







Empathy paper knowledge transfer

Sanjeev Namjoshi
12/16/25

Aims for active inference model

- Cooperation as emergent behavior through empathy
- Infer beliefs of other agent in response to your actions (theory of mind)
- Factor in emotional state of other agent (via valence/arousal calculated from VFE/EFE)
- Empathy parameter controlling degree to which other agent's beliefs are used for decision-making

Prisoners' dilemma

Prisoners' dilemma		prisoner B			
		confess 		remain silent 	
prisoner A	confess 	 5 years 5 years 0 year 20 years			
	remain silent 	 20 years 0 year 1 year 1 year			

© 2010 Encyclopædia Britannica, Inc.

An analytical model of active inference in the Iterated Prisoner's Dilemma

Daphne Demekas^{*1,2}, Conor Heins^{†1,3,4,5}, and Brennan Klein^{‡1,5,6}

Generative model

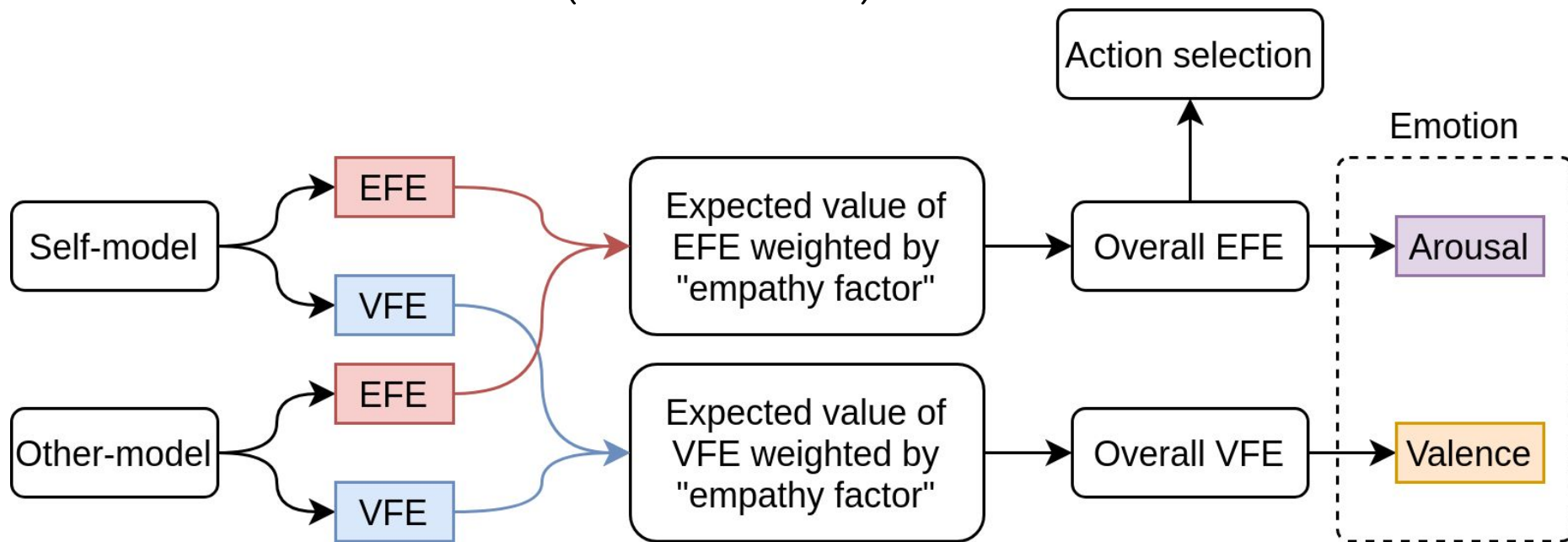
Variable Name	Notation
Hidden States	$\mathbf{s} \in \{\text{CC}, \text{CD}, \text{DC}, \text{DD}\}$
Observations	$\mathbf{o} \in \{\text{CC}, \text{CD}, \text{DC}, \text{DD}\}$
Actions	$\mathbf{u} \in \{u^C, u^D\}$
Observation Model	$P(\mathbf{o}_t \mathbf{s}_t; A) = \text{Cat}(\mathbf{A})$
Transition Model	$P(\mathbf{s}_{t+1} \mathbf{s}_{t-1}, \mathbf{u}_{t-1}; B) = \text{Cat}(\mathbf{B})$
Transition Model Parameter	$P(B) = \prod_{ju} P(B_{\bullet ju}), \quad P(B_{\bullet ju}) = \text{Dir}(\mathbf{b}_{\bullet ju})$
Initial State Prior	$P(\mathbf{s}_1; D) = \text{Cat}(\mathbf{D})$
‘Biased’ State Prior (Reward)	$\tilde{P}(\mathbf{s}; C) = \text{Cat}(\mathbf{C}), \quad \text{s.t. } \ln \mathbf{C} = [3, 1, 4, 2]$

Table 2: Generative model variables and notation.

