



Deposited via The University of Leeds.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/241021/>

Version: Accepted Version

Article:

Barfuss, W., Tittel, P. and Mann, R.P. (Accepted: 2026) Learning Together, Better, Faster - On the Timescales of Learning Social Learning Strategies. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*. ISSN: 1364-503X (In Press)

This is an author produced version of an article accepted for publication in *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, made available via the University of Leeds Research Outputs Policy under the terms of the Creative Commons Attribution License (CC-BY), which permits unrestricted use, distribution and reproduction in any medium, provided the original work is properly cited.

Reuse

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here: <https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.

LEARNING TOGETHER, BETTER, FASTER

ON THE TIMESCALES OF LEARNING SOCIAL LEARNING STRATEGIES

Wolfram Barfuss 

Transdisciplinary Research Area Sustainable Futures, University of Bonn, 53115 Bonn, Germany
Center for Development Research, University of Bonn, 53113 Bonn, Germany
Institute for Food & Resource Economics, University of Bonn, 53115 Bonn, Germany
wbarfuss@uni-bonn.de

Paula Tittel

Faculty of Mathematics and Physics, Leibniz University Hannover, 30167 Hannover, Germany

Richard P. Mann 

Department of Statistics, School of Mathematics, University of Leeds, Leeds, LS2 9JT, United Kingdom

April 24, 2026

ABSTRACT

1 Social learning strategies are the foundational mechanism for human-machine cultural evolution.
2 For decades, they have been studied either as fixed evolutionary heuristics or as rational Bayesian
3 solutions, yet neither approach explains their remarkable flexibility nor how boundedly rational
4 agents could acquire them. Recent work has begun to explain social learning as emerging from
5 domain-general reward-based learning, yet, left open to address the timescales required to acquire
6 them. Here, we investigate how quickly agents learn to rely on social information under different
7 levels of environmental uncertainty, using boundedly rational reinforcement learning dynamics. We
8 find that learning together not only yields better results under environmental uncertainty but can also
9 significantly speed up learning, provided that agents have a reliable internal model of the environment.
10 When such an internal model is not available, learning together still yields better results, but the
11 learning speed advantage is lost. However, even model-free learning does not take significantly
12 longer than independent learning to achieve most of the final reward that an independently learning
13 agent achieves. Our findings suggest a concrete design principle for hybrid collective intelligence:
14 the quality of the internal world models determines whether collaborative learning accelerates or
15 decelerates.

16 **Keywords** Social learning strategies • Internal timescales • Environmental uncertainty • Reinforcement learning •
17 Nonlinear dynamics

18 1 Introduction

19 The increasing integration of advanced artificial agents into everyday human activities presents a fundamental challenge
20 to our understanding of collective behavior [Rahwan et al., 2019, Brinkmann et al., 2023]. For example, when deciding
21 on a restaurant to eat at, a hotel to stay in, or a product to buy, humans often consult online reviews and ratings from
22 other users. These reviews nowadays are frequently summarized, generated, and moderated by AI systems on a much
23 faster timescale. Likewise, they are also used by advanced AI agents to make and execute decisions on behalf of their
24 human users, creating complex feedback loops. As humans and AI begin to learn, cooperate, and make decisions
25 in shared environments, we must re-examine the principles that underpin sociality itself. How do adaptive social
26 learning strategies emerge and evolve when humans and AI agents must learn from one another to navigate uncertain
27 worlds? In this article, we show that these rules governing when and from whom to learn emerge endogenously through
28 domain-general reward learning, and we quantify for the first time how quickly this emergence occurs under different
29 environmental and cognitive conditions.

30 To understand how populations adapt, the field of cultural evolution offers a robust theoretical framework for modeling
 31 the transmission and refinement of behaviors [Boyd and Richerson, 1988][Rogers, 1988]. This approach is invaluable
 32 for studying the evolution of sociality, as it provides a mathematical basis for analyzing how individuals choose between
 33 learning from their own experience and learning from others using various social learning strategies [Acerbi et al.,
 34 2022]. These are assumed to be fixed heuristics tailored to particular situations that individuals employ in learning
 35 from others. At least 25 distinct social learning strategies, such as ‘copy the majority’ or ‘copy the successful’, have
 36 been documented in various studies [Laland, 2004, Rendell et al., 2010, Morgan et al., 2011, Kendal et al., 2018]. This
 37 approach has transformed the study of cultural evolution from a descriptive analogy to a rigorous field of theoretical
 38 and empirical research [Deffner et al., 2024]. However, this fixed heuristics account treats social learning and individual
 39 (or ‘asocial’) learning as conceptually distinct, alternative processes. As such, it cannot explain the origins of these
 40 strategies, i.e., the cognitive and environmental factors that give rise to them [Heyes, 2016, Tump et al., 2024].

41 Rational decision-making models address this limitation by proposing that social learning arises from a single inferential
 42 process [Perreault et al., 2012]. Instead of applying a fixed heuristic, individuals use Bayesian inference to rationally
 43 update their beliefs about the world, combining social and nonsocial cues [Pérez-Escudero and de Polavieja, 2011,
 44 Arganda et al., 2012, Bikhchandani et al., 1992][Banerjee, 1992]. This allows agents to rationally weigh social
 45 information according to its reliability and predictive value, leading to social learning strategies that emerge from
 46 fundamental information asymmetries [Arganda et al., 2012, Mann, 2018, 2020]. For example, individuals place a
 47 heavy weight on social cues when the environment changes slowly or when its state cannot be well predicted using
 48 nonsocial cues [Perreault et al., 2012, Mann, 2020]. Rational models of Bayesian inference also predicted people’s
 49 social learning behaviour better than simpler fixed heuristics [Taylor-Davies et al., 2025]. However, these rational
 50 models still have a significant limitation. They excel at explaining an optimal social learning strategy at equilibrium,
 51 but do not address the process of how the agent initially acquired it. They therefore cannot answer questions about the
 52 timescale for acquiring social learning strategies or the necessary conditions for their emergence. Thus, they cannot
 53 account for the adaptability of social learning to changing environmental contexts, as recent experimental findings show
 54 [Efferson et al., 2008, Toelch et al., 2014, Deffner et al., 2020, Toyokawa and Gaissmaier, 2022].

55 To address this gap, an emerging consensus points to the underlying mechanism of domain-general reward learning
 56 [Heyes, 2012, Heyes and Pearce, 2015]. This view posits that seemingly complex social learning strategies are not
 57 fixed rules but emerge from a more fundamental process: associating social cues with expected rewards. Indeed,
 58 evidence from neuroscience [Behrens et al., 2008, Olsson et al., 2020, Zhang and Gläscher, 2020], computational
 59 cognitive science [Suganuma et al., 2025, Najar et al., 2020, Bergerot et al., 2024, Witt et al., 2024, Schultner et al.,
 60 2025a][Danwitz and von Helversen, 2025], and machine learning research [Borsa et al., 2019, Ndousse et al., 2021,
 61 Ha and Jeong, 2023] supports the view that reward learning is central to the emergence of social learning strategies.
 62 This body of work strongly suggests that the sophisticated heuristics of social learning are not pre-programmed but are
 63 themselves learned and refined through reward. The key advantage of this approach is its adaptation timescale. Learning
 64 from experience can happen flexibly within an individual’s lifetime. By contrast, the fixed-heuristic approach typically
 65 models the evolutionary success of a limited set of fixed heuristics under natural selection [Rendell et al., 2010, Aoki
 66 and Feldman, 2014, Turner et al., 2023, Sigalou and Mann, 2023]. As learning operates on a much faster timescale than
 67 natural selection, human social learning strategies are supposedly more flexible than is commonly assumed [Schultner
 68 et al., 2025b].

69 However, a systematic investigation of the adaptation timescales remains elusive. It is well known that environmental
 70 uncertainty promotes social learning because social cues provide valuable information that complements ambiguous
 71 environmental signals [Perreault et al., 2012, Mann, 2018][Wu et al., 2025]. Concerning learning timescales, Suganuma
 72 et al. [2025] recently found differences in the internal learning speed of different social reinforcement learning
 73 algorithms, such as value shaping and decision biasing, that rely invariably on social information. But how does
 74 environmental uncertainty affect the rate of acquiring social learning strategies compared to independent learning
 75 from a pure domain-general reward learning framework, where agents can flexibly choose whether or not to integrate
 76 social information [Schultner et al., 2025a]? Evidence from machine learning research suggests that the mere presence
 77 of an expert teacher agent not only consistently improves the final policy of a learning agent but also accelerates
 78 the convergence of learning [Borsa et al., 2019]. However, what if there is no expert teacher to begin with? A
 79 reinforcement learning agent that incorporates social cues must learn more complex strategies than one that relies solely
 80 on environmental cues. When two reinforcement learning agents learn simultaneously from environmental signals and
 81 from each other, achieving higher performance at convergence may come at the cost of slower learning.

82 In this paper, we demonstrate that this intuition is largely wrong. To do so, we investigate how environmental uncertainty
 83 affects both the final performance and the learning speed of emergent social learning strategies within a reward-based
 84 learning framework. Specifically, we ask: (1) Under what environmental conditions do learned social learning strategies
 85 outperform independent learning? (2) What types of social learning strategies emerge from reward learning under

86 different levels of uncertainty? (3) How quickly can agents acquire effective strategies, and does learning together
87 accelerate or hinder this process?

88 To address these questions, we develop a mathematical framework based on partially observable stochastic games, where
89 agents must make decisions based on noisy observations of an environmental state. We operationalize environmental
90 uncertainty in two dimensions: i) observational uncertainty and ii) state uncertainty. We compare single-agent scenarios,
91 in which individuals learn only from environmental feedback, with two-agent scenarios, in which a second agent can
92 observe and potentially learn from the first agent’s choices. We examine both model-based learning dynamics (in which
93 agents behave as if they possess an internal model of the environment) and model-free dynamics (in which agents learn
94 purely through trial-and-error experience).

95 Our analysis reveals four key findings. First, we demonstrate that social learning strategies consistently emerge and
96 outperform independent learning across a wide range of environmental uncertainties, with benefits reaching up to 14%
97 in reward performance. Second, we show that the interplay between observational and environmental state uncertainty
98 gives rise to a rich diversity of adaptive social learning strategies—from simple ‘copy-when-uncertain’ heuristics
99 to complex conditional strategies that depend on which environmental state is more likely. Third, we find that the
100 transitions between these learned social learning strategies are critical, marked by a significant slowdown in learning
101 at these points. Fourth, we find that away from these critical points, agents with internal environmental models learn
102 approximately twice as fast when paired together as when learning alone, whereas agents without such models show
103 the opposite pattern. Yet remarkably, the time required for learning together to match almost the final performance of
104 independent learners remains comparable, even without internal models.

105 These findings have important implications for understanding collective behavior in hybrid human-AI populations.
106 As AI systems increasingly mediate human decision-making and learning through content curation, recommendation
107 algorithms, and automated decision support, our results suggest that the speed and efficiency of collective adaptation
108 will depend critically on whether these systems can build reliable internal models of their environments. More broadly,
109 our work demonstrates that the sophisticated social learning strategies observed across human and animal societies need
110 not be innate or evolved over generations—they can emerge rapidly within individual lifetimes through domain-general
111 reward-learning mechanisms.

112 2 Methods

113 2.1 Mathematical framework

114 We formulate our model in the general framework of partially observable stochastic games $G = \langle \mathcal{N}, \mathcal{S}, \mathcal{A}, \mathcal{O}, T, R, O \rangle$.
115 There are $\mathfrak{N} \in \mathbb{N}$ agents. The choice environment consists of $\mathfrak{S} \in \mathbb{N}$ states, $\mathcal{S} = (\mathcal{S}_1, \dots, \mathcal{S}_{\mathfrak{S}})$. In each state s , each
116 agent $i \in \mathcal{N} = (1, \dots, N)$ has $\mathfrak{A} \in \mathbb{N}$ available actions $\mathcal{A}^i = (\mathcal{A}_1^i, \dots, \mathcal{A}_{\mathfrak{A}}^i)$ to choose from. We denote the joint action
117 set by $\mathcal{A} = \otimes_i \mathcal{A}^i$, where \otimes_i is a cross-product sign. A joint action by all agents is denoted by $a = (a^1, \dots, a^N) \in \mathcal{A}$,
118 and $a^{-i} = (a^1, \dots, a^{i-1}, a^{i+1}, \dots, a^N)$ denotes the joint action by all agents except agent i . Agents choose their actions
119 simultaneously.

