

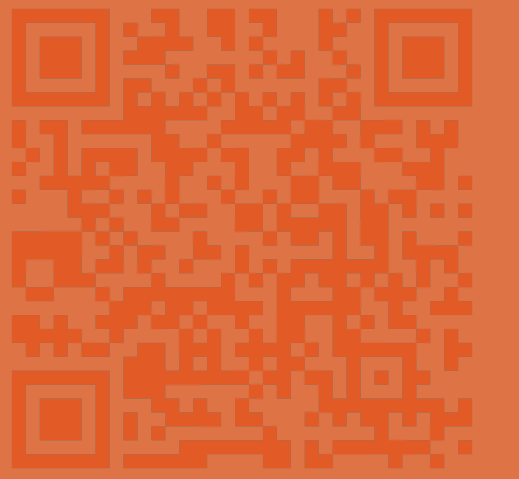
Factorised Active Inference for Strategic Multi-Agent Interactions

Jaime Ruiz-Serra, Patrick Sweeney, Michael Harré



THE UNIVERSITY OF
SYDNEY

Centre for Complex Systems,
The University of Sydney

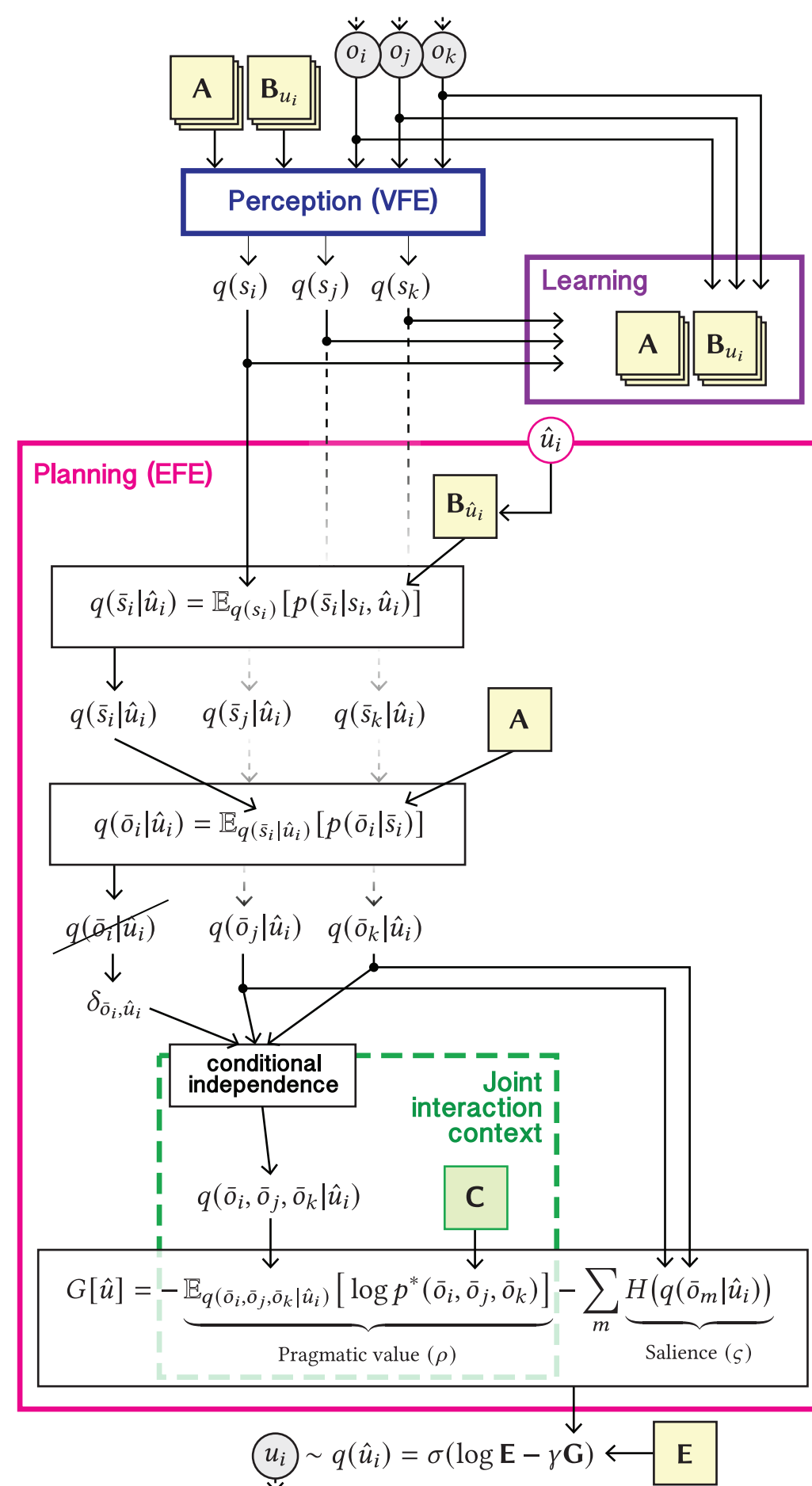


TL;DR

By factorising the generative model of AIF agents, they can maintain **individual beliefs** about others while planning strategically in a **joint context**. We employ **game transitions** to induce non-stationarity in agent preferences and study the resultant adaptive behaviour at the individual and collective levels. Bridging cognitive process models (AIF) with economic/MAS models (game theory) shows potential for understanding collective intelligence and designing interventions.

Methods

Factorised Beliefs: Ego maintains individual beliefs $q(s_n)$ about the hidden state s_n (e.g., propensity to cooperate, 'type') of each agent $n \in \{i, j, k\}$ as a separate factor.



Perception: Beliefs for each factor are updated based on observed actions $o = (o_i, o_j, o_k)$ by minimising VFE (negative ELBO)

$$F[q, o] = D_{KL}[q(s) || p(s|o)] - \log p(o)$$

Planning: ego evaluates counterfactual actions (\hat{u}_i) by calculating their EFE (pragmatic value, salience, novelty)

$$G[\hat{u}_i] = -\rho[\hat{u}_i] - \varsigma[\hat{u}_i] - \eta[\hat{u}_i]$$

Preferences are derived from the game payoff matrix (joint interaction context), $p^*(o_i, o_j, o_k) = \sigma(g(o_i, o_j, o_k))$