Payoff matrix

		<u>Player 2</u>	
		Cooperate (C)	Defect (D)
<u>Player 1</u>	C	(3, 3)	(1, 4)
	D	(4, 1)	(2, 2)

Active inference model (current form)



Pattisapu et al. 2024

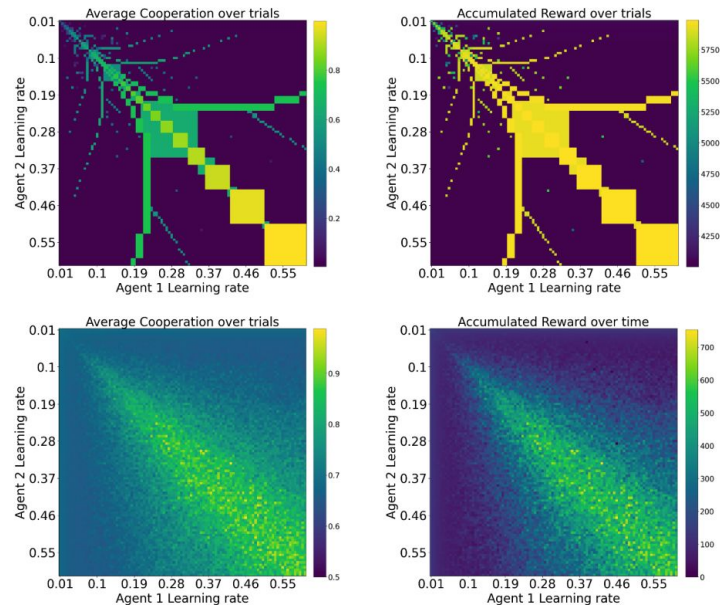
- Stochastic A matrix
- Learning for B matrix enabled
- Learning for C matrix (maybe)

Active inference model - major goals

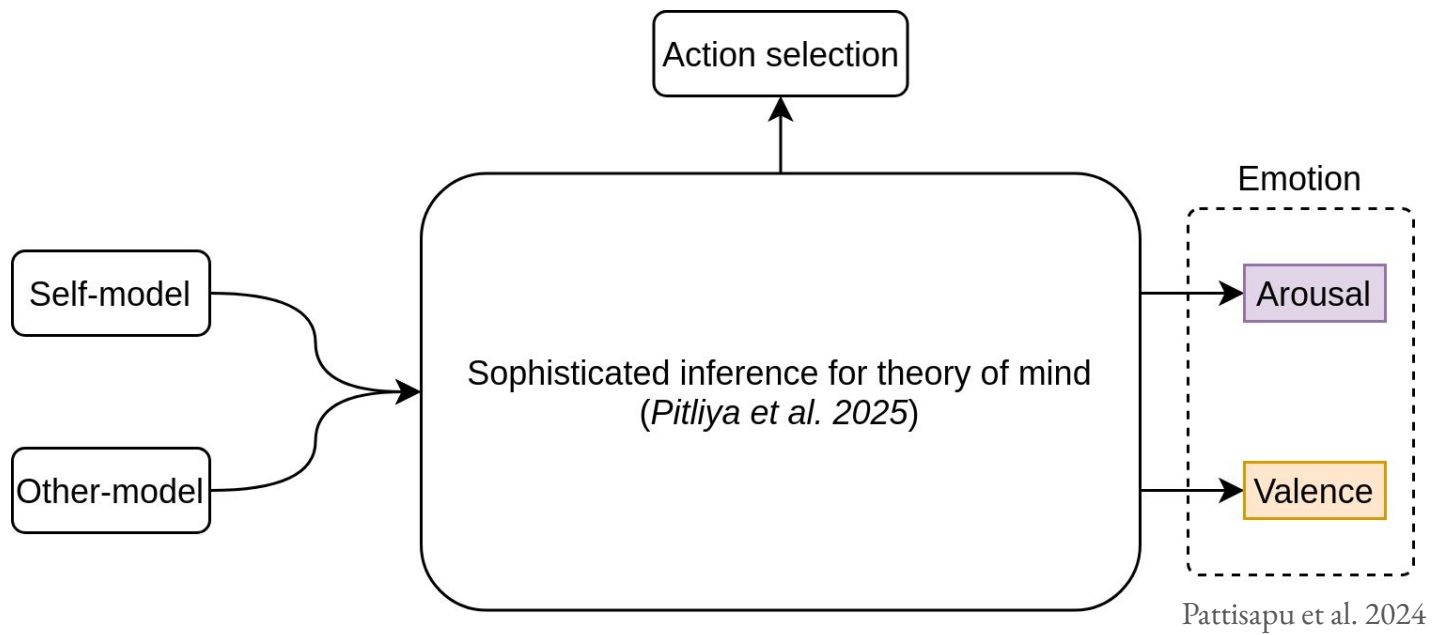
- Grid search over empathy parameter setting
- Find the conditions under which cooperation emerges