120 We use the same number of actions, \mathfrak{A} , and observations, \mathfrak{D} , (see below) across all agents and states for notational and
121 computational convenience. Doing so is computationally convenient, as it allows us to formulate the environmental
122 state transition, reward, and observation functions, T , R , and O , as tensors, thereby enabling accelerated linear algebra
123 (XLA) operations¹.

124 The transition tensor $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ determines the probability of a state change. $T^{s,a,\acute{s}}$ is the probability to
125 transition from the current state s under the joint action a to the next state \acute{s} .

126 The reward tensor $R : \mathcal{N} \times \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ encodes the immediate rewards the agents receive. $R^{i,s,a,\acute{s}}$ is the reward
127 of agent i when the environment transitions from state s , under the joint action a , to state \acute{s} .

128 Instead of observing the states $s \in \mathcal{S}$ directly, each agent i observes one of $\mathfrak{D} \in \mathbb{N}$ observations $\mathcal{O}^i = (\mathcal{O}_1^i, \dots, \mathcal{O}_{\mathfrak{D}}^i)$
129 according to its observation tensor $O^i : \mathcal{S} \times \mathcal{O}^i \rightarrow [0, 1]$. O^{i,s,o^i} is the probability that agent i observes observation
130 $o^i \in \mathcal{O}^i$ when the environment is in state s . The joint observation set is $\mathcal{O} = \otimes_i \mathcal{O}^i$ and $O = \otimes_i O^i : \mathcal{N} \times \mathcal{S} \times \mathcal{O} \rightarrow$
131 $[0, 1]$ is the joint observation tensor.

132 Agents select actions according to their strategies $X^i : \mathcal{O}^i \times \mathcal{A}^i \rightarrow [0, 1]$ (also known as policies). X^{i,o^i,a^i} is the
133 probability that agent i selects action a^i when it observes observation o^i . The joint strategy tensor is $X = \otimes_i X^i :$
134 $\mathcal{N} \times \mathcal{O} \times \mathcal{A} \rightarrow [0, 1]$.

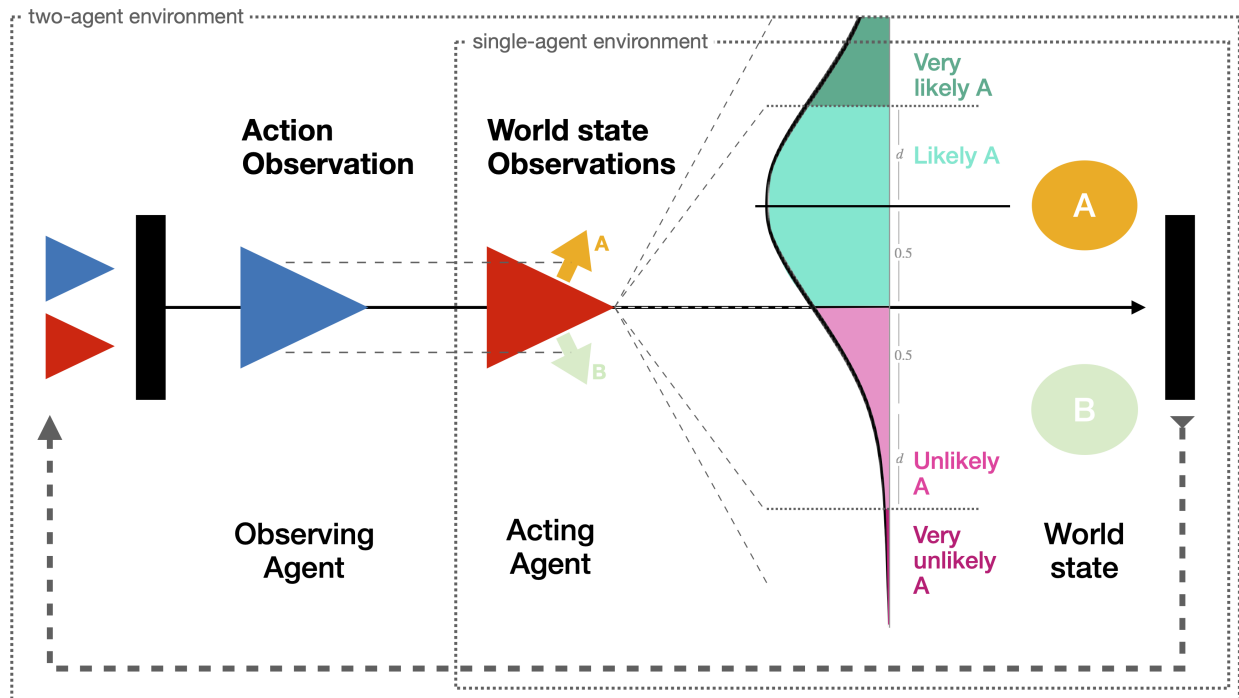


Figure 1: **Choice environments.** In the single-agent environment, the agent observes a noisy signal about the state of nature, A or B. If the agent acts in accordance with the state of nature, it receives a high reward; otherwise, it gets a low reward. The two-agent environment builds upon the single-agent environment. A second agent observes the choice taken by the first agent. Thus, the second agent can learn to condition its own action choices on the other agent’s choices. Which agent is first is decided by a coin flip after each trial.

135 2.2 Choice environment

136 We consider two choice environments: a single-agent environment and a two-agent environment (Figure 1). In the
 137 single-agent environment, the agent observes a noisy signal about the state of nature. If the agent acts in accordance
 138 with the state of nature, it receives a high reward; otherwise, it gets a low reward. The two-agent environment builds
 139 upon the single-agent environment. A second agent observes the choice taken by the first agent. Thus, the second agent
 140 can learn to condition its own action choices on the other agent’s choices.

141 **Single-agent environment.** Specifically, the single-agent environment consists of two states of nature, $\mathcal{S}_{\text{single}} = \{A, B\}$,
 142 and the agent can choose between two actions $\mathcal{A}^1 = \{A, B\}$. If the agent chooses action A (B), when the environment
 143 is in state A (B), the agent receives a reward $R^{1,A,A,\acute{s}} = R^{1,B,B,\acute{s}} = 1.0$, regardless of the next state \acute{s} . Otherwise, the
 144 agent receives a reward $R^{1,A,B,\acute{s}} = R^{1,B,A,\acute{s}} = -1.0$. However, the agent cannot directly observe the state of nature.

145 **Observational uncertainty.** The agent observes one of four possible signals, which we denote by $\mathcal{O}^1 = \{A, a, b, B\}$.
 146 Here, the signals A and a denote that nature is more likely to be in state A, whereas signals B and b denote that nature
 147 is more likely to be in state B. Capital letters indicate more certainty than lower-case letters regarding the respective
 148 signal-state pairs. Thus, observing A indicates that nature is in A is *very likely*, whereas observing a means that nature
 149 is only *likely* in A. For state B, the analogous applies. Note that the agent has no capacity to decipher the meaning of
 150 these letters. It only learns via observation-reward associations.

151 We model these observation probabilities by discretizing a normal distribution $\mathcal{N}_{\mu,\sigma}$, which is centered around the true
 152 state of nature ($\mu = 0$). The other state is imagined to reside one unit length apart. Thus, if the true state of nature is A
 153 (B),

- 154 • the agent observes signal A (B) with probability $O^{1,A,A} = O^{1,B,B} = \int_{-\infty}^{-d} \mathcal{N}_{0,\sigma}(x) dx$;
- 155 • it observes signal a (b) with probability $O^{1,A,a} = O^{1,B,b} = \int_{-d}^{0.5} \mathcal{N}_{0,\sigma}(x) dx$;

¹<https://openvla.org/>

- 156 • it observes signal b (a) with probability $O^{1,A,b} = O^{1,B,a} = \int_{0.5}^{1+d} \mathcal{N}_{0,\sigma}(x) dx$; and last
- 157 • it observes signal B (A) with probability $O^{1,A,B} = O^{1,B,A} = \int_{1+d}^{\infty} \mathcal{N}_{0,\sigma}(x) dx$,

158 where $d > 0$ is the distance from the center of nature’s true state at which the agent believes to observe the state *very*
 159 likely. Without loss of generality, we set $d = 0.5$ throughout this paper. We use the standard deviation of the normal
 160 distribution, σ , to represent the observational uncertainty of our model.

161 **State uncertainty.** We model the next state of the world as independent of the current state and actions. It is determined
 162 by a parameter p_A , denoting the probability of transitions to nature state A, i.e., $T^{s,a,A} = p_A$ and $T^{s,a,B} = 1 - p_A$. We
 163 will investigate the range $p_A > 0.5$. Thus, we denote the state uncertainty of our model by the probability that nature
 164 is in state B, $p_B = 1 - p_A$. This way, when $p_B = 1 - p_A = 0.5$, the environment is maximally uncertain, as both
 165 states occur equally likely. By contrast, when p_B approaches 0, nature becomes more predictable, and state uncertainty
 166 decreases. For example, taking the bus might be the best way to get to work on 9 out of 10 days. But you could also
 167 walk if you receive a strong enough signal that the bus will be severely delayed, e.g., when you see snow and ice on the
 168 streets. Note that most other works, e.g., [Mann, 2018, 2020, Sigalou and Mann, 2023], have not really explored this
 169 kind of uncertainty and always assumed $p_A = p_B = 0.5$ for convenience.

170 These two dimensions map onto distinct features of empirical social learning environments, and their interplay
 171 determines when observing others is most valuable. Consider a traveler exploring a new city whose culinary scene spans
 172 Chinese, Greek, Italian, Mexican, and Vietnamese cuisine, among others. Each day, the traveler narrows the choice to
 173 two restaurants. The hidden state of the world captures which of the two is the genuine better fit for the traveler. This is
 174 determined by the traveler’s personal tastes and the intrinsic quality of each restaurant. State uncertainty (p_B) reflects
 175 how predictable this fit is: when the traveler has a strong preference for one cuisine over the other, one option is reliably
 176 the better choice (low p_B); when the two cuisines are equally appealing, the better choice is unpredictable ($p_B \approx 0.5$).
 177 Observational uncertainty (σ) captures how clearly available cues (e.g., aromas drifting from the kitchen, the displayed
 178 menu, the restaurant’s atmosphere) signal which option is the better fit on this day. Low observational uncertainty
 179 means these cues are strong and diagnostic. High observational uncertainty means they are weak or ambiguous. The
 180 two dimensions have distinct consequences for social learning: when cues are ambiguous (high σ), a fellow traveler’s
 181 restaurant choice provides a valuable independent signal (assuming similar tastes). When one option almost always fits
 182 better (low p_B), even imperfect private cues tend to identify the right option, reducing the marginal value of copying
 183 others.

184 **Two-agent environment.** The two-agent environment adds a second agent to the uncertain decision problem described
 185 above. The natural state, $\mathcal{S}_{\text{single}} = \{A, B\}$, is complemented by a social component that describes the sequence of
 186 agents’ decision points. Agents choose sequentially. The first-acting agent can rely only on the natural observation signal
 187 to decide on an action. The subsequent agent observes this action along with the natural observation signal to decide on
 188 its action. After the second choice is made, both agents receive their rewards, and the process restarts, with the agents’
 189 order determined at random. The social state set describes this agent sequence, $\mathcal{S}_{\text{social}} = \{12., .1\underline{A}, .1\underline{B}, 21., .2\underline{A}, .2\underline{B}\}$.
 190 The first letter of each state identifier denotes the non-acting agent $i \in \mathcal{N} = \{1, 2\}$. The second letter describes the
 191 acting agent. The last letter indicates the action choice of the first-acting agent $a \in \{\underline{A}, \underline{B}\}$. If there is no such action,
 192 or no non-acting agent, a ‘.’ signifies an empty spot. Thus, the combined social and natural state set of the two-agent
 193 environment consists of 12 states,

$$\mathcal{S}_{\text{two}} = \mathcal{S}_{\text{single}} \times \mathcal{S}_{\text{social}} = \{A12., A.1\underline{A}, A.1\underline{B}, A21., A.2\underline{A}, A.2\underline{B}, B12., B.1\underline{A}, B.1\underline{B}, B21., B.2\underline{A}, B.2\underline{B}\}.$$

194 The observation set of the two-agent environment combines the noisy observations of the natural state with the
 195 observation about the first-acting agent plus a dummy observation, ‘.’, for the non-acting agent,

$$\mathcal{O}_{\text{two}}^i = \{A., a., b., B., \underline{AA}, \underline{aA}, \underline{bA}, \underline{BA}, \underline{AB}, \underline{aB}, \underline{bB}, ..\} \text{ for } i = 1, 2,$$

196 where the first letter of each observation identifier denotes the natural observational signal, and the second letter
 197 represents the action taken by the first-acting agent. If there is no such action because it’s the first-acting agent’s turn, a
 198 ‘.’ signifies an empty spot.