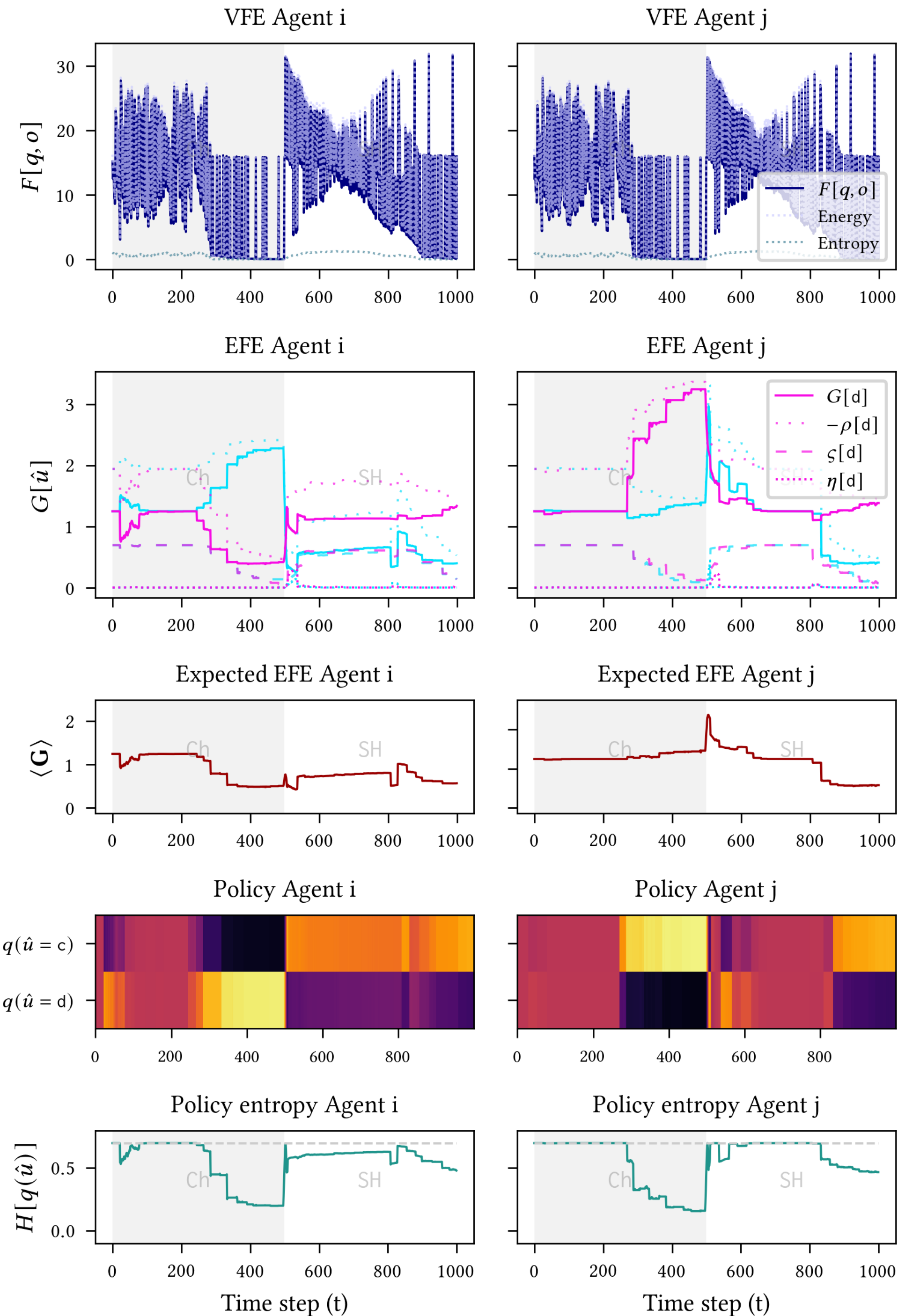
Requires predicting joint outcomes, $q(\bar{o}_i, \bar{o}_j, \bar{o}_k | \hat{u}_i)$

Pragmatic value becomes expected utility under predicted opponent actions,

$$\rho[\hat{u}_i] = \mathbb{E}_{q(\bar{o}_j)q(\bar{o}_k)}[\log p^*(\hat{u}_i, \bar{o}_j, \bar{o}_k)]$$

