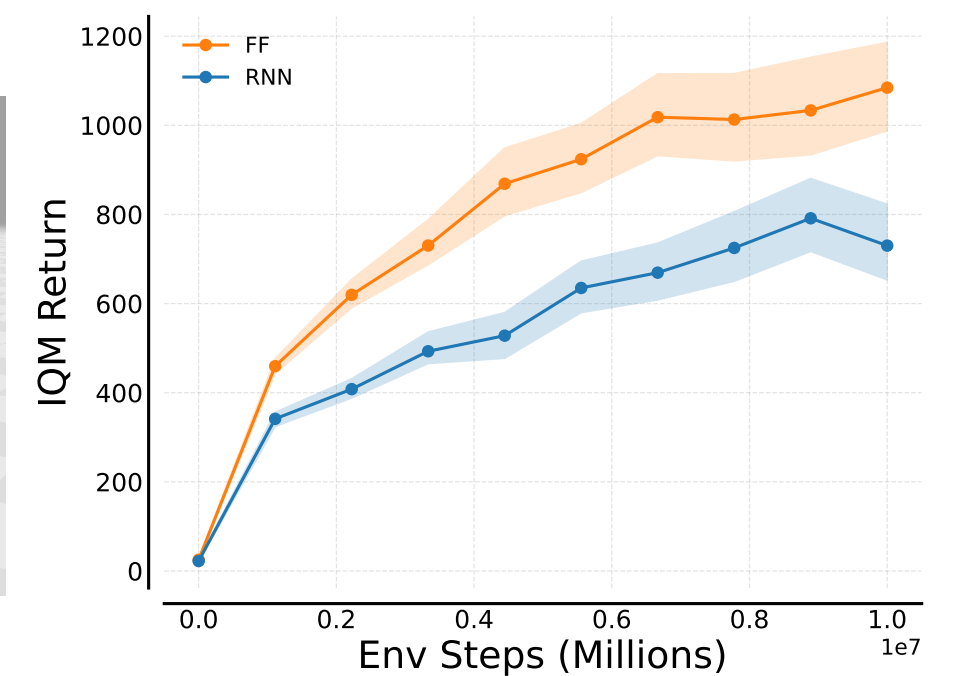
$\mathbb{I}(O; A)_{max} \approx 3$



$\mathbb{I}(H; A)_{max} \approx 3$



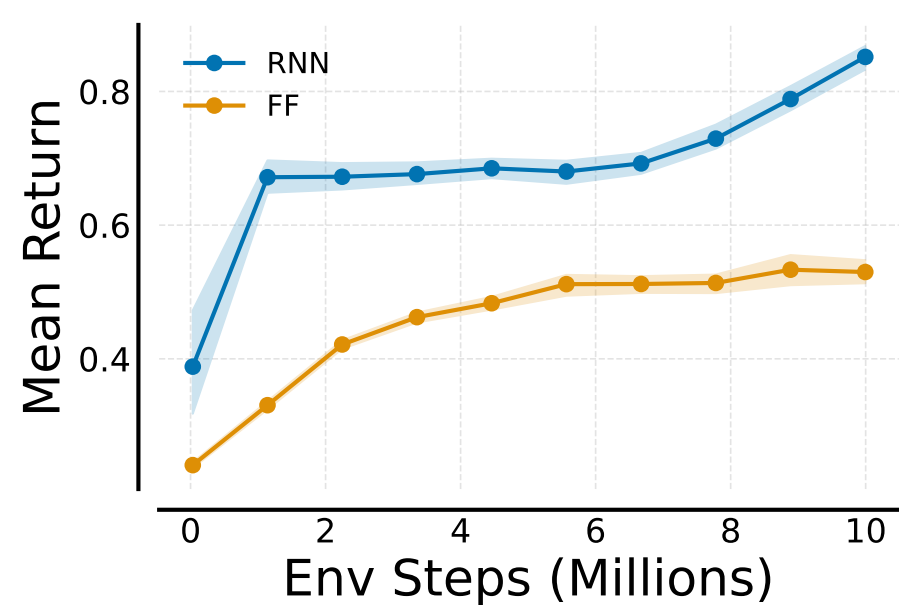
MaBrax



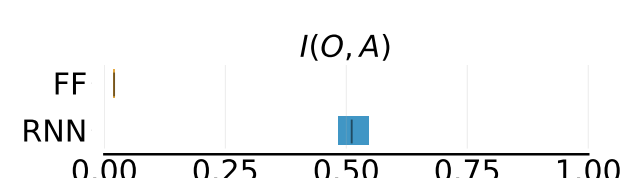
MaBrax Results (Across 5 settings)



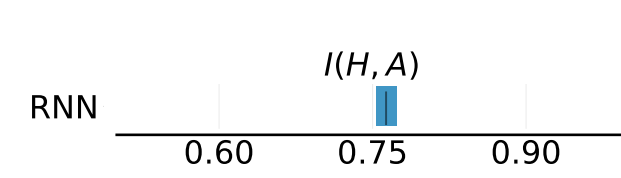
SMAX



SMAX-v2 style (5 Units).



$\mathbb{I}(O; A)_{max} \approx 2.30$



$\mathbb{I}(H; A)_{max} \approx 2.30$

SUMMARY

Environment	Obs-Grounded	Memory-Based
Hanabi	✓	✗
MaBrax	✗	✗
SMAX	✓ / ✗	✓

TAKEAWAYS

- Modern model-free recurrent MARL methods **can** learn robust cooperative behaviour when environments **necessitate** this.
- Yet current MARL environments may inadvertently allow success through **alternative means** (e.g. blind conventions, memoryless coordination) rather than genuine cooperation.

We therefore advocate for new cooperative environments built upon two core principles: (1) **behaviours grounded in observations** and (2) **memory-based reasoning about other agents**, ensuring success requires *genuine multi-agent cooperative reasoning*.