199 In the two-agent environment, each round consists of two sequential actions before rewards are distributed, so each
 200 agent receives reward once per round rather than once per action. To make average rewards comparable across the two
 201 environments, we multiply rewards in the two-agent environment by 2 to compensate for the sparser feedback. This
 202 rescaling encodes the assumption that social observation is temporally free: the time the second agent spends waiting
 203 within a round carries no real-time cost. This is the natural baseline in many real-world social learning contexts, where

204 social cues are largely ubiquitous and available passively without dedicated observation effort. Under this convention,
 205 iteration counts of the learning rule measure cognitive convergence cost rather than individual sample cost.

206 2.3 Learning dynamics

207 We consider deterministic reinforcement learning dynamics [Bloembergen et al., 2015, Barfuss and Mann, 2022], which
 208 provide a transparent and analytically tractable framework for studying the collective behavior of individual learners
 209 [Barfuss et al., 2025]. For example, they allow for a straightforward definition of convergence, enabling the analysis
 210 of learning timescales. With a formalism similar to that of evolutionary game theory [Börgers and Sarin, 1997, Tuyls
 211 et al., 2003, Sato and Crutchfield, 2003][Bernasconi et al., 2025], our reinforcement learning model boils down to the
 212 exponential replicator dynamics in discrete time [Hofbauer and Sigmund, 1998],

$$X_{t+1}^a \propto X_t^a \exp(\alpha F_{X_t}^a).$$

213 The probability of choosing action a at time step t , X_t^a , is updated proportionally to its exponential fitness $F_{X_t}^a$ under the
 214 current strategy profile X_t , scaled by the learning rate α . Without loss of generality, we set the learning rate $\alpha = 0.25$
 215 for all agents and environments throughout this article. In the Supplementary Information, we demonstrate that our main
 216 results are robust to varying $\alpha \in \{0.125, 0.25, 0.5\}$. As one would expect, the smaller α , the longer the convergence
 217 time becomes. This dynamic captures how agents adjust their action preferences based on the rewards they receive,
 218 leading to the evolution of their strategies over time. Action probabilities of high-performing actions increase, whereas
 219 those of underperforming actions decrease. High-performing actions are *replicated* more than underperforming ones.

220 As we seek to investigate the learning timescales across different levels of environmental uncertainty, we consider two
 221 types of agents: model-based and model-free reinforcement learning agents. Model-based agents operate as if they
 222 possess an internal model of the environment as they perceive it, allowing them to plan their actions by simulating
 223 future observations and rewards. In contrast, model-free agents learn to associate actions only through experienced
 224 rewards via trial-and-error interactions with the environment [Sutton and Barto, 2018]. Whether an agent has an internal
 225 model is likely to influence how quickly it can acquire effective (social learning) strategies.

226 As our deterministic dynamics operate on a higher level of cognitive abstraction than typical reinforcement learning
 227 algorithms [Barfuss, 2022], we do not explicitly model the internal processes of planning or value function estimation.
 228 Instead, we capture the essence of model-based and model-free learning by directly computing fitness values that drive
 229 the learning updates. Our two dynamics should, therefore, be seen as idealized representations of model-based and
 230 model-free learning processes, respectively, where the model-based dynamics operate as if they have a perfect internal
 231 model of the current environment they perceive.

232 **Model-based learning dynamics.** The update of the agents’ strategies of our model-based learning dynamics yields,

$$X_{t+1}^{i,o^i,a^i} = \frac{X_t^{i,o^i,a^i} \exp(\alpha^i R_{X_t}^{i,o^i,a^i})}{\sum_{b \in \mathcal{A}^i} X_t^{i,o^i,b} \exp(\alpha^i R_{X_t}^{i,o^i,b})},$$

233 where α^i denotes agent i ’s learning rate and $R_{X_t}^{i,o^i,a^i}$ is the average reward agent i receives when it observes observation
 234 o^i and takes action a^i , given that all agents behave according to the current strategy X_t . Thus, the ‘fitness’ of action
 235 a^i under observation o^i is simply the expected reward for taking that action in that observational context, even if that
 236 context is rarely observed or that action is seldom taken. Therefore, we call this update rule *model-based*.

237 To perform the learning update, we need to compute the average rewards $R_{X_t}^{i,o^i,a^i}$. However, the environment’s rewards
 238 $R^{i,s,a,s}$ depend on the states s , not the observations o^i . To obtain the average observation-action rewards $R_{X_t}^{i,o^i,a^i}$,
 239 we require a mapping from observations to states to account for the states s underlying a given observation o^i . The
 240 observation tensor O^{i,s,o^i} accounts for the observations o^i following from a state s . Following Barfuss and Mann
 241 [2022], we can invert the observation tensor into a belief tensor using Bayes’ rule,

$$B_X^{i,o^i,s} = \frac{O^{i,s,o^i} P_X^s}{\sum_s O^{i,s,o^i} P_X^s},$$

242 where P_X^s is the stationary state distribution given the strategy X . Thus, $B_X^{i,o^i,s}$ is the belief of agent i , that the
 243 environment is in state s , given it observes observation o^i . The stationary state distribution, P_X^s , is the left eigenvector
 244 of the strategy-average transition matrix $T_X^{s,\hat{s}}$. It can be obtained via

$$T_X^{s,\hat{s}} = \sum_{a^j \in \mathcal{A}^j} \prod_{j \in \mathcal{N}} Y_X^{j,s,a^j} T^{s,a,\hat{s}},$$

245 where Y_X^{j,s,a^j} are the effective state-action strategies. They can be obtained by averaging out the observations,

$$Y_X^{j,s,a^j} = \sum_{o^j \in \mathcal{O}^j} O^{j,s,o^j} X^{j,o^j,a^j}.$$

246 With Y_X^{j,s,a^j} we can obtain the effective transition matrix, $T^{s,\hat{s}}$, the stationary state distribution as its left eigenvector,
 247 P_X^s , and from that, the belief tensor $B^{i,o^i,s}$. With $B^{i,o^i,s}$, we can finally compute the strategy-average observation-action
 248 rewards R_X^{i,o^i,a^i} .

249 Whenever agent i observes observation o^i , the environment is state s with probability $B^{i,o^i,s}$. In s , all other agents
 250 $j \neq i$ behave according to Y_X^{j,s,a^j} . The environment transitions to state \hat{s} with probability $T^{s,a,\hat{s}}$. And agent i receives
 251 reward $R^{i,s,a,\hat{s}}$. Thus, the strategy-average observation-action reward for action a^i under observation o^i reads,

$$R_X^{i,o^i,a^i} = \sum_{s \in \mathcal{S}} \sum_{a^j \in \mathcal{A}^j} \sum_{\hat{s} \in \mathcal{S}} \sum_{j \neq i} B_X^{i,o^i,s} Y_X^{j,s,a^j} T^{s,a,\hat{s}} R^{i,s,a,\hat{s}}.$$

252 **Model-free learning dynamics.** In addition to the model-based learning dynamics described above, we also consider
 253 model-free learning dynamics. In model-free learning, agents do not maintain an internal model of the environment's
 254 transition, reward, and observation tensors. Hence, they cannot use the average rewards R_X^{i,o^i,a^i} to update their strategies.
 255 They have to be taken proportionally to the agents' experience. The agents' experience is captured by multiplying the
 256 average rewards with the stationary observation distribution P_X^{i,o^i} , and the strategy X^{i,o^i,a^i} . Thus, the reinforcement
 257 signal for an action a^i under an observation o^i is proportional to the average reward of that action times the frequency
 258 with which the action is taken under the observation o^i times the frequency of observing o^i . Therefore, the model-free
 259 learning dynamics read,

$$X_{t+1}^{i,o^i,a^i} = \frac{X_t^{i,o^i,a^i} \exp\left(\alpha^i P_{X_t}^{i,o^i} X_t^{i,o^i,a^i} R_{X_t}^{i,o^i,a^i}\right)}{\sum_{b \in \mathcal{A}^i} X_t^{i,o^i,b} \exp\left(\alpha^i P_{X_t}^{i,o^i} X_t^{i,o^i,b} R_{X_t}^{i,o^i,b}\right)}.$$

260 The stationary observation distribution P_X^{i,o^i} can be obtained by averaging the observation tensor O^{i,s,o^i} over the
 261 stationary state distribution P_X^s ,

$$P_X^{i,o^i} = \sum_{s \in \mathcal{S}} O^{i,s,o^i} P_X^s.$$

262 **Relation to standard reinforcement learning algorithms.** The terminology model-based and model-free spans a
 263 wide algorithmic space, including planning methods such as Dyna and temporal-difference methods such as Q-learning
 264 [Sutton and Barto, 2018], and our dynamics should not be equated with any specific algorithm within these families.
 265 The property our dynamics *do* capture is the core informational distinction: whether the agent uses an internal model of
 266 the environment's transition, reward, and observation structure when computing its update signal. Our model-based
 267 dynamics do so explicitly, constructing Bayesian beliefs about the hidden state (see below) and computing expected
 268 rewards over the full joint dynamics. Our model-free dynamics instead weight updates only by empirically determined
 269 observation visit frequencies P_X^{i,o^i} , as an agent without access to the underlying state structure would. The properties
 270 our dynamics do *not* capture are algorithm-specific implementation details. The dynamics are deterministic, computing
 271 exact expectations over the stationary distribution rather than sampling individual transitions; this corresponds to the
 272 mean-field or population limit of stochastic RL algorithms [Börger and Sarin, 1997, Tuyls et al., 2003]. Neither
 273 dynamics explicitly maintain value functions, perform multi-step planning, or use temporal-difference bootstrapping.
 274 This idealization isolates the informational mechanism from implementation details, providing a clean theoretical
 275 reference point for what each learning paradigm can achieve in this environment.

276 3 Results

277 We will first investigate how the final performance differences between individual and emergent social learning depend
 278 on environmental uncertainty and the learning model. Second, we analyze the learned strategies that underlie these

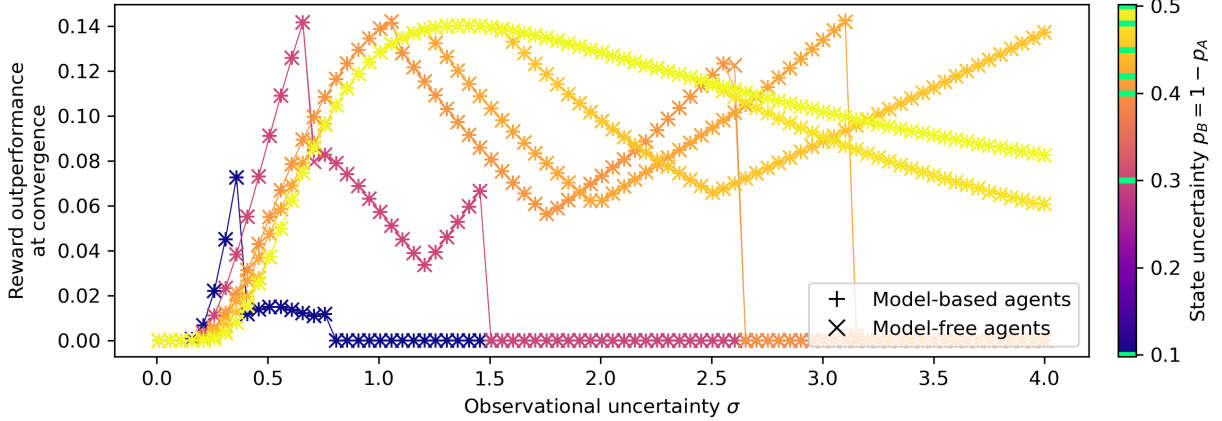


Figure 2: **Final Reward Outperformance** versus observational uncertainty σ for different state uncertainties p_B (shown in color, from dark blue at $p_B = 0.1$ to yellow at $p_B = 0.5$) and both learning dynamics: model-free (x) and model-based (+). Outperformance is defined as the relative reward gain of two-agent over single-agent learning, $(R_{\text{social}} - R_{\text{single}})/R_{\text{single}}$; positive values indicate that social learning yields higher rewards. Two-agent learning consistently outperforms independent learning across all tested levels of environmental uncertainty, with the largest benefit at intermediate observational uncertainty.

279 performance differences. Lastly, we study the learning timescales of learning together versus learning alone. We define
 280 a learning process as converged when the change in the strategy space, composed of the action probabilities across all
 281 agents, observations, and actions, falls below a threshold $\epsilon = 10^{-6}$. In the Supplementary Information we show that
 282 our main results are robust to the choice of convergence threshold by varying $\epsilon \in \{10^{-5}, 10^{-6}, 10^{-7}\}$. As one would
 283 expect, the smaller ϵ , the longer the convergence time becomes. The shape of the final reward outperformance curves is
 284 unaffected by this choice..