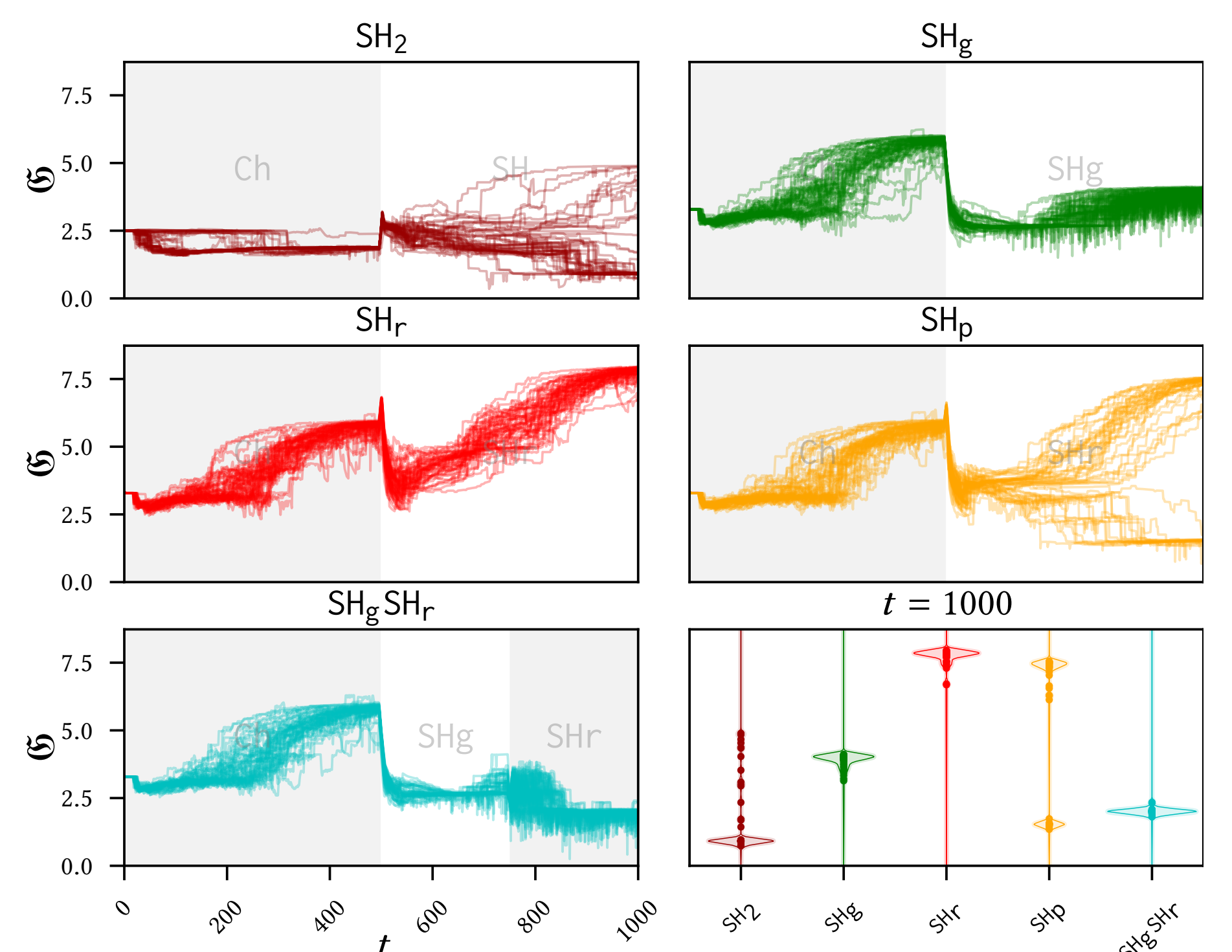
Results

VFE and strategic uncertainty



Ensemble-level EFE and equilibrium selection

$$\mathcal{G} = \sum_i \langle G \rangle^{(i)} = \sum_i \mathbb{E}_{q(\hat{u}_i)}[G[\hat{u}_i]]$$



Conclusions

1. \mathcal{G} dynamics characterise equilibria and attractor basins. Lower \mathcal{G} generally indicates 'better' collective outcomes (NE aren't always socially optimal).
2. Bifurcations in \mathcal{G} show convergence to different equilibria (e.g., payoff-dominant vs risk-dominant in SH) across trials. Shows relative basin size.
3. Game structure significantly impacts equilibrium selection (e.g., SH_g vs SH_r vs SH_p). Paradoxical results observed (requiring more cooperation sometimes led to less).
4. Strategic intervention possibility: Transitioning through a trust-building game (SH_g) can steer the collective to a better equilibrium more effectively than penalizing defection (SH_p).

Code: github.com/RuizSerra/factorised-MA-AIF

Contact: Jaime.RuizSerra@sydney.edu.au