Future directions:

- Examine emotional state over time
- Compare results with and without B learning
- Examine how preference learning is affected by the empathy parameter
- Non-identity A matrix
- Expansion to n-player game



Active inference model (revised)



Active inference model (revised)

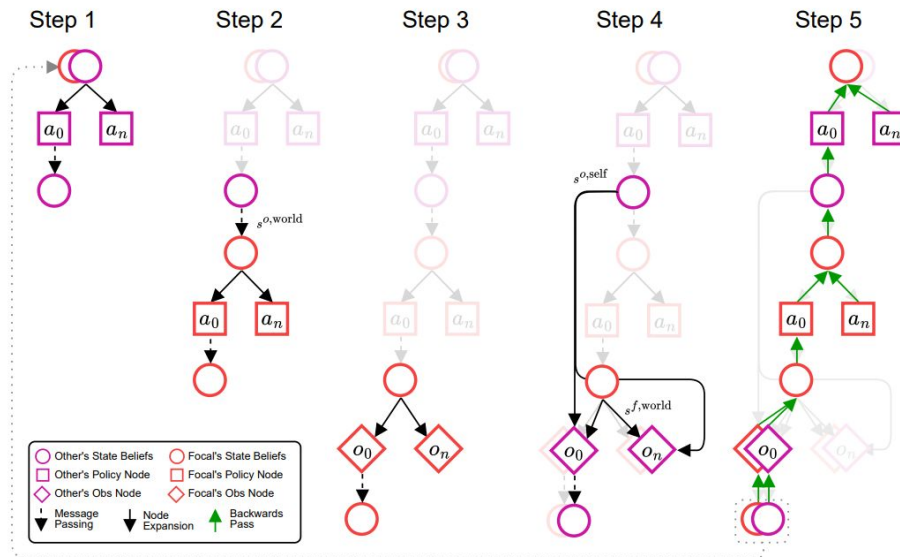


Fig. 1: **Recursive Planning Tree for Theory of Mind.** Red and purple represent the focal and other agent's nodes respectively. Circles indicate the agent's beliefs, squares indicate evaluated actions and diamonds indicated expected observations. For a detailed description of each step, see Section 2.3.

Active inference model (revised empathy factor)

$$\begin{aligned} G(o_{\tau}^f, o_{\tau}^o, a_{\tau}^f, a_{\tau}^o) = & \mathbb{E}_{Q(o_{\tau+1}^f, o_{\tau+1}^o | a_{\leq \tau}^f, a_{\leq \tau}^o)} \\ & [-\ln P(o_{\tau+1}^f | C^f) - \mathbb{D}_{\text{KL}}[Q(s_{\tau+1}^f | o_{\tau+1}^f) || Q(s_{\tau+1}^f)]] \\ & + \mathbb{E}_{Q(a_{\tau+1}^f | o_{\tau+1}^f) Q(a_{\tau+1}^o | o_{\tau+1}^o) Q(o_{\tau+1}^f, o_{\tau+1}^o | a_{\leq \tau}^f, a_{\leq \tau}^o)} \\ & [G(o_{\tau+1}^f, o_{\tau+1}^o, a_{\tau+1}^f, a_{\tau+1}^o)] \end{aligned} \quad (2)$$

- *Option 1*: Encode empathy factor by weighting EFE in the recursive expansion based on contribution from focal or other agent
- *Option 2*: Encode empathy factor through precision on observation/action expansion

Scoping options/difficulties

- Access to theory of mind code from Pitliya et al. 2025?
- Stop at emergence of cooperation under simplest scenario?
 - 2 player game, iterated
 - Empathy parameter
 - Identity A matrix (KL-control)
 - B matrix learning
 - No emotion component analysis
- Add any of the following?
 - Examine emotional state over time
 - Compare results with and without B learning
 - Examine how preference learning is affected by the empathy parameter
 - Non-identity A matrix
 - Expansion to n-player game

References

- **Demakes D, Heins C, and Klein B.** 2023. An analytical model of active inference in the Iterated Prisoner's Dilemma. arXiv:2306.15494v2
- **Pattisapu C, Verbelen T, Pitliya RJ, Kiefer AB, Albarracin M.** 2024. Free Energy in a Circumplex Model of Emotion. arXiv:2407.02474v1
- **Pitliya RJ, Catal O, Van de Maele T, Pezzato C, Verbelen T.** 2025. Theory of Mind Using Active Inference: A Framework for Multi-Agent Cooperation. arXiv:2508.00401v2