285 3.1 Final performance differences

286 With more environmental uncertainty, even optimally learned strategies yield lower rewards because the environment
 287 is less predictable. Thus, we are not measuring reward performance in absolute terms. Instead, we measure the
 288 *outperformance* of learning social learning strategies relative to independent learning strategies. We define outper-
 289 formance as the relative difference in average rewards between agents that can learn from social and environmental
 290 cues, R_{social} versus agents that can learn only from environmental cues, R_{single} , i.e., the reward outperformance amounts
 291 to $(R_{\text{social}} - R_{\text{single}})/R_{\text{single}}$. Positive outperformance indicates that learning social learning strategies yields higher
 292 average rewards than independent learning strategies, whereas negative outperformance indicates the opposite.

293 Figure 2 establishes that learning social learning strategies (both model-based and model-free) consistently outperform
 294 independent learning strategies across a wide range of environmental uncertainties. Model-based and model-free
 295 learning achieve identical final reward outperformance levels, indicating that their differences do not affect their final
 296 outcomes. They achieve an outperformance level of up to 14% across most levels of observational and state uncertainty.
 297 Only for very low state uncertainty $p_B \leq 0.1$, the maximal outperformance drops below 10%. When state uncertainty
 298 is very low, there is little benefit from social learning, as the environment is already quite predictable from the agent’s
 299 own observations. Importantly, they never underperform compared to independent learning strategies.

300 The yellow line in Figure 2 depicts the reward outperformance curve at convergence when state uncertainty is maximal
 301 ($p_B = 0.5$). At very low observational uncertainty, social learning provides minimal advantage because the environment
 302 is already quite predictable from individual observations alone. The outperformance rises sharply between $\sigma \approx 0.3$
 303 and peaks at intermediate observational uncertainty levels around $\sigma \approx 1.4$. Beyond this peak, outperformance
 304 declines gradually as observational uncertainty increases further and social cues become increasingly unreliable. Thus,
 305 an intermediate level of observational uncertainty provides the optimal balance for social learning to outperform
 306 independent learning strategies.

307 When state uncertainty decreases ($p_B < 0.5$), the final reward outperformance curve as a function of observational
 308 uncertainty σ reveals its characteristic shape (e.g., for $p_B = 0.4$). This shape persists across all levels of state uncertainty,
 309 but becomes increasingly compressed as state uncertainty p_B decreases. It is characterized by an initial sharp rise in
 310 outperformance at low observational uncertainty levels, reaching a maximum at intermediate levels, followed by a

311 gradual decline, then a secondary increase, before eventually dropping to zero at high observational uncertainty. This
 312 compression indicates that as the environment becomes more predictable from individual observations, the range of
 313 observational uncertainty levels for which learning together provides a significant advantage narrows. As the state of
 314 nature becomes more certain, short-term observations about it become less relevant.

315 Interestingly, the outperformance shape for maximal state uncertainty ($p_B = 0.5$, yellow line in Figure 2) serves as a
 316 reference point for understanding how reduced state uncertainty affects learned social learning benefits. For moderately
 317 reduced state uncertainty (with p_B between 0.48 and 0.42), the outperformance curves closely follow the reference line
 318 at low observational uncertainty levels. The curves begin to diverge only as observational uncertainty increases further.
 319 In contrast, when state uncertainty is substantially reduced (such as $p_B \leq 0.4$), the outperformance curves exhibit a
 320 more pronounced initial rise at low observational uncertainty than the reference line suggests. Yet after reaching their
 321 peak outperformance, these curves decline sharply and continue right below the reference line.

322 Overall, these results suggest that the final reward outperformance achieved through learned social learning strategies
 323 is robust across different learning models and levels of environmental uncertainty. The characteristic shape of the
 324 outperformance curve indicates a fundamental relationship between observational and state uncertainty in determining
 325 the benefits of learned social learning. To better understand the origins of this shape, we will next analyze the final
 326 strategies the agents have acquired at convergence.

327 3.2 Final strategies

328 What social learning strategies are learned? Figure 3 reveals that unexpectedly complex strategies emerge from reward
 329 learning, especially in environments where the state of nature is not fully uncertain ($p_B < 0.5$).

330 When the state of nature is maximally uncertain ($p_B = 0.5$), a single type of social learning strategy emerges consistently
 331 across intermediate to high levels of observational uncertainty ($\sigma \gtrsim 0.19$), as shown by the straight vertical gray line in
 332 Figure 3 A. In the absence of social information, i.e., when agents observe only the natural environmental signals (A_{\cdot} ,
 333 a_{\cdot} , b_{\cdot} , B_{\cdot}), they rely exclusively on these private observation signals to guide their decisions, as indicated by the light
 334 colors. However, when social information becomes available through observing the first-acting agent’s choice, a distinct
 335 behavioral pattern emerges. The second-acting agent learns to disregard its own natural observation signal (shown by
 336 dark colors) when that signal is weak (i.e., a or b) and contradicts the first-acting agent’s choice. At these specific
 337 observations ($a\bar{B}$ and $b\bar{A}$), the second-acting agent defers to the first-acting agent’s decision rather than following
 338 its own uncertain signal. We characterize this strategy as a ‘copy-when-uncertain’ strategy, in which agents rely on
 339 social information only when their private environmental cues are ambiguous. This strategy emerges robustly under
 340 both model-based and model-free learning dynamics, suggesting it represents a fundamental adaptive response to
 341 environmental uncertainty rather than an artifact of a particular learning mechanism.

342 At very low observational uncertainty ($\sigma \lesssim 0.19$), some observational signals occur with extremely low probability,
 343 which creates notable differences between model-based and model-free learning dynamics (Figure 3 A). When
 344 observational uncertainty is minimal, the strong signals A and B almost never occur. Consequently, agents select actions
 345 randomly at these observations, as indicated by the intermediate colors in Figure 3. However, this random behavior
 346 has no meaningful impact on final reward outperformance, precisely because such observations are so rare in practice.
 347 Model-free learning dynamics are more sensitive to observation probabilities than their model-based counterparts. As a
 348 result, the random behavior at strong signals A and B persists across a broader range of observational uncertainty levels
 349 for model-free agents compared to model-based agents. Interestingly, model-free agents also randomize their actions at
 350 observations $a\bar{B}$ and $b\bar{A}$, the very observations where, at higher uncertainty levels, they learn to completely disregard
 351 their natural signals and instead copy the first-acting agent’s choice. Model-based agents, by contrast, learn to rely
 352 fully on their private observation signals at these low uncertainty levels. Only within a narrow band of observational
 353 uncertainty ($0.09 \lesssim \sigma \lesssim 0.13$, dashed gray lines in Figure 3) do model-free agents display clear signatures of the
 354 ‘copy-when-uncertain’ strategy, a pattern that holds consistently across all levels of state uncertainty p_B . Despite
 355 these differences between learning dynamics, the overall behavioral patterns at very low observational uncertainty
 356 remain qualitatively similar across all levels of state uncertainty, suggesting that extreme information scarcity produces
 357 comparable learning outcomes regardless of environmental predictability.

358 When the environment is not fully uncertain ($p_B < 0.5$), a rich diversity of social learning strategies emerges, as
 359 illustrated in Figure 3 B. These strategies can be categorized into three primary regimes based on the behavior of the
 360 first-acting agent. i) In the first regime, the first-acting agent follows its natural observation signals. ii) The second
 361 regime occurs when the first-acting agent begins to ignore weak signals associated with the less-likely environmental
 362 state. iii) In the third regime, this behavior extends further, as the agent ignores both weak and strong signals in favor of
 363 the more likely state. We now examine the specific social learning strategies that arise within these three regimes across
 364 the range of increasing observational uncertainty ($\sigma \gtrsim 0.13$).

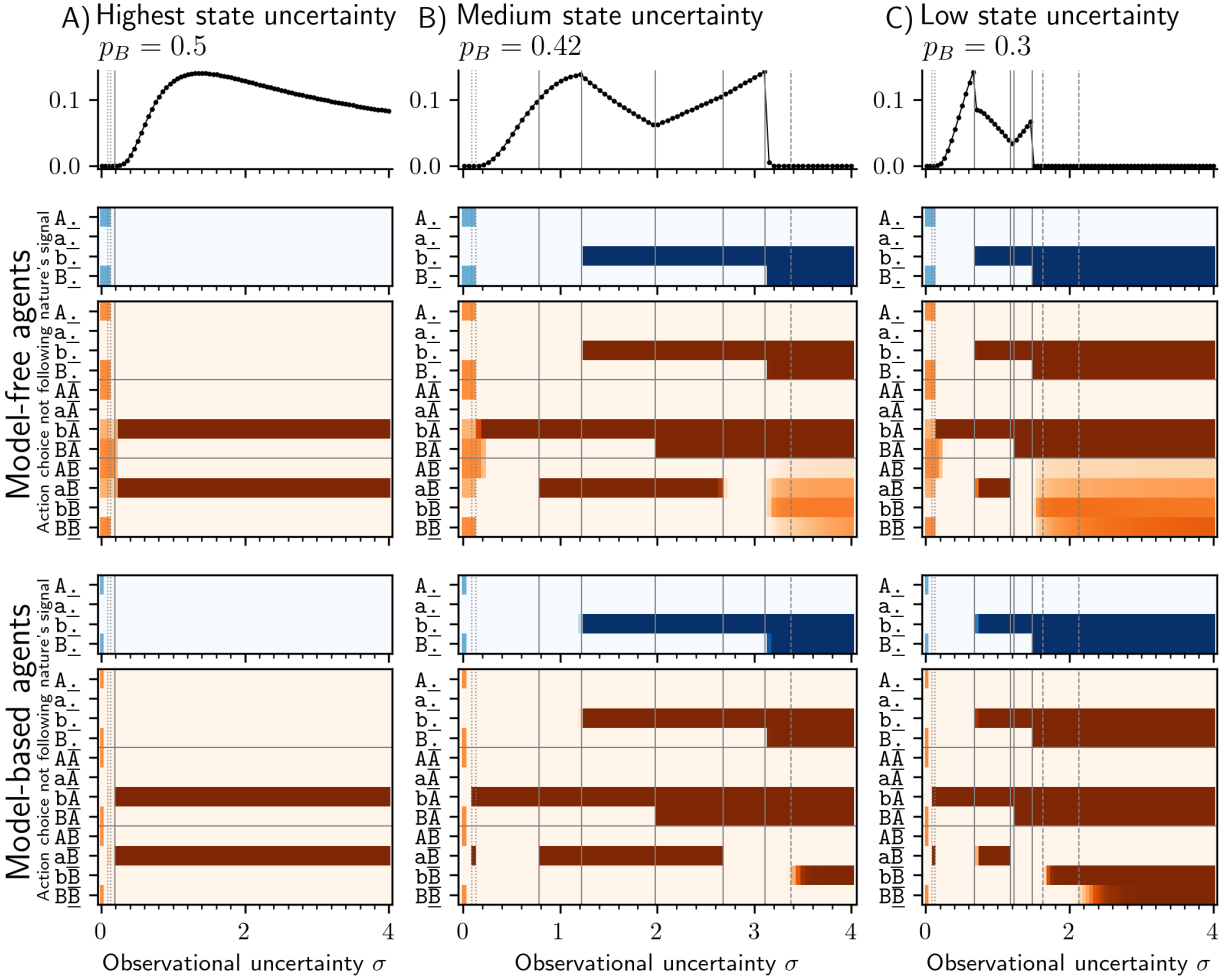


Figure 3: **Learned strategies** at convergence across different levels of observational uncertainty (σ) and state uncertainty $p_B = 0.5$ (A), $p_B = 0.42$ (B), $p_B = 0.3$ (C). The top panels show the reward outperformance curves for reference. The heatmaps display the final learned behaviors across different levels of observational uncertainty (σ) and state uncertainty (p_B). Blue colormaps represent the independent agent in the single-agent environment, while orange colormaps depict behavior in the two-agent environment. The color intensity encodes the final behavioral strategy: light colors indicate that agents follow their natural observation signal, while dark colors indicate that they deviate from it in favor of the alternative action. For example, for observations A_{\cdot} , $a_{\underline{B}}$ (where the first letter indicates the observation signal), light colors denote choosing action A, whereas dark colors denote choosing action B. Results for the model-based (model-free) agents are averaged over 25 (10) random initial strategies. At high observational uncertainty, agents robustly learn a 'copy-when-uncertain' strategy; richer strategy diversity emerges as state uncertainty decreases.

365 In the first regime, where observational uncertainty remains relatively low, the second-acting agent develops a strategy
 366 we term ‘copy-when-uncertain-about-the-less-likely-state.’ In this context, the first-acting agent consistently follows
 367 its natural observation signals. The second-acting agent, however, learns to selectively ignore its own weak signals
 368 when they point toward the less-likely environmental state (indicated by the dark colors at observation $b_{\underline{A}}$). In these
 369 instances, the agent defers to the first-acting agent’s choice. Conversely, when its private signal suggests the more-likely
 370 environmental state (as seen at observation $a_{\underline{B}}$), the agent continues to rely on its own information to guide its decision.
 371 Only as observational uncertainty increases further, specifically reaching a threshold of approximately $\sigma \approx 0.78$ in
 372 Figure 3 B, the second-acting agent transitions fully to the ‘copy-when-uncertain’ strategy. Thus, decreased state
 373 uncertainty shifts the emergence of this strategy to higher levels of observational noise. Notably, while this transition
 374 marks a clear change in behavioral patterns, it does not significantly impact the shape of the final reward outperformance
 375 curve, compared to the fully uncertain environment ($p_B = 0.5$) discussed earlier.

376 In the second regime, beginning around $\sigma \approx 1.22$ in Figure 3 B, changes in the first-acting agent’s strategy cause
 377 a decline in reward outperformance. While the second-acting agent initially persists with a ‘copy-when-uncertain’
 378 approach, it eventually transitions at higher levels of observational uncertainty ($\sigma \approx 1.98$) to a ‘follow-the-leader-in-the-
 379 more-likely-state-but-copy-when-uncertain-in-the-less-likely-state’ strategy. Under this rule, the second agent adopts
 380 the first agent’s choice when it aligns with the more probable state of the world, effectively ignoring its own private
 381 signals, yet only defers to the leader in the less likely state when its own observations are ambiguous. This transition
 382 marks a reversal, with reward outperformance beginning to rise again. Eventually, at $\sigma \approx 2.68$ in Figure 3 B, the agent
 383 adopts a ‘follow-the-leader-in-the-more-likely-state’ strategy, where it relies on its own natural observation signals for
 384 the less likely state regardless of the leader’s action. Interestingly, this final stage is less social than the previous strategy
 385 despite increased observational uncertainty. This might be due to the combination of high observational uncertainty and
 386 low state uncertainty, which makes the first agent’s choice of the less likely state rare and unreliable as a social cue.
 387 Note that in the first and second regimes, the learned strategies are consistent across both model-based and model-free
 388 learning dynamics.

389 In the third regime, which begins around $\sigma = 3.11$ in Figure 3 B, the first-acting agent fully ignores its natural
 390 observation signal for the less-likely environmental state (indicated by dark colors at $B_{\underline{.}}$ and $b_{\underline{.}}$), at which point the
 391 final reward outperformance drops back to zero. During this phase, the second-acting agent continues to follow the
 392 leader when they choose the naturally more-likely state (shown by dark colors at $b_{\underline{A}}$ and $B_{\underline{A}}$). When the second-acting
 393 agent observes the first-acting agent choosing the less-likely environmental state (B), the model-free learning dynamics
 394 randomize their strategy again, since such observations practically never occur. The model-based learning dynamics,
 395 by contrast, learn to ignore both their natural observation signal and the first-acting agent’s choice (first evident by
 396 dark colors at $b_{\underline{B}}$ from around $\sigma = 3.38$ in Figure 3 B) in order to choose the more-likely environmental state action
 397 A . These are the only transitions of the model-based agent that happen gradually rather than sharply. We interpret
 398 this behavior as an artifact of our model-based learning dynamics, which can learn from experience they practically
 399 never encounter. For this reason, we do not further classify these strategies, though interestingly, a signature of these
 400 strategies is also observed in the randomized strategy of the model-free learning dynamics at these same observations.

401 For lower state uncertainty, these transitions become compressed together, though their topology, i.e., the identity and
 402 ordering of transitions, may change (Figure 3 C). For instance, in the first regime, the second-acting agent does not
 403 re-adopt the ‘copy-when-uncertain’ strategy. Instead, when the second regime begins, and the first-acting agent starts
 404 disregarding its less likely individual observation of the less likely state ($b_{\underline{.}}$), the second-acting agent responds by
 405 adopting the ‘copy-when-uncertain’ strategy. Subsequently, at $\sigma \approx 1.18$ in Figure 3 C, the second-acting agent switches
 406 from this ‘copy-when-uncertain’ strategy to a new ‘copy-when-uncertain-about-the-less-likely-state’ strategy, before
 407 finally transitioning to the ‘follow-the-leader-in-the-more-likely-state’ strategy at $\sigma \approx 1.23$.

408 Overall, these results demonstrate that a rich diversity of social learning strategies can emerge from simple reward-based
 409 learning dynamics, particularly when environmental state uncertainty is reduced. The specific strategies that agents
 410 adopt depend critically on the interplay between observational and state uncertainty. Notably, reward outperformance
 411 increases when the second-acting agent learns to disregard its private signal about the less-likely environmental state and
 412 instead follows the first-acting agent’s choice. Conversely, reward outperformance decreases when the first-acting agent
 413 itself begins to ignore its private signals for the less-likely state, thereby reducing the reliability of social information.
 414 The information-theoretic basis of this pattern is straightforward: the second-acting agent’s observation combines its
 415 own noisy private signal with the first-acting agent’s choice, which itself integrates that agent’s private signal about the
 416 same environment. Combining two imperfect signals reduces residual state uncertainty more than either alone, and
 417 this benefit is greatest precisely when individual signals are weakest, which is why the copy-when-uncertain strategy
 418 emerges under high observational uncertainty. Having characterized these emergent behavioral patterns, we now turn to
 419 analyzing the timescales over which agents acquire these social learning strategies.

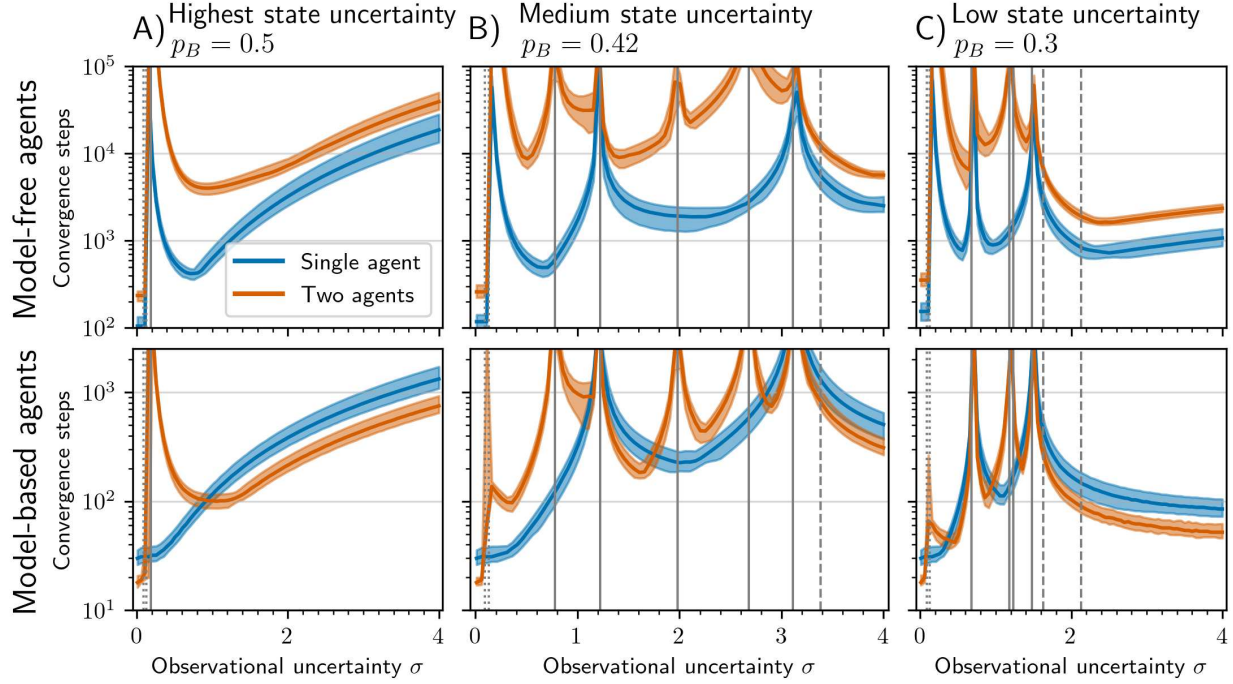


Figure 4: **Learning times** until convergence of the behavioral rules versus observational uncertainty σ for three different state uncertainties $p_B = 0.5$ (A), $p_B = 0.42$ (B), $p_B = 0.3$ (C). The upper panels show model-free learning dynamics, the lower panels show the model-based ones. Blue (vermillion) lines show the single-agent (two-agent) environment. Shaded regions indicate the space between 5- and 95-percentiles from 25 (10) random initial strategies for the model-based (model-free) agents. Convergence times spike at strategy transition points, reflecting a critical slowing down. Outside these points, model-based agents learn approximately twice as fast in the two-agent environment, while model-free agents learn faster alone.

420 3.3 Learning timescales

421 How long do learning agents require to converge to their final strategies? As we expected, model-based learning
 422 dynamics converge faster than model-free learning. Figure 4 reveals that the difference in learning speed is about 1-2
 423 orders of magnitude, as shown by the extent of the logarithmic y-axis in the plots. For most levels of environmental
 424 uncertainty, model-free learning in both single-agent and two-agent environments requires between 10^2 and 10^5 time
 425 steps, whereas model-based learning occurs between 10^1 and 10^3 - 10^4 time steps.

426 Interestingly, we observe critical transitions at the strategy transition points identified in the previous section. Here,
 427 learning slows dramatically: the convergence timescales peak, signaling a critical slowing down. This effect appears
 428 for both model-based and model-free dynamics and reflects agents being caught between two competing strategies.
 429 When the expected benefits of alternative actions are similar, updates become small, and convergence is prolonged.
 430 For example, a second-acting agent deciding between “copy-when-uncertain” and relying on its private signal will
 431 transition only slowly from choosing A at observation $a\bar{B}$ and B at $b\bar{A}$ to the opposite choices. Model-free learners
 432 show an additional peak at very low observational uncertainty σ because strong signals (A and B) occur so rarely that
 433 experience at those observations is sparse, which further increases indecision and slows convergence.

434 Generally, model-free agents learn faster alone, whereas model-based agents learn faster together. Only at the critical
 435 points do the model-based agents learn more slowly together than alone, as they have to undergo a critical transition
 436 to a different social learning strategy. Outside these critical points, two model-based agents learn about twice as fast
 437 as a single model-based agent, whereas a single model-free agent learns about twice as fast as two model-free agents
 438 – as, e.g., the parallel lines in Figure 4 A) for $\sigma > 1.5$ indicate. This model-free speed disadvantage arises from the
 439 stationary observation frequency weighting (P_X^{i,σ^2}) in the model-free update rule: in the two-agent environment, each
 440 acting observation is visited less frequently in stationarity (because each agent is inactive during approximately half of
 441 all steps) so updates are smaller per step. Model-based agents are insulated from this effect because their updates use
 442 analytically computed expected rewards regardless of how often each observation is visited. The model-based speed
 443 advantage has a distinct origin. Under the assumption that social observation is temporally free, each social round

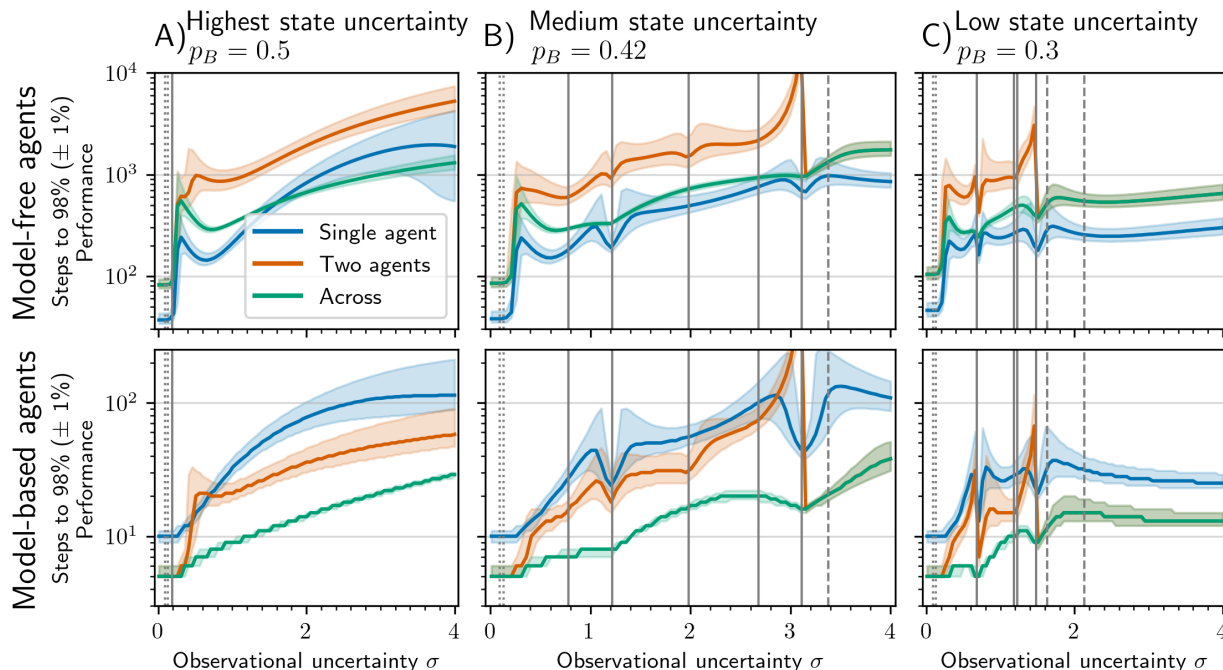


Figure 5: **Learning time steps** required until 98% ($\pm 1\%$) of the final reward at strategy convergence is reached versus observational uncertainty σ for three different state uncertainties $p_B = 0.5$ (A), $p_B = 0.42$ (B), $p_B = 0.3$ (C). The upper panels show model-free learning dynamics, the lower panels show the model-based ones. All trajectories started from uniformly random strategies, as we saw in Figure 4 that multiple initial conditions do not have a significant impact on the results. Blue (vermillion) colors indicate the single-agent (two-agent) environment. The green color indicates the time steps required for the two learners to reach the respective performance level of the single agent. Most of the time spent on strategy convergence is devoted to achieving outperformance over independent learning; reaching near-optimal reward takes an order of magnitude fewer steps than full strategy convergence.

444 (in which both agents act sequentially) occupies the same wall-clock time as a single-agent step yet generates two
 445 reward-earning interactions. An agent with a world model can compute expected rewards over the full joint dynamics,
 446 directly exploiting the other agent’s choice, and thus, achieving an approximately two-fold increase in effective learning
 447 rate. Without a world model, an agent updates only in proportion to its own individual observation encounters, which
 448 are not doubled by the social structure. The collective richness of the social round remains cognitively out of reach.
 449 Thus, when having a reliable internal model of the environment, learning together not only leads to better results, but it
 450 is also faster.

451 However, what if such an internal model is not available? How long does it take until a decent amount of the final
 452 reward is achieved? Figure 5 shows the time steps required to achieve 98% of the final reward (dark lines). The shaded
 453 regions show the times required to achieve between 99% (at the top) and 97% (at the bottom) of the final reward. Blue
 454 colors indicate single-agent learning, vermillion colors indicate two-agent learning, and green colors indicate the time
 455 steps required for two-agent learning to outperform 98% ($\pm 1\%$) of the final reward achieved by single-agent learning.

456 Overall, it takes all learning dynamics about an order of magnitude less time to achieve most of the final reward
 457 outperformance than to fully converge in strategy (compare the y-scales of Figure 5 and Figure 4). Even some of
 458 the critical transitions in learning time scales identified in the strategy convergence analysis are smoothed out when
 459 considering the time steps required to achieve most of the final reward outperformance. Most interestingly, however, is
 460 the time required for two-agent learning to achieve most of the performance of single-agent learning (green lines in
 461 Figure 5). The model-based agents require only between 3 and 30 steps across all levels of environmental uncertainty.
 462 The model-free agents require longer, but in the order of magnitude comparable to how much the single model-free
 463 agents require to achieve most of their final reward. Thus, most of the time spent converging in strategy is devoted to
 464 achieving outperformance over independent learning.

465 Taken together, we find that learning together not only yields better results under environmental uncertainty but can also
 466 significantly speed up the learning process – provided that agents have a reliable internal model of the environment.

467 When such an internal model is not available, learning together still yields better results, but the learning speed advantage
 468 is lost. However, even model-free learning does not take significantly longer than independent learning to achieve most
 469 of the final reward that an independently learning agent achieves.

470 4 Discussion

471 In this work, we have systematically investigated how environmental uncertainty shapes the timescales of learning social
 472 learning strategies through pure reward-based learning. We find that learning together through reward-based learning
 473 not only yields results up to 14% better than independent learning across wide ranges of environmental uncertainty,
 474 regardless of whether agents possess internal world models. If agents possess reliable internal world models, learning
 475 together typically occurs also faster than learning alone, despite the increased complexity of the strategy the agents
 476 must learn. Without such internal world models, learning together is slower than learning alone. Yet, it takes about
 477 the same amount of time to achieve most of an independent agent’s final performance. Thus, learning together poses
 478 no significant disadvantage in terms of learning speed. These findings highlight an often-overlooked dimension of
 479 the remarkable success of social learning strategies: they not only can yield better outcomes but also can accelerate
 480 the learning process itself. This finding helps explain why social learning is so widespread across diverse species and
 481 ecological contexts [Boyd et al., 2011].

482 One of the findings we were most surprised about is the rich diversity of social learning strategies that emerged in our
 483 comparably idealized model from pure reward-based learning. We designed our model with the intention of being the
 484 ‘simplest’ model that could show the emergence of a ‘copy-when-uncertain’ strategy. Indeed, this is what we find in the
 485 typical case of full environmental state uncertainty. However, when the environment is not fully uncertain, we find that
 486 more nuanced social learning strategies emerge, such as a ‘copy-when-uncertain-about-the-less-likely-state’ strategy
 487 and a ‘follow-the-leader-in-the-more-likely-state’ strategy. These strategies are not pre-programmed heuristics but
 488 emerge from our domain-general reward-learning mechanism, highlighting the flexibility and adaptability of learned
 489 social-learning strategies [Kendal et al., 2018]. Reinforcement learning was originally conceived as a single-agent
 490 learning process [Sutton and Barto, 2018]. Thus, there is an ongoing search for the best way to incorporate ‘the
 491 social’ into it, as evidenced by previous works devising various learning processes that incorporate sociality either
 492 via stimulus enhancement [Galef, 1988], value shaping [Najar et al., 2020], decision biasing [Toyokawa et al., 2019],
 493 social generalization [Witt et al., 2024], or social features [Schultner et al., 2025a]. Our learning dynamics most
 494 closely resemble the model by Schultner et al. [2025a]. Thus, our results suggest that no explicit incorporation of
 495 ‘the social’ into a reinforcement learning update is required beyond observing social features or signals. As we have
 496 shown, complex social learning strategies can emerge out of pure reward associations, resonating with recent calls in
 497 evolutionary game theory to move beyond hand-picked strategy sets toward learning-based strategy generation [Garcia
 498 and Traulsen, 2025].

499 The transitions between these learned social learning strategies also robustly reveal another interesting phenomenon:
 500 critical slowing down at the transition points. This is consistent with the general theory of critical transitions [Scheffer
 501 et al., 2009]: near a bifurcation point, the dominant eigenvalue of the system’s Jacobian approaches zero, and
 502 convergence times diverge. We hypothesize that each strategy transition corresponds to such a bifurcation of the
 503 learning dynamics, and we leave a formal eigenvalue analysis for future work. At these points, learning timescales
 504 peak, indicating that agents take multiple orders of magnitude longer to converge to their final strategies. However,
 505 this time is spent achieving only the last few percent of the final rewards. Thus, from a practical perspective, these
 506 critical slowing down points may not pose significant hindrances to effective decisions, as agents can achieve most of
 507 their performance benefits relatively quickly. In contrast, being in a critical transition offers the advantage of flexible
 508 adaptation to changing environmental conditions, as the strategies are not yet fully stabilized [Beggs, 2007, Klamser
 509 and Romanczuk, 2021]. Importantly, this phenomenon is not an artifact of our deterministic dynamics. The bifurcation
 510 structure is a property of the underlying reward landscape, and stochastic RL algorithms would exhibit the same
 511 qualitative behavior near strategy transitions, manifesting as longer expected convergence times and increased variance
 512 in learning trajectories.

513 These findings provide insights for designing and understanding natural, artificial, and hybrid multi-agent systems,
 514 especially regarding the cognitive requirements for social learning. The fact that model-free agents – without an internal
 515 environmental model and only simple associative learning – can acquire sophisticated, multi-layered social learning
 516 strategies suggests that complex social behavior need not require complex cognition. Even organisms with relatively
 517 simple nervous systems might develop effective social learning strategies through basic reward associations [Leadbeater
 518 and Chittka, 2007, Chittka and Rossi, 2022]. This could explain why adaptive social learning appears in taxonomically
 519 diverse groups, from insects to mammals. However, our finding that model-based learners acquire social strategies
 520 much faster than model-free learners suggests a cognitive advantage for organisms that can build internal models of
 521 their environments. Mechanistically, this advantage is specific to the social setting. A world model enables an agent to
 522 leverage the collective dynamics of social interaction, pooling the doubled reward signal that each social round provides,

rather than being confined to learning from its own directly experienced observations at individual encounter rates. In rapidly changing environments or during critical developmental periods, the speed of learning could be decisive for survival. This provides a potential adaptive explanation for why more cognitively sophisticated organisms evolved the capacity for model-based learning: not just to achieve better final performance, but to adapt quickly enough to survive in dynamic environments [Suzuki and Arita, 2004, Fernando et al., 2018]. Regarding AI systems, this also means that their architecture has profound implications for collective learning speed. AI systems with explicit world models [Ha and Schmidhuber, 2018] should accelerate collective learning when they interact with humans or other AI agents. By contrast, purely reactive systems without internal models – analogous to our model-free learners – may slow down collective adaptation even as they eventually reach good solutions. This finding is supported by empirical research in multi-agent reinforcement learning. Purely model-free agents have been shown to fail to adopt social learning strategies without an auxiliary model-based learning component, whereas agents with world models develop generalised social learning policies and transfer them to novel environments [Ndousse et al., 2021]. This suggests that investments in model-based AI architectures may pay off not just in individual AI performance but also in AI agents’ capacity for cultural accumulation [Cook et al., 2024] and emergent social behaviour [Köster et al., 2025]. More broadly, our findings contribute a micro-founded learning mechanism to emerging accounts of machine culture and machine behaviour in hybrid human-AI populations [Rahwan et al., 2019, Brinkmann et al., 2023], showing how agents embedded in shared environments develop and transmit social learning strategies through pure reward associations, without any explicit programming of cultural norms.

Methodologically, our work demonstrates that evolutionary and learning-dynamics approaches based on replicator-like dynamics [Hofbauer and Sigmund, 1998, Fudenberg and Levine, 2016, Barfuss, 2020a] can investigate the emergence of sociality from pure reward-based learning. Specifically, our environment model can serve as a milestone from where further questions can be investigated, such as different preference structures, asymmetric information, noisy or delayed social observations, larger group sizes including varied network topologies, and non-stationary changes in more complex environments. To our knowledge, this work is the first to apply the concepts of model-free and model-based learning to replicator-like learning dynamics [Barfuss, 2020b]. Classic reinforcement learning dynamics were using the model-based variant [Tuyls et al., 2003, Sato and Crutchfield, 2003, Barfuss et al., 2019][Leslie and Collins, 2005], without being aware of that interpretation [Kaisers and Tuyls, 2011]. More recent works incorporate choice frequencies into the reinforcement signal, as our model-free version [Hu et al., 2022, Goll et al., 2024], but have not made the connection to model-free learning explicit. Thus, our work contributes to the bridges between the fields of multi-agent reinforcement learning and evolutionary game theory, opening new avenues for cross-disciplinary research [Madhushani Sehwal et al., 2025, Barfuss et al., 2025]. Recent work in evolutionary game theory has investigated how rational decision-making strategies themselves evolve rather than being assumed a priori [Salahshour, 2025, Salahshour and Couzin, 2025], offering a complementary perspective to our reward-learning account of emergent strategies. Future work should endogenize the development of internal world models through experience rather than assuming perfect internal models from the start. While these deterministic dynamics provide analytical tractability, they abstract away potentially important stochastic effects [Galla, 2009, Barfuss and Meylahn, 2023, Rudd-Jones et al., 2025]. Further research is required to explore how stochasticity in learning processes influences the emergence of social learning strategies. We expect the qualitative speed ordering – model-based agents learning faster together, model-free agents learning faster alone – to be robust to stochastic sampling, since the underlying mechanism is informational rather than algorithmic: any model-based algorithm exploiting the joint dynamics will benefit from the doubled reward signal per social round, while any model-free algorithm updating in proportion to empirical visit rates will be penalized by the sparser individual observation frequencies in the two-agent environment. We also consider it conceivable that some of our findings, such as the existence of a critical slowing down in learning, could serve as hypotheses to be tested experimentally with human subjects in controlled laboratory settings.

In conclusion, our work demonstrates that sophisticated social learning strategies need not be inherited, rational, or preprogrammed — they can arise quickly whenever agents seek and learn from reward in uncertain environments.

Acknowledgments. WB acknowledges financial support from the Cooperative AI Foundation. RPM was supported by a UKRI Future Leaders Fellowship (MR/X036863/1) and Templeton World Charity Foundation Inc. (TWCF202120647).

Code availability. All computer code was written in Python and is contained in reproducible Supplementary Information. We will make it openly available upon acceptance of this manuscript.

References

- Alberto Acerbi, Alex Mesoudi, and Marco Smolla. *Individual-Based Models of Cultural Evolution*. 2022.
- Kenichi Aoki and Marcus W. Feldman. Evolution of learning strategies in temporally and spatially variable environments: A review of theory. *Theoretical Population Biology*, 91:3–19, February 2014. ISSN 0040-5809. doi: 10.1016/j.tpb.2013.10.004.

- 578 Sara Arganda, Alfonso Pérez-Escudero, and Gonzalo G. de Polavieja. A common rule for decision making in animal
579 collectives across species. *Proceedings of the National Academy of Sciences*, 109(50):20508–20513, December 2012.
580 doi: 10.1073/pnas.1210664109.
- 581 Abhijit V. Banerjee. A simple model of herd behavior. *The Quarterly Journal of Economics*, 107(3):797–817, 1992.
582 doi: 10.2307/2118364.
- 583 Wolfram Barfuss. Towards a unified treatment of the dynamics of collective learning. In *Challenges and Opportunities*
584 *for Multi-Agent Reinforcement Learning, AAAI Spring Symposium*, 2020a.
- 585 Wolfram Barfuss. Reinforcement Learning Dynamics in the Infinite Memory Limit. In *Proceedings of the 19th*
586 *International Conference on Autonomous Agents and MultiAgent Systems, AAMAS '20*, pages 1768–1770, Richland,
587 SC, May 2020b. International Foundation for Autonomous Agents and Multiagent Systems.
- 588 Wolfram Barfuss. Dynamical systems as a level of cognitive analysis of multi-agent learning. *Neural Computing and*
589 *Applications*, 34(3):1653–1671, February 2022. ISSN 1433-3058. doi: 10.1007/s00521-021-06117-0.
- 590 Wolfram Barfuss and Richard P. Mann. Modeling the effects of environmental and perceptual uncertainty using
591 deterministic reinforcement learning dynamics with partial observability. *Physical Review E*, 105(3):034409, March
592 2022. doi: 10.1103/PhysRevE.105.034409.
- 593 Wolfram Barfuss and Janusz M. Meylahn. Intrinsic fluctuations of reinforcement learning promote cooperation.
594 *Scientific Reports*, 13(1):1309, January 2023. ISSN 2045-2322. doi: 10.1038/s41598-023-27672-7.
- 595 Wolfram Barfuss, Jonathan F. Donges, and Jürgen Kurths. Deterministic limit of temporal difference reinforcement
596 learning for stochastic games. *Physical Review E*, 99(4):043305, April 2019. doi: 10.1103/PhysRevE.99.043305.
- 597 Wolfram Barfuss, Jessica Flack, Chaitanya S. Gokhale, Lewis Hammond, Christian Hilbe, Edward Hughes, Joel Z.
598 Leibo, Tom Lenaerts, Naomi Leonard, Simon Levin, Udari Madhushani Sehwan, Alex McAvoy, Janusz M. Meylahn,
599 and Fernando P. Santos. Collective cooperative intelligence. *Proceedings of the National Academy of Sciences*, 122
600 (25):e2319948121, June 2025. doi: 10.1073/pnas.2319948121.
- 601 John M Beggs. The criticality hypothesis: How local cortical networks might optimize information processing.
602 *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 366(1864):
603 329–343, August 2007. ISSN 1364-503X. doi: 10.1098/rsta.2007.2092.
- 604 Timothy E. J. Behrens, Laurence T. Hunt, Mark W. Woolrich, and Matthew F. S. Rushworth. Associative learning of
605 social value. *Nature*, 456(7219):245–249, November 2008. ISSN 1476-4687. doi: 10.1038/nature07538.
- 606 Clémence Bergerot, Wolfram Barfuss, and Pawel Romanczuk. Moderate confirmation bias enhances decision-making
607 in groups of reinforcement-learning agents. *PLOS Computational Biology*, 20(9):e1012404, September 2024. ISSN
608 1553-7358. doi: 10.1371/journal.pcbi.1012404.
- 609 Martino Bernasconi, Federico Cacciamani, Simone Fioravanti, Nicola Gatti, and Francesco Trovò. The evolutionary
610 dynamics of soft-max policy gradient in multi-agent settings. *Theoretical Computer Science*, 1027:115011, 2025.
611 doi: 10.1016/j.tcs.2024.115011.
- 612 Sushil Bikhchandani, David Hirshleifer, and Ivo Welch. A Theory of Fads, Fashion, Custom, and Cultural Change
613 as Informational Cascades. *Journal of Political Economy*, 100(5):992–1026, October 1992. ISSN 0022-3808,
614 1537-534X. doi: 10.1086/261849.
- 615 Daan Bloembergen, Karl Tuyls, Daniel Hennes, and Michael Kaisers. Evolutionary Dynamics of Multi-Agent
616 Learning: A Survey. *Journal of Artificial Intelligence Research*, 53:659–697, August 2015. ISSN 1076-9757. doi:
617 10.1613/jair.4818.
- 618 Tilman Börgers and Rajiv Sarin. Learning Through Reinforcement and Replicator Dynamics. *Journal of Economic*
619 *Theory*, 77(1):1–14, November 1997. ISSN 00220531. doi: 10.1006/jeth.1997.2319.
- 620 Diana Borsa, Nicolas Heess, Bilal Piot, Siqi Liu, Leonard Hasenclever, Remi Munos, and Olivier Pietquin. Observational
621 Learning by Reinforcement Learning. In *Proceedings of the 18th International Conference on Autonomous Agents*
622 *and MultiAgent Systems, AAMAS '19*, pages 1117–1124, Richland, SC, May 2019. International Foundation for
623 Autonomous Agents and Multiagent Systems. ISBN 978-1-4503-6309-9.
- 624 Robert Boyd and Peter J. Richerson. *Culture and the Evolutionary Process*. University of Chicago Press, June 1988.
625 ISBN 978-0-226-06933-3.

- 626 Robert Boyd, Peter J. Richerson, and Joseph Henrich. The cultural niche: Why social learning is essential for human
627 adaptation. *Proceedings of the National Academy of Sciences*, 108(supplement_2):10918–10925, June 2011. doi:
628 10.1073/pnas.1100290108.
- 629 Levin Brinkmann, Fabian Baumann, Jean-François Bonnefon, Maxime Derex, Thomas F. Müller, Anne-Marie Nuss-
630 berger, Agnieszka Czaplicka, Alberto Acerbi, Thomas L. Griffiths, Joseph Henrich, Joel Z. Leibo, Richard McElreath,
631 Pierre-Yves Oudeyer, Jonathan Stray, and Iyad Rahwan. Machine culture. *Nature Human Behaviour*, 7(11):
632 1855–1868, November 2023. ISSN 2397-3374. doi: 10.1038/s41562-023-01742-2.
- 633 Lars Chittka and Natacha Rossi. Social cognition in insects. *Trends in Cognitive Sciences*, 26(7):578–592, July 2022.
634 ISSN 1364-6613, 1879-307X. doi: 10.1016/j.tics.2022.04.001.
- 635 Jonathan Cook, Chris Lu, Edward Hughes, Joel Z. Leibo, and Jakob Foerster. Artificial Generational Intelligence:
636 Cultural Accumulation in Reinforcement Learning. In *Advances in Neural Information Processing Systems 37*,
637 volume 37, pages 59689–59715, December 2024. doi: 10.52202/079017-1907.
- 638 Lukas Danwitz and Bettina von Helversen. Observational learning of exploration-exploitation strategies in bandit tasks.
639 *Cognition*, 259:106124, 2025. doi: 10.1016/j.cognition.2025.106124.
- 640 Dominik Deffner, Vivien Kleinow, and Richard McElreath. Dynamic social learning in temporally and spatially
641 variable environments. *Royal Society Open Science*, 7(12):200734, December 2020. ISSN 2054-5703. doi:
642 10.1098/rsos.200734.
- 643 Dominik Deffner, Natalia Fedorova, Jeffrey Andrews, and Richard McElreath. Bridging theory and data: A computa-
644 tional workflow for cultural evolution. *Proceedings of the National Academy of Sciences*, 121(48):e2322887121,
645 November 2024. doi: 10.1073/pnas.2322887121.
- 646 Charles Efferson, Rafael Lalive, Peter J. Richerson, Richard McElreath, and Mark Lubell. Conformists and mavericks:
647 The empirics of frequency-dependent cultural transmission. *Evolution and Human Behavior*, 29(1):56–64, January
648 2008. ISSN 1090-5138. doi: 10.1016/j.evolhumbehav.2007.08.003.
- 649 Chrisantha Fernando, Jakub Sygnowski, Simon Osindero, Jane Wang, Tom Schaul, Denis Teplyashin, Pablo Sprech-
650 mann, Alexander Pritzel, and Andrei Rusu. Meta-learning by the Baldwin effect. In *Proceedings of the Genetic and
651 Evolutionary Computation Conference Companion, GECCO '18*, pages 1313–1320, New York, NY, USA, July 2018.
652 Association for Computing Machinery. doi: 10.1145/3205651.3208249.
- 653 Drew Fudenberg and David K. Levine. Whither Game Theory? Towards a Theory of Learning in Games. *Journal of
654 Economic Perspectives*, 30(4):151–170, November 2016. ISSN 0895-3309. doi: 10.1257/jep.30.4.151.
- 655 Bennett G. Galef. Imitation in Animals: History, Definition, and Interpretation of Data from the Psychological
656 Laboratory. In *Social Learning*. Psychology Press, 1988.
- 657 Tobias Galla. Intrinsic Noise in Game Dynamical Learning. *Physical Review Letters*, 103(19):198702, November 2009.
658 ISSN 0031-9007, 1079-7114. doi: 10.1103/PhysRevLett.103.198702.
- 659 Julian Garcia and Arne Traulsen. Picking Strategies in Games of Cooperation. *Proceedings of the National Academy of
660 Sciences*, 122(25):e2319925121, June 2025. doi: 10.1073/pnas.2319925121.
- 661 David Goll, Jobst Heitzig, and Wolfram Barfuss. Deterministic Model of Incremental Multi-Agent Boltzmann
662 Q-Learning: Transient Cooperation, Metastability, and Oscillations, December 2024.
- 663 David Ha and Jürgen Schmidhuber. World Models. *World Models*, 1(1):e10, March 2018. doi: 10.5281/zenodo.12076
664 31.
- 665 Seungwoong Ha and Hawoong Jeong. Social learning spontaneously emerges by searching optimal heuristics with
666 deep reinforcement learning. In *Proceedings of the 40th International Conference on Machine Learning*, pages
667 12319–12338. PMLR, July 2023.
- 668 Cecilia Heyes. What’s social about social learning? *Journal of Comparative Psychology*, 126(2):193–202, 2012. ISSN
669 1939-2087. doi: 10.1037/a0025180.
- 670 Cecilia Heyes. Blackboxing: Social learning strategies and cultural evolution. *Philosophical Transactions of the Royal
671 Society B: Biological Sciences*, 371(1693):20150369, May 2016. ISSN 0962-8436. doi: 10.1098/rstb.2015.0369.
- 672 Cecilia Heyes and John M. Pearce. Not-so-social learning strategies. *Proceedings of the Royal Society B: Biological
673 Sciences*, 282(1802):20141709, March 2015. ISSN 0962-8452. doi: 10.1098/rspb.2014.1709.

- 674 Josef Hofbauer and Karl Sigmund. *Evolutionary Games and Population Dynamics*. Cambridge University Press, 1
675 edition, May 1998. ISBN 978-0-521-62365-0 978-0-521-62570-8 978-1-139-17317-9. doi: 10.1017/CBO9781139
676 173179.
- 677 Shuyue Hu, Chin-Wing Leung, Ho-fung Leung, and Harold Soh. The Dynamics of Q-learning in Population Games: A
678 Physics-inspired Continuity Equation Model. In *Proceedings of the 21st International Conference on Autonomous
679 Agents and Multiagent Systems, AAMAS '22*, pages 615–623, Richland, SC, May 2022. International Foundation for
680 Autonomous Agents and Multiagent Systems. ISBN 978-1-4503-9213-6.
- 681 Michael Kaisers and Karl Tuyls. FAQ-learning in matrix games: Demonstrating convergence near Nash equilibria, and
682 bifurcation of attractors in the battle of sexes. In *Proceedings of the 13th AAI Conference on Interactive Decision
683 Theory and Game Theory, AAAIWS'11-13*, pages 36–42. AAAI Press, January 2011.
- 684 Rachel L. Kendal, Neeltje J. Boogert, Luke Rendell, Kevin N. Laland, Mike Webster, and Patricia L. Jones. Social
685 learning strategies: Bridge-building between fields. *Trends in cognitive sciences*, 22(7):651–665, 2018.
- 686 Pascal P. Klamser and Pawel Romanczuk. Collective predator evasion: Putting the criticality hypothesis to the test.
687 *PLOS Computational Biology*, 17(3):e1008832, March 2021. ISSN 1553-7358. doi: 10.1371/journal.pcbi.1008832.
- 688 Raphael Köster, Edgar A. Duéñez-Guzmán, William A. Cunningham, and Joel Z. Leibo. Tabula Rasa Agents Display
689 Emergent In-group Behavior. *Proceedings of the National Academy of Sciences*, 122(25):e2319947121, June 2025.
690 doi: 10.1073/pnas.2319947121.
- 691 Kevin N. Laland. Social learning strategies. *Learning & Behavior*, 32(1):4–14, February 2004. ISSN 1543-4508. doi:
692 10.3758/BF03196002.
- 693 Ellouise Leadbeater and Lars Chittka. Social Learning in Insects — From Miniature Brains to Consensus Building.
694 *Current Biology*, 17(16):R703–R713, August 2007. ISSN 0960-9822. doi: 10.1016/j.cub.2007.06.012.
- 695 David S. Leslie and Edmund J. Collins. Individual q-learning in normal form games. *SIAM Journal on Control and
696 Optimization*, 44(2):495–514, 2005. doi: 10.1137/S0363012903437976.
- 697 Udari Madhushani Sehwal, Alex McAvooy, and Joshua B. Plotkin. Collective artificial intelligence and evolutionary
698 dynamics. *Proceedings of the National Academy of Sciences*, 122(25):e2505860122, June 2025. doi: 10.1073/pnas
699 .2505860122.
- 700 Richard P. Mann. Collective decision making by rational individuals. *Proceedings of the National Academy of Sciences*,
701 115(44):E10387–E10396, October 2018. doi: 10.1073/pnas.1811964115.
- 702 Richard P. Mann. Collective decision-making by rational agents with differing preferences. *Proceedings of the National
703 Academy of Sciences*, 117(19):10388–10396, May 2020. doi: 10.1073/pnas.2000840117.
- 704 T. J. H. Morgan, L. E. Rendell, M. Ehn, W. Hoppitt, and K. N. Laland. The evolutionary basis of human social learning.
705 *Proceedings of the Royal Society B: Biological Sciences*, 279(1729):653–662, July 2011. ISSN 0962-8452. doi:
706 10.1098/rspb.2011.1172.
- 707 Anis Najar, Emmanuelle Bonnet, Bahador Bahrami, and Stefano Palminteri. The actions of others act as a pseudo-
708 reward to drive imitation in the context of social reinforcement learning. *PLOS Biology*, 18(12):e3001028, December
709 2020. ISSN 1545-7885. doi: 10.1371/journal.pbio.3001028.
- 710 Kamal K. Ndousse, Douglas Eck, Sergey Levine, and Natasha Jaques. Emergent Social Learning via Multi-agent
711 Reinforcement Learning. In *Proceedings of the 38th International Conference on Machine Learning*, pages
712 7991–8004. PMLR, July 2021.
- 713 Andreas Olsson, Ewelina Knapska, and Björn Lindström. The neural and computational systems of social learning.
714 *Nature Reviews Neuroscience*, 21(4):197–212, April 2020. ISSN 1471-0048. doi: 10.1038/s41583-020-0276-4.
- 715 Alfonso Pérez-Escudero and Gonzalo G. de Polavieja. Collective Animal Behavior from Bayesian Estimation and
716 Probability Matching. *PLOS Computational Biology*, 7(11):e1002282, November 2011. ISSN 1553-7358. doi:
717 10.1371/journal.pcbi.1002282.
- 718 Charles Perreault, Cristina Moya, and Robert Boyd. A Bayesian approach to the evolution of social learning. *Evolution
719 and Human Behavior*, 33(5):449–459, September 2012. ISSN 1090-5138. doi: 10.1016/j.evolhumbehav.2011.12.0
720 07.
- 721 Iyad Rahwan, Manuel Cebrian, Nick Obradovich, Josh Bongard, Jean-François Bonnefon, Cynthia Breazeal, Jacob W.
722 Crandall, Nicholas A. Christakis, Iain D. Couzin, Matthew O. Jackson, Nicholas R. Jennings, Ece Kamar, Isabel M.

- 723 Kloumann, Hugo Larochelle, David Lazer, Richard McElreath, Alan Mislove, David C. Parkes, Alex ‘Sandy’
724 Pentland, Margaret E. Roberts, Azim Shariff, Joshua B. Tenenbaum, and Michael Wellman. Machine behaviour.
725 *Nature*, 568(7753):477–486, April 2019. ISSN 1476-4687. doi: 10.1038/s41586-019-1138-y.
- 726 L. Rendell, R. Boyd, D. Cownden, M. Enquist, K. Eriksson, M. W. Feldman, L. Fogarty, S. Ghirlanda, T. Lillicrap, and
727 K. N. Laland. Why Copy Others? Insights from the Social Learning Strategies Tournament. *Science*, 328(5975):
728 208–213, April 2010. doi: 10.1126/science.1184719.
- 729 Alan R. Rogers. Does biology constrain culture? *American Anthropologist*, 90(4):819–831, 1988. doi: 10.1525/aa.1
730 988.90.4.02a00030.
- 731 James Rudd-Jones, María Pérez-Ortiz, and Mirco Musolesi. Understanding Individual Decision-Making in Multi-Agent
732 Reinforcement Learning: A Dynamical Systems Approach, December 2025.
- 733 Mohammad Salahshour. Perceptual rationality: An evolutionary game theory of perceptually rational decision-making.
734 *Royal Society Open Science*, 12(10), 2025. doi: 10.1098/rsos.251125.
- 735 Mohammad Salahshour and Iain D. Couzin. Evolution of altruistic rationality provides a solution to social dilemmas
736 via rational reciprocity. *Physical Review Research*, 7(3):033211, 2025. doi: 10.1103/sz5b-j75y.
- 737 Yuzuru Sato and James P. Crutchfield. Coupled replicator equations for the dynamics of learning in multiagent systems.
738 *Physical Review E*, 67(1):015206, January 2003. ISSN 1063-651X, 1095-3787. doi: 10.1103/PhysRevE.67.015206.
- 739 Marten Scheffer, Jordi Bascompte, William A. Brock, Victor Brovkin, Stephen R. Carpenter, Vasilis Dakos, Hermann
740 Held, Egbert H. Van Nes, Max Rietkerk, and George Sugihara. Early-warning signals for critical transitions. *Nature*,
741 461(7260):53–59, 2009. doi: 10.1038/nature08227.
- 742 David Schultner, Lucas Molleman, and Björn Lindström. Feature-based reward learning shapes human social learning
743 strategies. *Nature Human Behaviour*, pages 1–16, July 2025a. ISSN 2397-3374. doi: 10.1038/s41562-025-02269-4.
- 744 David Schultner, Lucas Molleman, and Björn Lindström. Reward is enough for social learning. *Trends in Cognitive
745 Sciences*, 29(9):787–789, September 2025b. ISSN 1364-6613, 1879-307X. doi: 10.1016/j.tics.2025.06.012.
- 746 Anna Sigalou and Richard P Mann. Evolutionary stability of social interaction rules in collective decision-making.
747 *Physical Biology*, 20(4):045003, May 2023. ISSN 1478-3975. doi: 10.1088/1478-3975/acd297.
- 748 Hidezo Suganuma, Kentaro Katahira, Hisashi Ohtsuki, and Tatsuya Kameda. How social learning enhances—or
749 undermines—efficiency and flexibility in collective decision-making under uncertainty. *Proceedings of the National
750 Academy of Sciences*, 122(48):e2516827122, December 2025. doi: 10.1073/pnas.2516827122.
- 751 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. Adaptive Computation and Machine
752 Learning Series. The MIT Press, Cambridge, Massachusetts, second edition edition, 2018. ISBN 978-0-262-03924-6.
- 753 Reiji Suzuki and Takaya Arita. Interactions between learning and evolution: The outstanding strategy generated by the
754 Baldwin effect. *Bio Systems*, 77(1-3):57–71, November 2004. ISSN 0303-2647. doi: 10.1016/j.biosystems.2004.04
755 .002.
- 756 Max Taylor-Davies, Neil Bramley, and Christopher G. Lucas. A Rational Framework for Group-Based Selective Social
757 Learning. *Open Mind*, 9:677–708, May 2025. ISSN 2470-2986. doi: 10.1162/opmi_a_00205.
- 758 Ulf Toelch, Matthew J. Bruce, Lesley Newson, Peter J. Richerson, and Simon M. Reader. Individual consistency and
759 flexibility in human social information use. *Proceedings of the Royal Society B: Biological Sciences*, 281(1776):
760 20132864, February 2014. ISSN 0962-8452. doi: 10.1098/rspb.2013.2864.
- 761 Wataru Toyokawa and Wolfgang Gaissmaier. Conformist social learning leads to self-organised prevention against
762 adverse bias in risky decision making. *eLife*, 11:e75308, May 2022. ISSN 2050-084X. doi: 10.7554/eLife.75308.
- 763 Wataru Toyokawa, Andrew Whalen, and Kevin N. Laland. Social learning strategies regulate the wisdom and
764 madness of interactive crowds. *Nature Human Behaviour*, 3(2):183–193, February 2019. ISSN 2397-3374. doi:
765 10.1038/s41562-018-0518-x.
- 766 Alan N. Tump, Dominik Deffner, Timothy J. Pleskac, Pawel Romanczuk, and Ralf H. J. M. Kurvers. A Cognitive
767 Computational Approach to Social and Collective Decision-Making. *Perspectives on Psychological Science*, 19(2):
768 538–551, March 2024. ISSN 1745-6916. doi: 10.1177/17456916231186964.
- 769 Matthew A. Turner, Cristina Moya, Paul E. Smaldino, and James Holland Jones. The form of uncertainty affects selection
770 for social learning. *Evolutionary Human Sciences*, 5:e20, January 2023. ISSN 2513-843X. doi: 10.1017/ehs.2023.11.

- 771 Karl Tuyls, Katja Verbeeck, and Tom Lenaerts. A selection-mutation model for q-learning in multi-agent systems. In
772 *Proceedings of the Second International Joint Conference on Autonomous Agents and Multiagent Systems - AAMAS*
773 '03, page 693, Melbourne, Australia, 2003. ACM Press. ISBN 978-1-58113-683-8. doi: 10.1145/860575.860687.
- 774 Alexandra Witt, Wataru Toyokawa, Kevin N. Lala, Wolfgang Gaissmaier, and Charley M. Wu. Humans flexibly
775 integrate social information despite interindividual differences in reward. *Proceedings of the National Academy of*
776 *Sciences*, 121(39):e2404928121, September 2024. doi: 10.1073/pnas.2404928121.
- 777 Charley M. Wu, Dominik Deffner, Björn Kahl, Björn Meder, Mark K. Ho, and Ralf H. J. M. Kurvers. Adaptive
778 mechanisms of social and asocial learning in immersive collective foraging. *Nature Communications*, 16(1):3539,
779 2025. doi: 10.1038/s41467-025-58365-6.
- 780 Lei Zhang and Jan Gläscher. A brain network supporting social influences in human decision-making. *Science Advances*,
781 6(34):eabb4159, August 2020. doi: 10.1126/sciadv.abb4159.